1    **A review of uncertainty in in situ measurements and data sets of sea-surface**

2    **temperature**

3

4    John J. Kennedy, Met Office Hadley Centre, FitzRoy Road, Exeter, EX1 3PB, UK

5    (john.kennedy@metoffice.gov.uk)

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24 **Abstract**

25

26 Archives of in situ sea-surface temperature (SST) measurements extend back more than

27 160 years. Quality of the measurements is variable and the area of the oceans they sample

28 is limited, especially early in the record and during the two World Wars. Measurements

29 of SST and the gridded data sets that are based on them are used in many applications so

30 understanding and estimating the uncertainties are vital. The aim of this review is to give

31 an overview of the various components that contribute to the overall uncertainty of SST

32 measurements made in situ and of the data sets that are derived from them. In doing so, it

33 also aims to identify current gaps in understanding. Uncertainties arise at the level of

34 individual measurements with both systematic and random effects and, although these

35 have been extensively studied, refinement of the error models continues. Recent

36 improvements have been made in the understanding of the pervasive systematic errors

37 that affect the assessment of long-term trends and variability. However, the adjustments

38 applied to minimize these systematic errors are uncertain and these uncertainties are

39 higher before the 1970s and particularly large in the period surrounding the Second

40 World War owing to a lack of reliable metadata. The uncertainties associated with the

41 choice of statistical methods used to create globally complete SST data sets have been

42 explored using different analysis techniques but they do not incorporate the latest

43 understanding of measurement errors and they want for a fair benchmark against which

44 their skill can be objectively assessed. These problems can be addressed by the creation

45 of new end-to-end SST analyses and by the recovery and digitization of data and

46 metadata from ship log books and other contemporary literature.

47

## 1. Introduction

49

50  Measurements of the temperature of the sea surface have been made for more than 200

51  years for a wide variety of purposes. The earliest measurements of sea-surface

52  temperature (SST) in the eighteenth century were taken out of pure scientific interest.

53  Later, after the connection between SST and ocean currents was made, large numbers of

54  measurements were made for the construction of navigational charts. In the twentieth

55  century, the needs of weather forecasting and, to an extent, the need to produce marine

56  climate summaries determined the quantity and quality of observations. Most historical

57  SST measurements were not made by dedicated scientific vessels, but by voluntary

58  observing ships (VOS) on the basis that they would contribute to the safety of life at sea.

59  This is reflected in the geographical distribution of observations, which are largely

60  confined to major shipping lanes.

61

62  Nowadays, in situ measurements of SST – those made at the surface as opposed to those

63  made remotely by satellites or aircraft – are used in diverse applications. They are used

64  directly in calibration and validation of satellite retrievals and they are assimilated into

65  ocean analyses [*Roberts-Jones et al.*, 2012]. They are also used to construct data sets of

66  summaries of SST on regular grids and globally-complete SST fields are created using

67  statistical techniques to impute SSTs in regions where there are no observations. The SST

68  data sets and statistical SST 'reconstructions' or 'analyses' are widely used, for example

69  as an index of global climate change [*Morice et al.*, 2012], as a boundary condition for

70　climate simulations [*Folland*, 2005] and reanalyses [*Simmons et al.*, 2010], as initial

71　conditions for decadal forecasts [*Smith et al.*, 2007], in studies of hurricane formation

72　[*Saunders and Harris*, 1997] and in studies of the impact of climate change on marine

73　ecosystems [*Sheppard and Rayner*, 2002].

74

75　As the demands for SST measurements have changed, so have the instruments used to

76　make them, and so have the ships and other vessels from which the measurements were

77　made. The first systematic observations were made using buckets to collect a water

78　sample. Buckets made of wood, canvas, tin, leather, brass, rubber and plastic – of designs

79　as various as the materials employed in their construction – have all been used to measure

80　the temperature of the surface layers of the ocean. There are two problems with this

81　approach. The first is that during the collection and hauling, the temperature of the water

82　sample can be modified by the combined actions of latent and sensible heat transfer and

83　the warmth of the Sun. Even in the best conditions, an accurate measurement requires

84　diligence on the part of the sailor; that is the second problem. Improvements to minimize

85　the physical effects were made to bucket designs during the 1950s, but as ships became

86　larger and faster, the making of the measurements became not just thankless, but

87　dangerous.

88

89　After the advent of steam ships in the late nineteenth century, it was routine to measure

90　the temperature of the sea water that was circulated through the steam condenser.

91　Condenser inlet measurements and later, engine room inlet (ERI) measurements, were

92　often recorded in ship logbooks, but they were not entered into meteorological logs until

93    the 1930s. The convenience of using measurements that were made as a matter of routine,

94    and the attendant reduction in the risk of losing a bucket or sailor overboard, meant that

95    ERI measurements became the preferred method for measuring SST on board ships

96    during the latter half of the twentieth century. That is not to say that the method was

97    without its difficulties. Modification of the temperature of the water between inlet and

98    thermometer was still a problem and it was now compounded by the varying depth of the

99    measurements.

100

101   Since the 1970s, a growing number of ships have been fitted with dedicated sensors either

102   outside or inside the hull. These have been joined by a growing array of moored and

103   drifting buoys which make automated measurements that are relayed by satellite. At

104   present, around 90% of all SST observations come from buoys. In calm conditions

105   drifting buoys measure at a nominal depth of between 10 and 20 cm depending on their

106   size. However, wave motion means that in some conditions the buoy will be submerged

107   for part of the time and report temperatures that are representative of something like the

108   upper 2 m.

109

110   Moored buoys are fixed platforms, akin, in some ways, to meteorological stations on

111   land. They come in a variety of shapes and sizes. Most are a few meters in height and

112   width, but the largest in regular use are the 12 m Discus buoys designed to weather the

113   wilder climates of the northern oceans. There are two loose groupings of moored buoys:

114   the Global Tropical Moored Buoy Array (GTMBA) and a more diverse group of coastal

115   moorings mostly around the US. The GTMBA has regular arrays of moorings in the

116    tropical Pacific, Atlantic and Indian Oceans. The majority of moored buoys measure SST

117    at a nominal depth of 1 m. Some measure slightly deeper and some moorings make

118    measurements at a range of depths.

119

120    SST measurements from ships and buoys together with near-surface measurements made

121    by oceanographic cruises have been gathered in digital archives. The largest and most

122    comprehensive of these is the International Comprehensive Ocean-Atmosphere Data Set

123    (ICOADS, *Woodruff et al.* [2011]). The latest release of ICOADS, release 2.5, contains

124    individual marine reports from 1662 to 2007, but air and sea temperature measurements

125    only start to appear in the 19[th] Century. Metadata giving information about some of the

126    measurements and the ships that make them is also provided and is now complemented

127    by information from regular bulletins such as WMO publication 47

128    (http://www.wmo.int/pages/prog/www/ois/pub47/pub47-home.htm).

129

130    Other digital archives exist. Research vessel (RV

131    http://coaps.fsu.edu/RVSMDC/index.shtml) data are gathered at the Research Vessel

132    Surface Meteorology Data Center at Florida State University. Woods Hole

133    Oceanographic Institute (http://www.whoi.edu/) maintains an archive of research

134    mooring data and the OceanSites website (http://www.oceansites.org/data/index.html)

135    provides links to other mooring data. The Pacific Marine Environmental Laboratory

136    maintains an archive of water temperature measurements from the GTMBA at a range of

137    depths and time resolutions that are not available in ICOADS

138    (http://www.pmel.noaa.gov.tao/global/global.html). Near-surface measurements from

139  other sub-surface sources such as the Argo array of autonomous profiling floats also

140  exist.

141

142  Despite being comprehensive, ICOADS is incomplete. Large archives of paper records

143  exist around the world and many of these have yet to be digitized. It is not possible yet to

144  know exactly how many undigitized records remain because there is no definitive

145  catalogue of global archives. What is known is that many archives that have been

146  identified are far from being exhausted. The potential for reducing the uncertainty in SST

147  analyses as well as in reconstructions of other marine variables is clear, but funding,

148  particularly sustained funding for the efforts to identify, image and key the data has

149  proved difficult to find. Nonetheless, there have been some successes such as a project to

150  crowd source the keying of Royal Navy logbooks from the First World War. Volunteers

151  on the OldWeather.org project keyed pages from the logbooks online. In the three years

152  since the project started more than 1.6 million weather observations have been digitized,

153  by around 16,400 volunteers.

154

155  The observing network was not created with a single purpose in mind. It was certainly

156  not intended to meet the stringent criteria demanded for monitoring long-term

157  environmental change. Nonetheless, historical SST measurements have been widely used

158  in such studies. In a 2010 paper, *Jones and Wigley* [2010] identified uncertainties

159  associated with pervasive systematic errors in SST data sets as an important uncertainty

160  in the estimation of global temperature trends. The obvious gulf between the ideal and the

161  reality leads naturally to questions about the reliability of the SST record. Often this

162    question is couched as a yes/no dichotomy: "are SST records reliable?" A more useful

163    question is "How reliable are they?" Although historical measurements were not made for

164    climate research, or any single purpose, it does not mean that it is impossible to derive

165    from them a record that is useful to a particular end. However, it does mean that special

166    care must be taken in identifying and, as best as possible, quantifying uncertainties.

167

168    In using SST observations and the analyses that are based on them, it is important to

169    understand the uncertainties inherent in them and the assumptions and statistical methods

170    that have gone into their creation. In this review I aim to give an overview of the various

171    components that contribute to the overall uncertainty of SST measurements made in situ

172    and of the data sets that are derived from them. In doing so, I also aim to identify current

173    gaps in understanding.

174

175    Section 2 provides a classification of uncertainties. The classifications are not definitive,

176    nor are they completely distinct. They do, however, reflect the way in which uncertainties

177    have been approached in the literature and provide a useful framework for thinking about

178    the uncertainties in SST data sets. The uncertainties have been tackled in ascending order

179    of abstraction from the random errors associated with individual observations to the

180    generic problem of unknown unknowns. In this review quoted uncertainties represent one

181    standard deviation of the relevant distribution unless otherwise stated. Section 3 applies

182    this framework to analyze progress and understanding under each of the headings. Some

183    shortcomings of the presentation of uncertainties are discussed in section 4 along with

184    possible solutions. Section 5 reviews how some analyses have used knowledge of likely

185 errors in SST data sets to minimize their exposure to uncertainty. Section 6 briefly

186 discusses SST retrievals from satellites and how these have been used to understand the

187 in situ record. The review concludes with a summary of possible future directions.

188

189 **2. General Classification of Uncertainties**

190

191 Throughout this review the distinction will be made between an *error* and an *uncertainty*.

192 The distinction between the two loosely follows the usage in the Guide to the Expression

193 of Uncertainty in Measurement (GUM) [*BIPM*, 2008]. The *error* in a measurement is the

194 difference between some idealized "true value" and the measured value and is

195 unknowable. The GUM defines the uncertainty of a measurement as the "parameter,

196 associated with the result of a measurement, that characterizes the dispersion of the

197 values that could reasonably be attributed to the measurand". This is the sense in which

198 uncertainty is generally meant in the following discussion. This is not necessarily the

199 same usage as is found in the cited papers. It is common to see the word error used as a

200 synonym for uncertainty such as in the commonly used phrases standard error and

201 analysis error.

202

203 Broadly speaking, errors in individual SST observations have been split into two

204 groupings: random observational errors and systematic observational errors. Although

205 this is a convenient way to deal with the uncertainties, errors in SST measurements will

206 generally share a little of the characteristics of each.

207

208    *Random observational errors* occur for many reasons: misreading of the thermometer,

209    rounding errors, the difficulty of reading the thermometer to a precision higher than the

210    smallest marked gradation, incorrectly recorded values, errors in transcription from

211    written to digital sources and sensor noise among others. Although they might confound a

212    single measurement, the independence of the individual errors means they tend to cancel

213    out when large numbers are averaged together. Therefore, the contribution of random

214    independent errors to the uncertainty on the global average SST is much smaller than the

215    contribution of random error to the uncertainty on a single observation even in the most

216    sparsely observed years. Nonetheless, where observations are few, random observational

217    errors can be an important component of the total uncertainty.

218

219    *Systematic observational errors* are much more problematic because their effects become

220    relatively more pronounced as greater numbers of observations are aggregated.

221    Systematic errors might occur because a particular thermometer is mis-calibrated, or

222    poorly sited. No amount of averaging of observations from a thermometer that is mis-

223    calibrated such that it reads 1 K too high will reduce the error in the aggregate below this

224    level save by chance. However, in many cases the systematic error will depend on the

225    particular environment of the thermometer and will therefore be independent from ship to

226    ship. In this case, averaging together observations from many different ships or buoys

227    will tend to reduce the contribution of systematic observational errors to the uncertainty

228    of the average.

229

230  In the 19th and early 20th century, the majority of observations were made using buckets

231  to haul a sample of water up to the deck for measurement. Although buckets were not

232  always of a standard shape or size, they had a general tendency under typical

233  environmental conditions to lose heat *via* evaporation or directly to the air when the air-

234  sea temperature difference was large. *Folland and Parker* [1995] provide a more

235  comprehensive survey of the problem which was already well known in the early 20th

236  Century (see, for example, the introduction to *Brooks* [1926]). *Pervasive systematic*

237  *observational errors* like the cold bucket bias are particularly pertinent for climate studies

238  because the errors affect the whole observational system and change over time as

239  observing technologies and practices change. The change can be gradual as old methods

240  are slowly phased out, but they can also be abrupt, reflecting significant geopolitical

241  events such as the Second World War [*Thompson et al.*, 2008]. Rapid changes also arise

242  because the digital archives of marine meteorological reports (ICOADS *Woodruff et al.*

243  [2011]) are themselves discontinuous.

244

245  Generally, systematic errors are dealt with by making adjustments based on knowledge of

246  the systematic effects. The adjustments are uncertain because the variables that determine

247  the size of the systematic error are imperfectly known. The atmospheric conditions at the

248  point where the measurement was made, the method used to make the measurement –

249  ERI or bucket – the material used in the construction of the bucket if one was used, as

250  well as the general diligence of the sailors making the observations have not in many

251  cases been reliably recorded. Part of the uncertainty can be estimated by allowing

252  uncertain parameters and inputs to the adjustment algorithms to be varied within their

253    plausible ranges thus generating a range of adjustments (e.g., *Kennedy et al.* [2011c]).

254    This *parametric uncertainty* gives an idea of the uncertainties associated with poorly

255    determined parameters within a particular approach, but it does not address the more

256    general uncertainty arising from the underlying assumptions. This uncertainty will be

257    dealt with later as *structural uncertainty*.

258

259    First, however, there are a number of other uncertainties associated with the creation of

260    the gridded data sets and SST analyses that are commonly used as a convenient

261    alternative to dealing with individual marine observations. The uncertainties are closely

262    related because they arise in the estimation of area-averages from a finite number of

263    noisy and often sparsely-distributed observations.

264

265    In *Kennedy et al.*, [2011b] two forms of this uncertainty were considered: *grid-box*

266    *sampling uncertainty* and *large-scale sampling uncertainty* (which they referred to as

267    coverage uncertainty). Grid-box sampling uncertainty refers to the uncertainty accruing

268    from the estimation of an area-average SST anomaly within a grid box from a finite, and

269    often small, number of observations. Large-scale sampling uncertainty refers to the

270    uncertainty arising from estimating an area-average for a larger area that encompasses

271    many grid boxes that do not contain observations. Although these two uncertainties are

272    closely related, it is often easier to estimate the grid-box sampling uncertainty, where one

273    is dealing with variability within a grid box, than the large-scale sampling uncertainty,

274    where one must take into consideration the rich spectrum of variability at a global scale.

275

276    Although some gridded SST data sets contain many grid boxes which are not assigned an

277    SST value because they contain no measurements, other SST data sets – oftentimes

278    referred to as SST analyses – use a variety of techniques to fill the gaps. They use

279    information gleaned from data-rich periods to estimate the parameters of statistical

280    models that are then used to estimate SSTs in the data voids, often by interpolation or

281    pattern fitting. There are many ways to tackle this problem and all are necessarily

282    approximations to the truth. The correctness of the *analysis uncertainty* estimates derived

283    from these statistical methods are conditional upon the correctness of the methods, inputs

284    and assumptions used to derive them. No method is correct therefore analytic

285    uncertainties based on a particular method will not give a definitive estimate of the true

286    uncertainty. To gain an appreciation of the full uncertainty it is necessary to factor in the

287    lack of knowledge about the correct methods to use, which brings the discussion back to

288    structural uncertainty.

289

290    There are many scientifically defensible ways to produce a data set. For example, one

291    might choose to fill gaps in the data by projecting a set of Empirical Orthogonal

292    Functions (EOFs) onto the available data. Alternatively, one might opt to fill the data

293    using simple optimal interpolation. Both are defensible approaches to the problem, but

294    each will give different results. In the process of creating any data set, many such choices

295    are made. *Structural uncertainty* [*Thorne et al.*, 2005] is the term used to understand the

296    spread that arises from the many choices and foundational assumptions that can be (and

297    have to be) made during data set creation. The character of structural uncertainty is

298    somewhat different to the other uncertainties considered so far. The uncertainty

299     associated with a measurement error, for example, assumes that there is some underlying

300     distribution that characterizes the dispersion of the measured values. In contrast, there is

301     generally no underlying "distribution of methods" that can be used to quantify the

302     structural uncertainty. Furthermore, the diverse approaches taken by different teams

303     might reflect genuine scientific differences about the nature of the problems to be tackled.

304     Consequently, structural uncertainty is one of the more difficult uncertainties to quantify

305     or explore efficiently. It requires multiple, independent attempts to resolve the same

306     difficulties, it is an ongoing commitment, and it does not guarantee that the true value

307     will be encompassed by those independent estimates. Nevertheless, the role that the

308     creation of multiple independent estimates and their comparison has played in

309     uncovering, resolving, and quantifying some of the more mystifying uncertainties in

310     climate analyses is unquestionable. The most obvious – one might say, notorious –

311     examples are those of tropospheric temperature records made using satellites and

312     radiosondes [*Thorne et al.*, 2011] and sub-surface ocean temperature analyses [*Lyman et*

313     *al.*, 2010; *Abraham et al.,* 2013].

314

315     Which leads finally to *unknown unknowns*. On February 12[th] 2002, at a news briefing at

316     the US Department of Defense, Donald Rumsfeld memorably divided the world of

317     knowledge into three quarters:

318

319         *"There are known knowns. These are things we know we know. We also know there*

320     *are known unknowns. That is to say, we know there are some things we do not know. But*

321     *there are also unknown unknowns, the ones we don't know we don't know."*

322

323  In the context of SST uncertainty, unknown unknowns are those things that have been

324  overlooked. By their nature, unknown unknowns are unquantifiable; they represent the

325  deeper uncertainties that beset all scientific endeavors. By deep, I do not mean to imply

326  that they are necessarily large. In this review I hope to show that the scope for revolutions

327  in our understanding is limited. Nevertheless, refinement through the continual evolution

328  of our understanding can only come if we accept that our understanding is incomplete.

329  Unknown unknowns will only come to light with continued, diligent and sometimes

330  imaginative investigation of the data and metadata.

331

332  **3. The Current State of Uncertainty in in situ SST Analyses**

333

334  The classification of uncertainties outlined in section 2 will now be used as a framework

335  to assess uncertainties in the global data sets based on in situ measurements. Preliminary

336  to this it will be helpful to define what exactly is meant by sea-surface temperature.

337

338  **3.1 Defining Sea-surface Temperature**

339

340  Traditionally, in situ SST analyses have been considered representative of the upper ten

341  or so meters of the ocean. However, the near-surface temperature structure of the ocean

342  can be rather complex. Under conditions of low wind speed and high insolation, a stable

343  stratified layer of warm water can form near the surface. For a recent review see *Kawai*

344  *and Wada* [2007].  The diurnal temperature range of the sea-surface can, under certain

345     conditions, exceed 5 K and, somewhat attenuated, penetrate to many tens of meters

346     [*Prytherch et al.,* 2013]. This can lead to strong temperature gradients in the upper few

347     meters of the ocean and consequently measurements made at the same time and location

348     but at different depths can record quite different temperatures. Temperatures measured at

349     the same depth but at different times of day can also differ markedly.

350

351     *Donlon et al.* [2007] proposed that the depth of the measurement be recorded along with

352     the temperature as a first step to reconciling measurements made at different depths and

353     different times of day. *Donlon et al.* [2007] also introduced the concept of an SST

354     foundation (SST$_{fnd}$) temperature. The current definition (https://www.ghrsst.org/ghrsst-

355     science/sst-definitions/) of "SST$_{fnd}$, is the temperature free of diurnal temperature

356     variability, i.e., SST$_{fnd}$ is defined as the temperature at the first time of the day when the

357     heat gain from the solar radiation absorption exceeds the heat loss at the sea surface." It is

358     generally assumed that the upper few meters of the ocean are of approximately constant

359     temperature at this point. SST$_{fnd}$ has proved a practical reference point for comparing and

360     combining satellite observations [*Roberts-Jones et al.*, 2012] and was intended to provide

361     "a more precise, well-defined quantity than the previous loosely-defined bulk SST"

362     *Donlon et al.* [2007].

363

364     Unfortunately, such niceties of definition are not readily applicable to historical SST

365     measurements and the effect of the interaction between measurement depth and water

366     temperature on SST measurements in in situ archives is not clear. For many ships that

367     measure the temperature of water drawn in below the surface, the depth of the

368  measurements is not known and is likely to have changed depending on how heavily the

369  ship was loaded. Nor is it clear to what extent any warm surface layer is mixed with

370  cooler subsurface water by the passage of the ship or by the interaction of wind, water,

371  Sun and hull [*Amot*, 1954; *Stevenson*, 1964]. Similar interactions have been noted closer

372  to the surface with moored buoys [*Kawai and Kawamura*, 2000]. *James and Fox* [1972]

373  found that ERI measurements from ships became progressively warmer relative to

374  simultaneous bucket observations as the depth of the ERI measurement increased, a

375  similar pattern to that seen by *Kent et al.* [1993]. *Reynolds et al.* [2010] found that

376  measurements made by ships, which were largely ERI measurements in their study

377  period, were on average warmer than nearby drifting buoy observations made nearer to

378  the surface.

379

380  Nonetheless, the concept of the foundation SST can be used to get an idea of how

381  changing measurement depth might have affected SST trends in the absence of other

382  considerations. Figure 1 shows an upper estimate of the potential size of the effect of

383  changing measurement depth on global average SST over time (for calculation details see

384  Appendix A). The assumption is that buckets and buoys measure in the upper 30 cm and

385  engine room measurements are measuring $SST_{fnd}$. The estimated global average bias

386  (relative to the 1961-1990 average) is less than 0.1 K at all times and from 1945 onwards

387  is less than 0.05 K. The bias is largest in the early record when all measurements were

388  made using buckets which sample in the upper meter of the water column. In the more

389  recent period, the blend of buckets, ERI measurements and buoys leads to a smaller,

390   time-varying bias. Although the size of the effect is modest at a global level, locally the

391   average diurnal warming can exceed 0.5K, which would imply a larger effect.

392

393   A related problem is that changing times of observation could potentially interact with the

394   diurnal cycle of temperature leading to spurious trends in the data. *Kent et al.* [2010] note

395   "*The implicit assumption is that the sampling of conditions is regular enough that no*

396   *regional or time-varying bias is introduced into the datasets by neglecting such effects.*"

397   Ships currently make SST observations at regular intervals throughout the day, typically

398   every four or six hours, which is sufficient to minimize the aliasing of diurnal cycles,

399   particularly if the measurements are made at depth. During earlier periods when buckets

400   were widely used, there were systematic changes in the time of observation that might

401   have a more pronounced effect on average SSTs but this has not been quantified.

402

403   Even when the measurement depth is known, there are potential problems. Metadata in

404   WMO Publication 47 show that ships measure water temperatures through a wide range

405   of depths from the near surface down to around 25 m [*Kent et al.*, 2007]. Although the

406   average depth was typically less than 10 m, the deepest measurements could be sampling

407   water that is colder than the $SST_{fnd}$. How large this effect might be is not yet well

408   understood.

409

410   *Chiodi and Harrison* [2006] identified large-scale warm surface features using SST

411   retrievals from microwave satellite instruments that persisted for several days. The warm

412   layer was observed at night suggesting that the effect was independent from diurnal

413     warming and they hypothesized that the multi-day warming might have been confined to

414     a relatively shallow layer between 1 and 5 m thick. The implication is that the depth of

415     the SST foundation temperature can vary rapidly and that it can be much shallower than

416     the deepest in situ SST measurements. During a two week cruise, *Matthews and*

417     *Matthews* [2013] found persistent temperature difference between the surface and 3 m

418     depth in the tropical Pacific. Similar warm layers can be seen in data from moored buoys.

419     Figure 2 shows time series from several moorings showing multi-day near-surface warm

420     layers that do not penetrate down to 10 m and in some cases do not reach 5 m.

421     Climatologies of mixed layer depth (MLD, see for example *de Boyer Montégut* [2004])

422     indicate large areas – in regions of upwelling and in the summer hemisphere – where the

423     average MLD is shallower than 30 m, implying measurable temperature gradients within

424     the depth range of ship SST measurements. *Grodsky et al.* [2008] also found differences

425     between SST and temperatures in the mixed layer, which were largest in areas of

426     persistent upwelling – most notably the eastern Pacific – but they did not consider the

427     possible confounding effects of systematic errors in SST or other measurements.

428

429     To isolate the specific effect of multi-day or persistent temperature stratification of the

430     near-surface waters would require regular measurements of near-surface waters at a range

431     of depths. Such an analysis is now possible thanks to the network of Argo floats [*Castro*

432     *et al.,* 2013]. In what follows, it should be noted that variations in depth will contribute to

433     the variance of measurements and will therefore be partly, or wholly, counted in

434     estimates of random and systematic measurement errors.

435

**3.2 Individual Observational Errors**

The general quality of raw SST measurements recorded in digital archives is mixed. Consequently, all SST analyses perform a stage of pre-screening, or quality control (QC) in order to remove observations of low quality and minimize the number of egregious errors. The size of the uncertainties of individual measurements will depend to a certain extent on the QC that is applied but the effects of differences in QC have not been assessed systematically.

**3.2.1 Random Measurement Errors**

Many estimates of random observational error uncertainty have been made. Although thermometers issued to ships by many port meteorological officers are calibrated, such calibration information is not routinely published, nor is there any guarantee that the temperature of a water sample measured by a well calibrated thermometer is equal to the actual SST when the sample has spent time in a bucket, or passed through the pipe work of a ship. Consequently, estimates of measurement uncertainty from the literature are empirical estimates derived from considerations of the variance of the data: for example, spatial [*Lindau,* 2003; *Kent and Challenor*, 2006; *Emery et al.*, 2001] and temporal [*Stubbs*, 1965] semivariograms, by comparing collocated observations [*O'Carroll et al.*, 2008], by resampling [*Shen et al.*, 2007], by using the variation of the variance with the number of observations [*Rayner et al.*, 2006], or by comparison with a background field [*Kent and Berry*, 2008; *Xu and Ignatov*, 2010; *Ingleby*, 2010; *Kennedy et al.*, 2011a;

459    *Atkinson et al.*, 2013]. Some of the analyses did not distinguish between random

460    observational errors and systematic observational errors, tending to combine them into

461    one estimate. In addition it is not always easy to separate the effects of spatial sampling

462    from measurement errors particularly in regions of high SST variability [*Castro et al.*,

463    2012].

464

465    A single SST measurement from a ship has a typical combined random and systematic

466    error uncertainty of around 1 K to 1.5 K. Results from individual analyses are

467    summarized in Table 1. The studies are mostly based on data from 1970 onwards.

468

469    Measurements are not all of identical quality. *Kent and Challenor* [2006] showed that in

470    the period 1970-1997 the uncertainties of measurements from ships varied with location,

471    time, measurement method and the country that recruited the ship. Uncertainties were

472    estimated to be larger in the mid-1970s probably due to data being incorrectly transmitted

473    in real time in the early days of the Global Telecommunication System. Their estimated

474    uncertainty for engine room measurements was larger than for bucket measurements.

475    *Tabata* [1978a] noted that bucket measurements *could* be accurate to 0.15 K, but that ERI

476    measurements were nearly an order of magnitude worse (1.16 K). *Ingleby* [2010]

477    estimated uncertainties for different subsets of the data and noted that manual VOSclim

478    (a high-quality subset of the VOS fleet) measurements and automated measurements

479    were of slightly higher quality than manual ship measurements in general. *Beggs et al.*

480    [2012] showed that Australia Integrated Marine Observing System ships had

481    uncertainties comparable to those from data buoys. Analyses that have looked at statistics

482    for individual ships and buoys have found that some ships and buoys take much higher

483    quality measurements than others [*Kent and Berry*, 2008; *Brasnett*, 2008; *Kennedy et al.,*

484    2011a; *Atkinson et al.*, 2013]. The subset of ships (around 40-50% of ship observations)

485    that passed the more stringent quality control procedures of *Atkinson et al.* [2013] had

486    significantly lower measurement uncertainties assessed using the method of *Kennedy et*

487    *al.* [2011a] than did the full fleet of ships. Early results on hull sensors reported by *Emery*

488    *et al.* [1997] indicated the potential for these sensors to make accurate measurements.

489    Indeed, *Kent et al.* [1993] found that hull sensors installed on ships in the Voluntary

490    Observing Ships Special Observing Project for the North Atlantic (VSOP-NA) gave

491    consistent measurements during the two year observing period.

492

493    Drifting buoy measurements are generally more accurate and consistent than ship

494    measurements, but there is a greater relative spread between the estimates which are

495    summarized in Table 2. In part these differences are likely to arise from the level of pre-

496    screening that is applied to the observations. Where quality control is more stringent,

497    estimated uncertainties are likely to be lower and, where the error variance of the

498    observations is low already, the effects of quality control and processing choices are

499    likely to be more pronounced [*Xu and Ignatov*, 2012]. *Castro et al.* [2012] considered

500    differences between drifting buoys and two different satellite products and found that

501    there was little difference between buoys produced by different manufacturers. There is

502    some evidence that the quality of drifting buoy observations has improved slightly over

503    time [*Merchant et al.*, 2012], but this has not been conclusively demonstrated. As a

504     comparison, temperature measurements from Argo have been reckoned to have an

505     uncertainty of around 0.002K [*Abraham et al.,* 2013].

506

507     Moored buoys have received less attention. Estimates of the measurement uncertainties

508     are summarized in Table 3. The two studies [*Kennedy et al.*, 2011a; *Xu and Ignatov,*

509     2010] that examined moorings from the GTMBA separately from other moorings found

510     that they had lower measurement error uncertainties. *Castro et al.* [2012] found that the

511     standard deviations of differences between moorings and satellite data were lower for

512     tropical moorings than for coastal moorings. They noted that in coastal waters there can

513     be large local variations in temperature, which satellites cannot resolve. Some moorings

514     along coastlines are located in estuaries and river mouths and are therefore less likely to

515     be representative of open ocean areas. This is perhaps one reason why *Wilkerson and*

516     *Earl* [1990], who studied US coastal buoys, found such large standard deviations between

517     ships and moorings (Table 1). *Merchant et al.* [2012] found that few coastal moorings

518     met their required stability criteria.

519

520     As noted in section 2, random observational errors are of relatively minor importance in

521     large-scale averages (see Figure 8 and section 3.5), particularly in the modern period

522     when observations are numerous. For an uncertainty of 1.0 K for a single observation due

523     to random observational error, the resulting uncertainty of a global annual average based

524     on 10000 observations would be of order 0.01 K.

525

526     **3.2.2 Random and Systematic Measurement Errors**

527

528    *Kent and Berry* [2008] and *Kennedy et al.* [2011a, 2011b] decomposed the observational

529    errors into random and systematic components. *Brasnett* [2008] and *Xu and Ignatov*

530    [2010] implicitly used the same error model – their analyses output the same statistics

531    produced by *Kent and Berry* [2008] – and the results are indeed very similar (Figure 3).

532    Estimates are summarized in Table 4. The possibility of correlated measurement errors is

533    also implicitly allowed for by *Ishii et al.* [2003] and *Hirahara et al.* [2013] who merge

534    observations from a single ship into a super observation before calculating uncertainties.

535    Adding the uncertainties in quadrature gives a combined observational uncertainty of

536    between 1 and 1.5 K, consistent with earlier estimates (Table 1) that did not differentiate

537    between the two.

538

539    In the studies listed in Table 4, the systematic component of the error was assumed to be

540    different for each ship, but this does not on its own capture the effects of pervasive

541    systematic errors. The data from *Kent and Berry* [2008], *Brasnett* [2008] and *Xu and*

542    *Ignatov* [2010] also show that the systematic observational error component for some

543    ships varies from month to month suggesting that the partitioning of systematic and

544    random effects is also a function of the time period considered.

545

546    The addition of a systematic component has a pronounced effect on the uncertainty of

547    large-scale averages comprising many observations. *Kennedy et al.* [2011b] estimated

548    that the effect of the correlations between errors was to increase the uncertainty of the

549    global annual average SST anomaly due to measurement error from 0.01 K (uncorrelated

550    case) to more than 0.05 K in the 19th Century and to more than 0.01 K even in the well-

551    observed modern period when millions of observations contribute to the annual global

552    average (see Figure 8). Systematic errors could also have a pronounced effect on

553    reconstructions when they project onto large-scale modes of variability, or on the

554    estimation of EOFs. However, because of the assumed independence of the errors

555    between ships, the correlated component of the uncertainty remains relatively

556    unimportant for the analysis of long-term trends of large-scale averages. Pervasive

557    systematic errors, which are correlated across a large proportion of the global fleet,

558    (section 3.2) are far more important from that point of view.

559

560    One of the difficulties with estimating the uncertainties associated with systematic errors

561    from individual ships is that not all observations in ICOADS can be associated with an

562    individual ship. Some of the reports have no more information than a location, time and

563    SST measurement. *Kennedy et al.* [2011b] had to make estimates of how the uncertainty

564    arising from systematic errors behaved as the number of observations increased by

565    considering the behavior at times when the majority of reports contained a ship name or

566    call-sign. They assumed that observations without call signs behaved in the same way.

567    *Kent and Berry* [2008] suggested that only ship reports with extant metadata be used in

568    climate analyses of the modern period to minimize such ambiguities. For earlier periods,

569    the gains in improved quantification of uncertainty would need to be balanced against the

570    increased uncertainty arising from reduced coverage.

571

572 Many gridded SST data sets and analyses, as well as the studies that depend on them,

573 assume that the observational errors are normally distributed, but this is not necessarily

574 the case for individual observations. *Kennedy et al.* [2011a] investigated the properties of

575 observations that had been quality controlled using the procedures described in *Rayner et*

576 *al.* [2006]. They found that in comparisons with satellite observations the distributions of

577 errors were 'fat-tailed' with the distribution of errors having a positive kurtosis. In the

578 creation of gridded data sets from SST observations, the effects of outliers can be

579 minimized somewhat by the use of resistant or robust statistics such as Winsorised, or

580 trimmed means (see e.g., *Rayner et al.* [2006]). The effect of outliers is further reduced in

581 large scale averages and the distribution of errors in large scale averages tends towards a

582 normal distribution as the number of observations increases [*Kennedy et al.*, 2011a].

583

584 **3.2.3 Summary of Individual Observational Errors**

585

586 Many estimates of uncertainties of ship and buoy SST measurements have been made. A

587 typical SST measurement made by a ship has an uncertainty of around 1-1.5K and a

588 drifting buoy observation a typical uncertainty of around 0.1-0.7K. More recent studies

589 split these uncertainties into random and systematic components, which better describe

590 the error characteristics of these platforms. However, a lack of metadata, most

591 particularly ship call signs, hampers the application of such an error model and it does not

592 capture behavior seen in SST measurements such as non-Normal distributions or

593 systematic errors that vary on time scales from months to years.

594

595 **3.3 Pervasive Systematic Errors and Biases**

596

597 *Kent et al.* [2010] conducted a review of literature on pervasive systematic errors (often

598 termed 'biases') in in situ SST measurements. Many studies have looked at the

599 differences in pervasive systematic errors between measurement methods, but fewer have

600 attempted to adjust SST records to minimize the effects of changes in instrumentation.

601

602 **3.3.1 Bias Adjustments 1850 to 1941**

603

604 The need for adjustments to minimize the cold bias associated with bucket measurements

605 in the period from 1850 to 1941 is well established. *Folland and Parker* [1995] calculated

606 adjustments using a simplified physical model of the buckets used to make SST

607 measurements combined with fields of climatological air-temperature, SST, humidity,

608 wind and solar radiation. Some parameters in their model were taken from literature and

609 others were estimated from the data. The length of time between the water sample leaving

610 the sea surface and the measurement was estimated by integrating their model until a

611 seasonal cycle in the SST was minimized. The fractional contributions of canvas and

612 wooden buckets were estimated by assuming a linear change over time from a mix of

613 wooden and canvas buckets to predominantly canvas buckets by 1920. The rate of this

614 change was estimated by minimizing the air-sea temperature difference in the tropics.

615 The same method was also used in *Rayner et al.* [2006] and *Kennedy et al.* [2011c].

616

617     *Smith and Reynolds* [2002] took an alternative approach. They adjusted SSTs based on

618     statistical relationships between Night Marine Air Temperature (NMAT) and SST. The

619     resulting adjustments were different to those produced by *Folland and Parker* [1995]

620     although the magnitude of the global average adjustment was similar. Both *Folland and*

621     *Parker* [1995] and *Smith and Reynolds* [2002] found a long term increase in the

622     magnitude of the adjustments – that is, an increasing cold bias – from the 1850s to 1941.

623

624     The methods employed by *Folland and Parker* [1995] and *Smith and Reynolds* [2002] are

625     not independent as they both rely on NMAT, which have their own particular pervasive

626     systematic errors [*Bottomley et al.,* 1990; *Rayner et al.*, 2003; *Kent et al.*, 2013]. The use

627     of NMAT to adjust SST data is, to an extent, unavoidable as the heat loss from a bucket

628     does depend on the air-sea temperature difference.

629

630     In data sets based on a ICOADS release 2.0 and later, the earlier bucket adjustments were

631     found to over-adjust SST in the period 1939-1941. *Rayner et al.* [2006] and *Smith et al.*

632     [2008] ramped the adjustments down to zero over this period. *Kennedy et al.* [2011c]

633     showed that the ramp-down corresponded to new data in that release of ICOADS that

634     included a large fraction of ERI measurements.

635

636     **3.3.2 Bias Adjustments 1941 to Present**

637

638     In the post-1941 period, *Folland and Parker* [1995], *Smith and Reynolds* [2003], *Smith*

639     *and Reynolds* [2005] and *Rayner et al.* [2006] opted not to adjust the data because they

640   found no clear evidence of the need for adjustments. However, *Rayner et al.* [2006] did

641   identify biases in Japanese and Dutch data after the Second World War. *Thompson et al.*

642   [2008] identified a discontinuity in global-average SST associated with a change in the

643   composition of ICOADS release 2.1 in late 1945. *Reynolds et al.* [2010] quantified a

644   relative bias between ship and drifting buoy measurements that they thought could lead to

645   an artificial cooling of the global average SST. *Kent et al.* [1999] applied adjustments to

646   ERI measurements, but removed the adjustment from later versions of their data set.

647

648   *Kennedy et al.* [2011c] and *Hirahara et al.* [2013] developed bias adjustments for the

649   period 1941 onwards. *Kennedy et al.* [2011c] used metadata from ICOADS, WMO

650   Publication 47, observer instructions, technical reports and scientific papers to estimate

651   biases for individual measurement types and to assign a measurement method to as many

652   observations as possible. *Hirahara et al.* [2013] used a narrower range of metadata. By

653   comparing subsamples of the data for which the metadata were known, they could

654   estimate appropriate metadata assignments for the remainder.

655

656   To estimate the bias adjustments for long-term analyses, an understanding is needed of

657   how biases varied for individual components of the observing system. Several studies

658   have examined ERI and bucket biases in ship data [*Brooks*, 1926; *Brooks*, 1928; *Lumby*,

659   1927; *Collins et al.*, 1975; *Wahl*, 1948; *Roll*, 1951; *Kirk and Gordon*, 1952; *Amot*, 1954;

660   *Perlroth*, 1962; *Saur*, 1963; *Walden*, 1966; *Knudsen*, 1966; *Tauber*, 1969; *James and*

661   *Fox*, 1972; *Tabata,* 1978a, 1978b; *Folland et al.*, 1993; *Kent et al.*, 1993] but only *Kent*

662   *and Kaplan* [2006] provide information that is time-resolved and traceable back to

663     ICOADS. There is a single study of pervasive systematic errors in hull sensor

664     measurements [*Kent et al.*, 1993], which analyzed data from a small number of ships over

665     a two year period and found that hull sensors were relatively unbiased and showed no

666     systematic change of bias with depth.

667

668     Few studies have looked at the long-term stability and calibration drifts of drifting buoys.

669     *Reverdin et al.* [2010] installed 16 drifters with high quality temperature sensors in

670     addition to their usual temperature sensors and found that the temperatures measured by

671     the drifters showed inaccuracies that were larger than the 0.1 °C target accuracy and that

672     they exhibited significant calibration drifts. This is consistent with the behavior seen by

673     *Atkinson et al.* [2013].

674

675     **3.3.3 Estimating Uncertainty in Bias Adjustments**

676

677     *Folland and Parker* [1995] did not explicitly estimate the uncertainties in their

678     adjustments. *Rayner et al.* [2006] explored the parametric uncertainty in the *Folland and*

679     *Parker* [1995] adjustments using a Monte-Carlo method. In *Smith and Reynolds* [2004]

680     the uncertainty in the bias adjustments was estimated by taking the mean-squared

681     difference between the *Smith and Reynolds* [2002] adjustments and the *Folland and*

682     *Parker* [1995] adjustments, a first-order estimate of the structural uncertainty.

683

684     *Kennedy et al.* [2011c] used a Monte-Carlo method to explore the parametric uncertainty

685     within their particular approach to bias adjustment. *Hirahara et al.* [2013] also provide

686 uncertainties on their adjustments that are a combination of analysis uncertainties and

687 regression uncertainty.

688

689 An important component of the uncertainty of adjustments for the effects of persistent

690 systematic errors arises from a lack of knowledge concerning how the measurements

691 were made. Metadata are often missing, incomplete or ambiguous and sometimes

692 different sources give conflicting information. *Kent et al.* [2007] assessed metadata from

693 ICOADS and WMO Publication 47. They found disagreement in around 20-40% of cases

694 where metadata were available from both sources. *Kennedy et al.* [2011c] allowed for up

695 to 50% uncertainty in metadata assignments based on the discrepancy between observer

696 instructions and measurement methods recorded in WMO Publication 47. *Hirahara et al.*

697 [2013] used differences between subsets of data to infer the fraction of observations made

698 using different methods.

699

700 Figure 6 compares estimated biases and metadata assignments from *Kennedy et al.*

701 [2011c] and *Hirahara et al.* [2013]. It shows that from 1945, the estimated biases agree

702 within their parametric uncertainty ranges (Figure 6a) and that the fractions of

703 measurement methods estimated by *Kennedy et al.* [2011c] from literature and other

704 metadata are consistent with the fractions inferred from the data by *Hirahara et al.* [2013]

705 (Figure 6b). However, there are two key differences that highlight the importance of

706 structural uncertainty for understanding the bias adjustments. The first difference is that

707 the phasing out of uninsulated buckets in *Hirahara et al.* [2013] happens earlier and

708 faster than allowed for in the parametric uncertainty analysis of *Kennedy et al.* [2011c]

709    (Figure 6c). In *Hirahara et al.* [2013] the changeover starts in the 1940s and is especially

710    rapidly in the early 1960s, being nearly complete by around 1962. The second difference

711    is that the estimated bias during the Second World War is higher in the analysis of

712    *Hirahara et al.* [2013] than in *Kennedy et al.* [2011c]. Further work is needed to

713    understand these differences and more complete, more reliable metadata would help

714    reduce uncertainty in SST records.

715

716    In the post-1941 period, *Smith and Reynolds* [2003] and *Smith and Reynolds* [2005]

717    estimated the uncertainty due to pervasive systematic errors by considering the difference

718    in estimated bias between measurements made in the engine rooms of the ships and

719    measurements from all ships between 1994 and 1997. They estimated a minimum 1-

720    sigma standard error in the global average of around 0.015 K. The range is similar to,

721    albeit slightly narrower than, that estimated by *Kennedy et al.* [2011c]. The difficulty

722    with the approach taken by *Smith and Reynolds* [2003], *Smith and Reynolds* [2005] and

723    *Smith et al.* [2008] is that the quoted uncertainty range is considered to be symmetric

724    whereas *Kennedy et al.* [2011c] and *Hirahara et al.* [2013] suggest that the true global

725    mean is consistently higher than *Smith et al.* [2008] in the period 1945-1960 (Figure 9).

726    It also suggests that the estimate of *Smith et al.* [2008] in the post World War 2 period

727    (1945-1950s) was slightly too conservative because it compared ERI measurements with

728    a mixture of ERI and insulated bucket measurements, whereas large numbers of

729    observations were made using buckets [*Kennedy et al.*, 2011c; *Hirahara et al.,* 2013].

730

731    **3.3.4 Refinements to Estimates of Pervasive Systematic Errors**

732  There are some factors that have not been explicitly considered in estimates of biases.

733  Refinements to the models of pervasive systematic errors will address with factors that

734  are implicitly included in random and systematic measurement uncertainties. If it is

735  possible to estimate the bias on a ship-by-ship, or observation-by-observation basis,

736  taking account of the conditions peculiar to that observation, then it might be expected

737  that uncertainties associated with random and systematic observational error will

738  decrease.

739

740  Both *Kennedy et al.* [2011c] and *Hirahara et al.* [2013] make simplifying assumptions

741  about the systematic errors associated with modern insulated buckets. Various bucket

742  designs have been used since the end of the Second World War, which are likely to have

743  different bias characteristics. Physical models could be developed for each type of bucket

744  similar to those used by *Folland and Parker* [1995], or statistical methods could be used

745  to estimate the biases as was done in *Kent and Kaplan* [2006].

746

747  Other simplifying assumptions used in all analyses include such things as assuming that

748  changes in the observing system happened linearly. Evidence suggests that changes in

749  measurement method were not always monotonic and sometimes happened abruptly (see

750  Figure 6). Improved metadata or more sophisticated statistical techniques could help

751  assess these uncertainties.

752

753  An uncertainty associated with pervasive systematic biases, which is not explicitly

754  resolved by current analyses, arises when the conditions at the time of the measurement

755    deviate from the climatological values assumed by the bias correction scheme. If, for

756    instance, the air sea temperature difference is larger than that assumed by the *Folland and*

757    *Parker* [1995] scheme, then there will be an additional systematic uncertainty that is

758    correlated strongly across synoptic spatial and temporal scales with a potential long-term

759    component where differences persist for months or years. Likewise conditions vary

760    during the day. Such discrepancies could be assessed by evaluating the systematic error

761    using local conditions. Such information could be taken from reanalyses, or an

762    appropriate bucket model could be explicitly included when SST observations are

763    assimilated into ocean-only and coupled reanalyses.

764

765    **3.3.5 Assessing the Efficacy of Bias Adjustments**

766

767    The efficacy of the bias adjustments and their uncertainties are difficult to assess. *Folland*

768    *and Parker* [1995] presented wind tunnel and ship board tests and also used their

769    adjustments to estimate the differences between bucket and ERI measurements in broad

770    latitude bands. These limited comparisons showed that their model could predict

771    experimental results to better than 0.2 K. *Folland and Salinger* [1995] presented direct

772    comparisons between air temperatures measured in New Zealand and SST measurements

773    made nearby. *Smith and Renyolds* [2002] used oceanographic observations to assess their

774    adjustments and those of *Folland and Parker* [1995]. In regions with sufficient

775    observations they found that the magnitude of the *Smith and Reynolds* [2002] adjustments

776    better explained the differences between SSTs and oceanographic observations, but the

777    phase of the annual cycle was better captured by *Folland and Parker* [1995]. *Hanawa et*

778    *al.* [2000] showed that the *Folland and Parker* [1995] adjustments improved the

779    agreement between Japanese ship data and independent SST data from Japanese coastal

780    stations in two periods: before and after the Second World War. However, the collection

781    of ship data (COADS and Kobe collections) used in *Hanawa et al.* [2000] might not have

782    had the same bias characteristics as assumed by *Folland and Parker* [1995] (based on the

783    Met Office Marine Data Bank) in developing their adjustments. Other long term coastal

784    records of water temperature exist. Some of these [*Hanna et al.,* 2006; *MacKenzie and*

785    *Schiedek,* 2007; *Cannaby and Hüsrevoğlu*, 2009] have been compared to open ocean SST

786    analyses (though not with the express intention of assessing bias adjustments), others

787    have not [*Maul et al.,* 2001; *Nixon et al.,* 2004; *Breaker et al.* 2005].

788

789    More recently, *Matthews* [2013] and *Matthews and Matthews* [2013] reported field

790    measurements of SST made using different buckets and simultaneous thermo-salinograph

791    measurements. They found negligible biases between different buckets, but their

792    experimental design involved larger buckets and shorter measurement times than were

793    used in *Folland and Parker* [1995]. Nevertheless, this highlights the potential for well-

794    designed field experiments to improve understanding of historical biases.

795

796    An analysis by *Gouretski et al.* [2012] compared SST observations with near-surface

797    measurements (0-20 m depth) taken from oceanographic profiles. It shows that the

798    overall shape of the global average is consistent between the two independent analyses,

799    but that there are differences of around 0.1 K between 1950 and 1970. These are most

800    likely attributable to residual biases, although, as noted above, actual physical differences

801    between the sea surface and the 0-20 m layer cannot be ruled out. Similar differences are

802    seen when comparing SST with the average over the 0-20 m layer of the analysis of

803    *Palmer et al.* [2007] (not shown).

804

805    Since the late 1940s, global and hemispheric average SST anomalies calculated

806    separately from adjusted bucket measurements and adjusted ERI measurements showed

807    consistent long-term and short-term changes [*Kennedy et al.*, 2011c]. From the 1990s,

808    there are also plentiful observations from drifting and moored buoys.

809

810    In contrast to the modern period, the period before 1950 is characterized by a much less

811    diverse observing fleet. During the Second World War, the majority of measurements

812    were ERI measurements. Before the war, buckets were the primary means by which SST

813    observations were made. This makes it very difficult to compare simultaneous

814    independent subsets of the data. In periods with fewer independent measurement types, it

815    might be possible to use changes in environmental conditions such as day-night

816    differences or air-sea temperature differences to diagnose systematic errors in the data.

817

818    Qualitative agreement between the long-term behavior of different global temperature

819    measures – including NMAT, SST and land temperatures – gives a generally consistent

820    picture of historical global temperature change (Figure 5), but a direct comparison is less

821    informative about uncertainty in the magnitude of the trends. *Kent et al.* [2013] showed

822    similar temporal evolution of NMAT and SST in broad latitude bands in the northern

823    hemisphere and tropics. However there are differences of up to 0.4 K in the band from

824   55°S to 15°S between 1940 and 1960. Studies such as that by *Folland* [2005] can be used

825   to make more quantitative comparisons. *Folland* [2005] compared measured land air

826   temperatures with land air temperatures from an atmosphere-only climate model that had

827   observed SSTs (with and without bucket adjustments) as a boundary forcing. He found

828   much better agreement when the SSTs were adjusted. Atmospheric reanalyses also use

829   observed SSTs along with other observed meteorological variables to infer a physically

830   consistent estimate of land surface air temperatures. *Simmons et al.* [2010] showed that

831   land air temperatures from a reanalysis driven by observed SSTs were very close to those

832   of CRUTEM3 [*Brohan et al.*, 2006] over the period 1973 to 2008. *Compo et al.* [2013]

833   showed similar results for the whole of the twentieth century although the agreement was

834   not quite so close. Although their intention was to show that land temperatures were

835   reliable, their results indicate that there is broad consistency between observed SSTs and

836   land temperatures.

837

838   **3.3.6 Summary of Pervasive Systematic Errors and Biases**

839

840   The need to adjust SST data prior to 1941 to account for a cold bias associated with the

841   use of canvas and wooden buckets is well established. There is also good evidence for the

842   need to adjust data after 1941. Adjustments for these pervasive systematic errors have

843   been developed. There are, at all times, two different estimates of the bias adjustments,

844   which are in general agreement and give a first indication of the structural uncertainty.

845   Evidence for the efficacy of the adjustments comes from wind tunnel tests, comparisons

846   with coastal sites and consistency with subsurface ocean temperatures, marine air

847    temperatures and land air temperatures. Contrary evidence comes from a recent field

848    experiment in the Pacific. Uncertainty could be better understood by: improvements in

849    metadata; carefully designed fields tests of buckets and other measurements methods; the

850    creation of new independent evaluations of the biases; and continued comparison

851    between SST and related variables.

852

853    **3.4 Sampling Uncertainty**

854

855    The magnitude of the grid-box sampling uncertainty depends on the correlation and

856    variability of SSTs within the grid box, on the number of observations contributing to the

857    grid-box average and where in the grid box they are located. High average correlations

858    within a grid box, low variability and large numbers of observations lead to lower

859    uncertainty estimates. Conversely areas of high variability or low average correlation,

860    such as frontal regions or western boundary currents, tend to have higher grid-box

861    sampling uncertainties as do grid-box averages based on smaller numbers of

862    observations. The estimation of uncertainties arising from the sparseness of observations

863    at scales from grid box level to global has been approached in a number of ways.

864

865    **3.4.1 Grid-box Sampling Uncertainty**

866

867    *Weare and Strub* [1981] counted the number of observations needed to minimize

868    sampling uncertainty in a 5°x5° grid box by ensuring that the observations were evenly

869    split between all areas of the grid box, month and diurnal cycle. From this, they

870    concluded that even sampling could not be achieved with fewer than eleven observations,

871    but that in practice more than eleven, sometimes many more, would be needed.

872

873    *Rayner et al.* [2006] estimated a combined measurement and grid-box sampling

874    uncertainty by considering how the variance of the grid-box average changed as a

875    function of the number of observations. The technique picked up spatial variations in

876    grid-box sampling uncertainty associated with regions of high variability. *Rayner et al.*

877    [2009] showed results from an unpublished analysis by Kaplan, in which spatially

878    complete satellite data were used to estimate the variability within 1°x1° grid boxes. The

879    same features were seen as in the *Rayner et al.* [2006] analysis, allowing for differences

880    in resolution, although the uncertainties estimated by Kaplan tended to be higher. *She et*

881    *al.* [2007] also used sub-sampling of satellite data to estimate grid-box sampling

882    uncertainty for the Baltic Sea and North Sea. *Kent and Berry* [2005] showed that

883    separately assessing measurement and sampling uncertainties can help to decide whether

884    more, or better, observations are needed to reduce the average uncertainty in an

885    individual grid box.

886

887    *Morrissey and Greene* [2009] developed a theoretical model for estimating grid-box

888    sampling uncertainty that accounted for non-random sampling within a grid box. This

889    was an extension of the method used to estimate sampling uncertainties in land

890    temperature data and global temperatures by *Jones et al.* [1997]. Land temperatures are

891    measured by stations at fixed locations that take measurements every day. Marine

892    temperature measurements are taken at fixed times, but the ships and drifting buoys move

893 during a particular month. *Morrissey and Greene* [2009] do not provide a practical

894 implementation of their approach, only a theoretical framework. *Kennedy et al.* [2011b]

895 extended the concept of the average correlation within a grid box developed in *Jones et*

896 *al.* [1997] to incorporate a time dimension. *Kent and Berry* [2008] used a temporal

897 autocorrelation model that took account of the days within the period that were sampled,

898 and the days which were not, to estimate the temporal sampling uncertainty. An

899 alternative to the *Jones et al.* [1997] method for land data was provided by *Shen et al.*

900 [2007], but it has not yet been applied in SST analyses.

901

902 It is possible that the locations visited by ships and drifting buoys are related and, to an

903 extent, dictated by meteorological and oceanographic conditions. Ships have long used

904 the prevailing currents in the Atlantic to speed their progress and it is in the interest of

905 almost all shipping to steer clear of hurricanes and other foul weather. Bad weather is

906 also likely to have influenced how and when observations were made. Conversely, the

907 conditions in which a sail ship might become becalmed could lead to over sampling of

908 higher SSTs. Drifting buoys drift, and a drifter trapped in an eddy might persistently

909 measure temperatures that are representative of only a very limited area. Drifters also

910 tend to drift out of areas of upwelling and congregate in other areas.

911

912 The effect of uneven sampling can be reduced by the creation of 'super observations'

913 during the gridding process [*Rayner et al.,* 2006], or data preparation stage [*Ishii et al.,*

914 2003], but such processes cannot readily account for the situations where no observations

915 are made at all.

916

917    As noted by *Rayner et al.* [2006], the grid-box sampling uncertainties are likely to be

918    uncorrelated or only weakly correlated between grid boxes so the effect of averaging

919    together many grid boxes will be to reduce the combined grid-box sampling uncertainty

920    by a factor proportional to the square root of the number of grid boxes. Consequently the

921    sampling component of the uncertainty will be of minor importance in the global annual

922    average (Figure 8).

923

924    **3.4.2 Large-scale Sampling Uncertainty**

925

926    Because *Rayner et al.* [2006] and *Kennedy et al.* [2011b] make no attempt to estimate

927    temperatures in grid boxes which contain no observations, an additional uncertainty had

928    to be computed when estimating area-averages. *Rayner et al.* [2006] used Optimal

929    Averaging (OA) as described in *Folland et al.* [2001] which estimates the area average in

930    a statistically optimal way and provides an estimate of the large-scale sampling

931    uncertainty. *Kennedy et al.* [2011b] subsampled globally complete fields taken from three

932    SST analyses and obtained similar uncertainties from each. The uncertainties of the

933    global averages computed by *Kennedy et al.* [2011b] were generally larger than those

934    estimated by *Rayner et al.* [2006]. *Palmer and Brohan* [2011] used an empirical method

935    based on that employed for grid-box averages in *Rayner et al.* [2006] to estimate global

936    and ocean basin averages of subsurface temperatures.

937

938    The *Kennedy et al.* [2011b] large-scale sampling uncertainty of the global average SST

939    anomaly is largest (with a 2-sigma uncertainty of around 0.15°C) in the 1860s when

940    coverage was at its worst (Figure 8). This falls to 0.03 °C by 2006. The fact that the

941    large-scale sampling uncertainty should be so small – particularly in the nineteenth

942    century – may be surprising. The relatively small uncertainty might simply be a reflection

943    of the assumptions made in the analyses used by *Kennedy et al.* [2011b] to estimate the

944    large-scale sampling uncertainty. Indeed, *Gouretski et al.* [2012] found that subsampling

945    an ocean reanalysis underestimated the uncertainty when the coverage was very sparse.

946    However, estimates made by *Jones* [1994] suggest that a hemispheric-average land-

947    surface air temperature series might be constructed using as few as a 109 stations. For

948    SST, the variability is typically much lower than for land temperatures though the area is

949    larger. It seems likely that the number of stations needed to make a reliable estimate of

950    the global average SST anomaly would not be vastly greater.

951

952    Another way of assessing the large-scale sampling uncertainty is to look at the effect of

953    reducing the coverage of well-sampled periods to that of the less-well-sampled nineteenth

954    century and recomputing the global average (see for example *Parker* [1987]). Figure 4

955    shows the range of global annual average SST anomalies obtained by reducing each year

956    to the coverage of years in the nineteenth century. So, for example, the range indicated by

957    the blue area in the upper panel for 2006 shows the range of global annual averages

958    obtained by reducing the coverage of 2006 successively to that of 1850, 1851, 1852... and

959    so on to 1899. The red line shows the global average SST anomaly from data that have

960    not been reduced in coverage. For most years, the difference between the sub-sampled

961 and more fully sampled data is smaller than 0.15K and the largest deviations are smaller

962 than 0.2K. For the large-scale sampling uncertainty of the global average to be

963 significantly larger would require the variability in the nineteenth century data gaps to be

964 different from that in the better-observed period.

965

966 **3.4.3 Summary of Sampling Uncertainty**

967

968 Uncertainties arising from under-sampling at a grid-box level are easy to assess if the

969 observations can be assumed to be randomly distributed within a grid box. However,

970 sampling is not random. The effect of this is reduced in most analyses by the calculation

971 of super-observations that combine nearby measurements; however, optimal methods to

972 minimize uncertainty are not generally applied. Simple estimates of large-scale sampling

973 uncertainty in the global-average SST from subsampling well-sampled periods suggest a

974 value of at most 0.2K even in poorly observed years. However, there are potential

975 limitations of these simple methods and they should be considered together with the

976 range of statistical reconstructions to get a more complete idea of uncertainty in large-

977 scale averages.

978

979 **3.5 Reconstruction Techniques and Other Structural Choices**

980

981 Creating global SST analyses is challenging because of the relative sparseness of

982 observations before the satellite era and the non-stationarity of the changing climate. A

983 large number of different SST data sets based on in situ data have been produced

984 employing a variety of statistical methods. The structural uncertainties associated with

985 estimating SSTs in data voids and at data-sparse times are therefore somewhat better

986 explored than structural uncertainties in the pervasive systematic errors. Data sets used in

987 this paper have been summarized in Table 5 and global averages for these data sets are

988 shown in Figure 5.

989

990 **3.5.1 Critique of Reconstruction Techniques**

991

992 The current generation of SST analyses are the survivors of an evolutionary process

993 during which less effective techniques were discarded in favor of better adapted

994 alternatives. It is worthwhile to ask how, as a group, they address the range of criticisms

995 that have arisen during that time.

996

997 One concern is that patterns of variability in the modern era which are used to estimate

998 the parameters of the statistical models might not faithfully represent variability at earlier

999 times [*Hurrell and Trenberth*, 1999]. The concern is allayed somewhat by the range of

1000 approaches taken. The method of *Kaplan et al.* [1998] which uses the modern period to

1001 define Empirical Orthogonal Functions (EOFs, see *Hannachi et al.*, [2007] for a review

1002 of the use of EOFs in the atmospheric sciences) tends to underestimate the long-term

1003 trend. This is particularly obvious in the nineteenth and early twentieth century. *Rayner et*

1004 *al.* [2003] extended the method by defining a low-frequency, large-scale EOF that better

1005 captured the long-term trend in the data. However, it is possible that a single EOF will

1006 fail to capture all the low-frequency changes. *Smith et al.* [2008] allow for a non-

1007     stationary low-frequency component in their analysis which contributes a large

1008     component of uncertainty in the early record, but their reconstruction reproduces less

1009     high-frequency variability at data-sparse epochs. *Ilin and Kaplan* [2009] and *Luttinen and*

1010     *Ilin* [2009, 2012] used algorithms that make use of data throughout the record to estimate

1011     the covariance structures and other parameters of their statistical models. The three

1012     algorithms use either large-scale patterns (VBPCA, GPFA) or local correlations (GP).

1013     Differences between the three methods are generally small at the global level, but they

1014     diverge during the 1860s when data are few. There is a caveat that despite using all

1015     available observations, such methods will still tend to give a greater weight to periods

1016     with more plentiful observations. *Ishii et al.* [2005] use a simply-parameterized local

1017     covariance function for interpolation. Their optimal interpolation (OI) method was

1018     assessed by *Hirahara et al.* [2013] to have larger analysis uncertainties and larger cross-

1019     validation errors than the EOF-based COBE-2 analysis. However, the use of a simple

1020     optimal interpolation method has the advantage that it makes fewer assumptions

1021     regarding the stationarity of large-scale variability.

1022

1023     Another concern is that methods that use EOFs to describe the variability might

1024     inadvertently impose spurious long-range teleconnections that do not exist in the real

1025     world [*Dommenget*, 2007]. *Smith et al.* [2008] explicitly limit the range across which

1026     teleconnections can act. *Ishii et al.* [2005] used a local covariance structure in their

1027     analysis. Analyses such as *Kaplan et al.* [1998] and *Rayner et al.* [2003] make the

1028     assumption that the EOFs retained in the analysis capture actual variability in the SST

1029     fields, but do not explicitly differentiate between variability that can be characterized

1030   purely in terms of local co-variability and large-scale teleconnections. *Karspeck et al.*

1031   [2012] note that there is not a clear separation of scales and that joint estimation of local

1032   and large scale covariances is the logical way to approach the problem.

1033

1034   Most, if not all, statistical methods have a tendency to lose variance either because they

1035   do not explicitly resolve small scale processes [*Kaplan et al.,* 1998; *Smith et al.*, 2008],

1036   because the method tends towards the climatological average in the absence of data [*Ishii*

1037   *et al.,* 2005; *Berry and Kent*, 2011], or because they tend to smooth the data. *Rayner et al.*

1038   [2003] used the method of *Kaplan et al.* [1998] but blended high-quality gridded

1039   averages back into the reconstructed fields to improve small scale variability where

1040   observations were plentiful. *Karspeck et al.* [2012] analyzed the residual difference

1041   between the observations and the analysis of *Kaplan et al.* [1998] analysis using local

1042   non-stationary covariances, and then drew a range of samples from the posterior

1043   distribution in order to provide consistent variance at all times and locations.

1044

1045   One assumption common to most of the above analysis methods is that SST variability

1046   can be decomposed into a small set of distinct patterns that can be combined linearly to

1047   describe any SST field. However, it is well known that phenomena such as El Niño and

1048   La Niña are not symmetric and that the equations that describe the evolution of SST are

1049   non-linear. Consequently, current analyses might not capture the full range of behavior in

1050   real SST fields [*Karnauskas,* 2013]. Current generation SST analyses are based on the

1051   assumption that individual measurement errors are uncorrelated and that errors are

1052   normally distributed. Analysis techniques that incorporate information about the

1053 correlation structure of the errors have not yet been developed. Such techniques are likely

1054 to be more computationally expensive and lead to larger analysis uncertainties.

1055

1056 **3.5.2 Other Structural Choices**

1057

1058 Analyses based on SST anomalies will also have an uncertainty associated with the

1059 climatological reference fields used to calculate the anomalies. Sub-surface analyses have

1060 been shown to be particularly sensitive to choice of base period [*Lyman et al.*, 2010], due

1061 in a large part to the relative sparseness of the data sets. Although the problem is likely to

1062 be less severe for the better-observed SST record, there are still regions – the Southern

1063 Ocean and Arctic Ocean – where observations are few. *Yasunaka and Hanawa* [2011]

1064 found that differences between long-term-average SSTs from different analyses were

1065 typically less than 0.5 K, but that they exceeded 1 K in places. The largest differences

1066 were at high latitudes and in regions with strong SST gradients. There are also likely to

1067 be pervasive systematic errors in the climatological averages [*Kennedy et al.*, 2011c].

1068

1069 Other structural differences arise from the way that SSTs are extended to the edge of the

1070 sea ice. SSTs can be estimated from measurements of sea-ice concentration [*Rayner et*

1071 *al.*, 2003; *Smith et al.*, 2008; *Hirahara et al.*, 2013]. Although their global impact is

1072 likely to be small, the uncertainties in these relationships and estimates need also to be

1073 factored into the uncertainty of SSTs in these regions. At the moment, the uncertainty

1074 associated with historical sea-ice concentrations is poorly understood.

1075

1076     **3.5.3 Comparisons of Reconstructions**

1077

1078     *Yasunaka and Hanawa* [2011] examined a range of climate indices based on seven

1079     different SST data sets. They found that the disagreement between data sets was marked

1080     before 1880, and that the trends in large scale averages and indices tend to diverge

1081     outside of the common climatology period. For the global average, the differences

1082     between analyses were around 0.2 K before 1920 and around 0.1-0.2 K in the modern

1083     period. Even for relatively well-observed events such as the 1925/26 El Niño, the detailed

1084     evolution of the SSTs in the tropical Pacific varied from analysis to analysis. The reasons

1085     for the discrepancies are not completely clear because each data set is based on a slightly

1086     different set of observations that have been quality controlled and processed in different

1087     ways, a problem that could be alleviated by running analyses on identical input data sets.

1088

1089     Combined with information about large-scale sampling uncertainties estimated in other

1090     ways, the spread between analyses suggests that the large-scale sampling uncertainty in

1091     global average SST anomaly is around 0.2 K in the late 19[th] century. For the large-scale

1092     sampling uncertainty of the global average to be much larger would require variability in

1093     the early record to have been different from variability in the modern period, which is a

1094     possibility. The resolution of such a question is most likely to be achieved via the

1095     digitisation of more observations from paper records.

1096

1097     Progress in assessing the differences between analysis techniques can also be made by

1098     studying the relative strengths and weaknesses of interpolation techniques on carefully

1099     prepared test data sets using synthetic data, or on 'withheld' data from well observed

1100     regions. By running each analysis on the same carefully-defined subsets and tests, it

1101     should be possible to isolate reasons for the differences between the analyses and assess

1102     the reliability of analysis uncertainty estimates. The International Surface Temperature

1103     Initiative (http://www.surfacetemperatures.org/) has been working on such benchmarking

1104     exercises for land surface air temperature data, building on work such as the COST

1105     ACTION project [*Venema et al.*, 2012].

1106

1107     **3.5.4 Summary of Reconstruction Techniques and Structural Uncertainty**

1108

1109     A range of reconstruction techniques exist to make globally-complete or near globally-

1110     complete SST analyses. The spread in global mean SST between analyses is at worst

1111     around 0.2K. The analyses are based on a variety of different statistical models

1112     suggesting that estimates of global average SST are not strongly dependent on such

1113     choices. However, current reconstruction techniques do not account for systematic errors

1114     in the data – they assume errors are random and uncorrelated – and assume that SST

1115     fields can be simply parameterized in terms of limited numbers of patterns or simple

1116     covariance relationships. Objective comparison of different reconstruction techniques and

1117     their associated uncertainty estimates would be aided by the creation of standard

1118     benchmark tests which mimic the distribution and character of observational data.

1119

1120     **3.6 Comparing Components of Uncertainty**

1121

1122     Figure 7 shows individual components of the overall uncertainty estimated for three

1123     months. The components include: estimates of structural uncertainty (in lieu of a formal

1124     way to estimate this, it is calculated as the standard deviation of seven near-globally-

1125     complete analyses: COBE, Kaplan, ERSSTv3, HadISST, GPFA, GP and VBPCA),

1126     sampling uncertainty, combined random and systematic measurement error uncertainty,

1127     bias uncertainty (estimated from a 200-member ensemble described in section 4) and

1128     analysis uncertainties from ERSSTv3b [*Smith et al.* 2008].

1129

1130     At a monthly, grid-box level, the parametric uncertainty in the *Kennedy et al.* [2011c]

1131     systematic error estimates is typically the smallest uncertainty and is nearly always less

1132     than 0.2 K. The sampling uncertainty and measurement uncertainty both depend on the

1133     number of observations, so they are larger in areas with fewer observations. Of the two,

1134     measurement uncertainty is typically larger.

1135

1136     In well-observed periods, the spread between the different analyses is roughly what one

1137     might expect: closer agreement in well-observed regions, poorer agreement in data-sparse

1138     regions, principally the Southern Ocean and Arctic Ocean. At more poorly-observed

1139     times, the spread between analyses is narrower than the climatological standard deviation

1140     suggesting that the reconstructions are skilful in the sense that they are providing useful

1141     information in data voids. However, the narrow spread is in contrast to those areas where

1142     there have been changes in the input observations (see, for example, the Indian Ocean in

1143     Figure 7(b) and Figure 7(h)). A small number of observations, which are available to one

1144     analysis but not another, lead to a larger spread than is seen in data-free regions implying

1145   that, while there is diversity in the approaches, there may still be too little for the best

1146   estimates alone to effectively bracket the true uncertainty range.

1147

1148   The ERSSTv3 analysis uncertainties are largest in regions where there are consistent data

1149   voids. They show a similar pattern to the structural uncertainty estimate in 1944 and

1150   2003, but there is marked difference in 1891, with the analysis uncertainty being larger

1151   than the structural uncertainty in the poorly-observed western Pacific.

1152

1153   Figure 8 shows time series of the different components of uncertainty at different spatial

1154   scales from global to grid box. The bias uncertainty is relatively constant and is the

1155   smallest component of uncertainty at the grid box level for much of the record. The

1156   sampling uncertainty for a grid box is larger than the bias uncertainty when observations

1157   are few, but in the recent record they are comparable. In this example, the measurement

1158   uncertainty is larger than bias and sampling uncertainties at the grid box level, even when

1159   observations are numerous. However, in other grid boxes, characterised by strong SST

1160   gradients or high variability, such as the western boundary currents, the sampling

1161   uncertainty could be larger.

1162

1163   As the size of the area increases and more observations are included in the average, the

1164   sampling and measurement uncertainties decrease. Two estimates of the measurement

1165   uncertainty are included. In one, correlations between individual errors are taken into

1166   account. In the other, measurement errors are assumed to be random and independent. In

1167   the latter case, the measurement uncertainties become small relative to other sources of

1168    uncertainty at a basin scale early in the 20<sup>th</sup> century. However, the effect of correlated

1169    errors is such that measurement uncertainty remains a major source of uncertainty at all

1170    scales until the 1980s when the global VOS fleet reached its peak and the deployment of

1171    drifting and moored buoys began.

1172

1173    The largest component at the scales shown here is the structural uncertainty. In the grid

1174    box shown, the structural uncertainty is, at times, larger than the combined uncertainty

1175    from other components suggesting that some or all of the analyses are losing information.

1176    At a global level, where estimated analysis uncertainties are available for COBE, COBE-

1177    2, Kaplan and ERSSTv3b data sets, the structural uncertainty is comparable to the

1178    estimated analysis uncertainties. For example, in 1900, the ERSSTv3b analysis

1179    uncertainty is 0.03K, the COBE analysis uncertainty is 0.06K, COBE-2 gives 0.05K and

1180    Kaplan is around 0.05K.

1181

1182    Because of the nature of the uncertainties arising from the adjustments for pervasive

1183    systematic errors, the uncertainties become relatively more important as the averaging

1184    scale increases. At a global scale, bias uncertainties are comparable to or larger than all

1185    other uncertainty components from the 1940s to the present. There is a caveat: because

1186    the SSTs are expressed as anomalies, the size of the bias uncertainty depends on the base

1187    period used to calculate the anomalies. In Figure 8, the period used is 1961-1990, which

1188    is why there is a local minimum in the bias uncertainty centred on that period.

1189

1190    **3.7 Estimates of Total Uncertainty**

1191

1192    *Smith and Reynolds* [2005] attempted to combine all the different uncertainties described

1193    above to get a total uncertainty estimate. They combined their analysis uncertainty with

1194    measurement uncertainty, bias uncertainty and structural uncertainty. Uncertainty

1195    associated with pervasive systematic errors and structural uncertainty in the adjustments

1196    were estimated by taking the mean squared difference between the *Smith and Reynolds*

1197    [2002] and *Folland and Parker* [1995] bias adjustments in the prewar period. After

1198    World War 2, the bias uncertainty was estimated by calculating the average difference

1199    between engine room measurements and all measurements. Structural uncertainties were

1200    estimated by analysing the spread of three SST analyses.

1201

1202    Figure 9 shows the total uncertainty estimate from the latest version of the ERSST

1203    analysis, ERSSTv3b, in blue. A similar estimate was made based on the HadSST3 data

1204    set in the following way. Measurement uncertainties, grid-box sampling uncertainties and

1205    large-scale sampling uncertainties were estimated using the method of *Kennedy et al.*

1206    [2011b, 2011c]. To estimate the uncertainty associated with pervasive systematic

1207    errors,an ensemble of 200 data sets, comprising the 100 original ensemble members from

1208    HadSST3 and a 100-member ensemble generated by replacing the *Rayner et al.* [2006]

1209    bucket-correction fields with the fields from *Smith and Reynolds* [2002]. The adjustment

1210    uncertainties on individual months were assumed to be correlated within a year, giving a

1211    greater uncertainty range than in *Kennedy et al.* [2011c], particularly before 1941. During

1212    the war years 0.2 K was added to reflect the additional uncertainty during that period as

1213    described by *Kennedy et al.* [2011c]. As above, structural uncertainties were estimated by

1214    taking the standard deviation of area-average time series from seven analyses.

1215

1216    The total uncertainty estimates from these two assessments are comparable between 1880

1217    and 1915. Between 1915 and 1941, the ERSSTv3b uncertainty estimate is larger because

1218    the estimated bias uncertainty is larger. The difference is most obvious in the northern

1219    hemisphere where the differences between the *Smith and Reynolds* [2002] and *Folland*

1220    *and Parker* [1995] bias adjustments are largest. From 1941 to present, the HadSST3-

1221    based uncertainty estimate is the larger because the bias uncertainty is larger than in

1222    ERSSTv3b.

1223

1224    The obvious question that arises is "do these assessments span the full uncertainty

1225    range?" In this case, it probably pays to err on the side of caution. Although the structural

1226    uncertainty is based on a range of methods for infilling missing data, there are still

1227    commonalities in the approaches taken and there is little diversity in the approaches to

1228    bias adjustment. The lack of diversity is troubling because the differences between the

1229    median estimates of HadSST3 and ERSSTv3b are greater than the estimated uncertainties

1230    of the ERSSTv3b analysis at times during the period 1950-1970 suggesting that the

1231    uncertainties may have been underestimated in the earlier assessment.

1232

1233    **4 Presentation of Uncertainty**

1234

1235 At present, some groups provide explicit uncertainty estimates based on their analysis

1236 techniques [*Kaplan et al.,* 1998; *Smith et al.*, 2008; *Kennedy et al.*, 2011b, 2011c, *Ishii et*

1237 *al.,* 2005; *Hirahara et al.,* 2013]. The uncertainty estimates derived from a particular

1238 analysis will tend to misestimate the true uncertainty because they rely on the analysis

1239 method and the assumptions on which it is based being correct.

1240

1241 Comparing uncertainty estimates provided with analyses can be difficult because not all

1242 analyses consider the same sources of uncertainties. Consequently, a narrower

1243 uncertainty range does not necessarily imply a better analysis. One way that data set

1244 providers could help users is to provide an inventory of sources of uncertainty that have

1245 been considered either explicitly or implicitly. This would allow users to assess the

1246 relative maturity of the uncertainty analysis.

1247

1248 There is a further difficulty in supplying and using uncertainty estimates: the traditional

1249 means of displaying uncertainties – the error bar, or error range – does not preserve the

1250 covariance structure of the uncertainties. Unfortunately, storing covariance information

1251 for all but the lowest resolution data sets can be prohibitively expensive. EOF-based

1252 analyses, like that of *Kaplan et al.* [1998], could in principle efficiently store the spatial-

1253 error covariances because only the covariances of the reduced space of principal

1254 components need to be kept. For *Kaplan et al.* [1998], based on a reduced space of only

1255 80 EOFs, this is a matrix of order $80^2$ elements for each time step as opposed to $1000^2$

1256 elements for the full-field covariance matrix. The difficulty with this approach is that not

1257    all variability can be resolved by the leading EOFs and excluding higher-order EOFs will

1258    underestimate the full uncertainty.

1259

1260    *Karspeck et al.* [2012] drew samples from the posterior probability produced by their

1261    analysis. Each sample provides an SST field that is consistent with the available

1262    observations and the estimated covariance structure. Sampling has the added advantage

1263    that it can be combined easily with Monte-Carlo samples from the measurement bias

1264    distributions. However, production of samples is not always computationally efficient.

1265    *Karspeck et al.* [2012] were able to do it for the North Atlantic region, but the

1266    computational costs of extending the analysis unchanged to the rest of the world could be

1267    prohibitive. *Kennedy et al.* [2011c] provided an ensemble of 100 interchangeable

1268    realizations of their bias-adjusted data set, HadSST3. The ensemble spans parametric

1269    uncertainties in their adjustment method.

1270

1271    By providing a set of plausible realizations of a data set, or alternatively by providing

1272    plausible realizations of typical measurement errors [*Mears et al.*, 2011], it can be

1273    relatively easy for users to assess the sensitivity of their analysis to uncertainties in SST

1274    data. For example, individual ensemble members of HadSST3 were used in *Tokinaga et*

1275    *al.* [2012], along with other SST analyses, to show that their results were robust to the

1276    estimated bias uncertainties in SSTs.

1277

1278    Another approach [*Merchant et al.* 2013] is to separate out components of the uncertainty

1279    that correlate at different scales. Random measurement errors, such as sensor noise, are

1280 uncorrelated. Some uncertainties, for example those related to water vapor in a satellite

1281 view, are correlated at a synoptic scale. Yet others are correlated at all times and places.

1282 Grouping uncertainties in this way allows users to propagate uncertainty information

1283 more easily.

1284

1285 **5 Minimizing Exposure to Uncertainty**

1286

1287 Alternative approaches to using the SST data in a way that is less sensitive to biases and

1288 other data errors have been made. The following approaches make use of knowledge

1289 concerning the types of errors and uncertainties found in SST data and have been adapted

1290 to account for them. They highlight the importance of combining understanding of the

1291 measurements and their potential errors, as well as understanding of the phenomenon

1292 being analyzed. Perhaps the simplest example is *Schell* [1959] who suggested discarding

1293 grid-box averages (in that case Marsden squares) based on small numbers of

1294 observations.

1295

1296 *Thompson et al.* [2008] identified an abrupt drop in the observed global average SST

1297 anomaly in late 1945, which they attributed to a rapid change in the composition of

1298 ICOADS 2.0 [*Worley et al.*, 2005] from mostly US ships immediately before the 1945

1299 drop to mostly UK ships immediately afterwards. This hypothesis was lent further weight

1300 by *Kennedy et al.* [2011c]. In a follow-up paper [*Thompson et al.*, 2010], a drop in

1301 northern-hemisphere SSTs was identified. In order to show that the drop was not an

1302 artifact of the change in measurement method, they divided the ICOADS data into

1303     distinct subsets based on the country of the ships making the measurements, considered a

1304     range of different SST analyses, and looked at related variables such as NMAT and land

1305     surface air temperatures. The probability of a drop being due to a coincident change in

1306     the way that all countries measured SST, simultaneous with a sudden change in NMAT

1307     and land temperature bias, is small. The fact that the drop was seen in all the different

1308     data sets implied that the drop was real. *Tokinaga et al.* [2012] took a similar approach,

1309     using bucket measurements from ICOADS as a quasi-homogeneous estimate of SST

1310     change over the period 1950 to 2009.

1311

1312     In detection and attribution studies it is common to reduce the coverage of the models to

1313     match that of the data. Doing so reduces the exposure of the study to uncertainties

1314     associated with interpolation techniques, but it does not avoid the problem of systematic

1315     biases. Recent studies [*Jones and Stott*, 2011] have explicitly used a range of data sets to

1316     start to map out the effects of structural uncertainties on detection and attribution studies.

1317

1318     SST data sets are routinely compared to the output of climate simulations. Bearing in

1319     mind the discussion in section 2 on the definition of SST it might be necessary to ensure

1320     that the modeled output and the measured SST correspond to the same quantity. Many

1321     climate models employ a surface mixed layer that is several meters thick. However,

1322     models have been run with greater resolution in the near-surface ocean [e.g., *Bernie et al.,*

1323     2008] in order to simulate diurnal variability.

1324

1325     Another common use of SST data for which an understanding of the limitations of the

1326     data is important is in the calculation and interpretation of EOFs. In many studies EOFs

1327     are calculated from globally complete SST analyses because the lack of missing data

1328     makes calculating EOFs easy. However, it seems wise to bear in mind that a good deal of

1329     statistical processing has already been applied to the SST analyses to make them globally

1330     complete. Extracting EOFs from (or applying any other analysis technique to) what are in

1331     some cases EOF analyses already, could lead to difficulties of interpretation on top of the

1332     more general problems [*Hannachi et al.,* 2007; *Dommenget,* 2007; *Karnauskas,* 2013].

1333     Techniques exist for estimating EOFs from gridded data sets with missing data and these

1334     can also incorporate uncertainty information though many assume that the errors are

1335     uncorrelated and will tend to underestimate uncertainty in the EOFs and their principal

1336     components. See for example, *Roweis* [1998], *Schneider* [2001], *Beckers and Rixen*

1337     [2003], *Rutherford et al.* [2004], *Houseago-Stokes and Challenor* [2004], *Kondrashov*

1338     *and Ghil* [2006], *Ilin and Kaplan* [2009] and *Luttinen and Ilin* [*2009*].

1339

1340     **6 Satellites**

1341

1342     Although the present review is principally concerned with in situ measurements of SST it

1343     is necessary to mention the important role that satellite data play in understanding SST

1344     variability and uncertainty. The advantages of satellite data are obvious; particularly the

1345     ability to measure large areas of the ocean using a single instrument, giving a more nearly

1346     global view of SST.

1347

1348    However, the first thing to note is that satellites monitor radiances and do not directly

1349    measure SSTs. The measured radiances are affected by the state and constituents of the

1350    atmosphere as well as variations in the state and temperature of the sea-surface. The

1351    wavelengths that are sampled are set by the design of the instrument. Retrieving SST

1352    from the radiances is a difficult inverse process and sensitive to biases and other errors

1353    [*Merchant et al.* 2008b]. The second thing to note is that satellite instruments are

1354    sensitive to the skin (upper few microns), or sub-skin (upper few millimeters)

1355    temperature depending on the wavelengths measured by the satellite. Because satellite

1356    instruments are sensitive to the topmost layer of the ocean, the diurnal range of retrieved

1357    SSTs is larger than for measurements made at depth. Thirdly, accurate SST retrievals

1358    from infra-red instruments are only possible when the view is not obscured by cloud.

1359    Microwave retrievals can penetrate cloud, but suffer from problems near to coastlines,

1360    and where precipitation rates are high. They also have coarser spatial resolution and

1361    higher measurement uncertainties [*O'Carroll et al.*, 2008].

1362

1363    The longest records of SST from satellite are derived from the AVHRR (Advanced Very

1364    High Resolution Radiometer) instruments. These instruments make nadir measurements

1365    using two infra-red channels. The retrievals are usually calibrated relative to in situ data.

1366    More recent re-processings use optimal estimation to obtain a retrieval that is

1367    independent of the in situ record [*Merchant et al.*, 2008b] but these have not yet been

1368    extended to calculating global averages. Furthermore, the AVHRR instrument is prone to

1369    systematic errors caused by aerosols in the atmosphere and the satellite orbits drift slowly

1370    altering the sampling of the diurnal cycle through time. Despite the numerous

1371    shortcomings of the AVHRR record, *Good et al.* [2007] showed that there was a long-

1372    term warming trend in SSTs as measured by AVHRR.

1373

1374    The Along-Track Scanning Radiometers (ATSR) [*Smith et al.*, 2012] were designed to

1375    meet the needs of climate monitoring. The satellite is a dual view instrument, taking nadir

1376    and forward views using three infra-red channels. The dual view configuration allows for

1377    more effective screening of aerosols and the three channels allow for accurate retrievals

1378    across a wider range of conditions. Furthermore, the onboard calibration system allows

1379    the stability of the radiance measurements from the instrument to be maintained. The

1380    ATSR data have been reprocessed in the ATSR Reanalysis for Climate (ARC) project

1381    [*Merchant et al.*, 2008a] and the resulting time series have been shown to have biases of

1382    less than 0.1 K and stability better than 5 mK/year since 1993 in the tropics where

1383    reliable long term moorings can be found [*Embury et al.*, 2012; *Merchant et al.*, 2012].

1384    The ARC reprocessing is almost independent of the in situ network therefore it can be

1385    used to corroborate trends seen in the in situ network. In a comparison between global

1386    average SST anomalies (at a nominal depth of 0.2 m) calculated using the ARC data and

1387    HadSST3, the two time series agree within the estimated HadSST3 uncertainties except

1388    for parts of the ATSR1 record in the early 1990s. The ATSR1 period is believed to be of

1389    lower quality as a result of the failure of one of the IR channels, failure of the satellite

1390    cooling system as well as the high stratospheric aerosol loadings following the eruption

1391    of Mount Pinatubo in 1991.

1392

1393     The nearly global, high-resolution view of the world's oceans provided by satellite

1394     instruments can be used as a way of improving and testing many aspects of SST analysis.

1395     By combining the more detailed fields produced by satellites with the long records of in

1396     situ measurements, more detailed reconstructions are possible over a wider area of the

1397     Earth [*Rayner et al.*, 2003; *Smith et al.*, 2008; *Hirahara et al.,* 2013]. Satellite data can

1398     also be used to assess the verisimilitude of reconstructions based on sparser in situ data.

1399

1400     **7 Concluding Remarks and Future Directions**

1401

1402     One of the chief difficulties in assessing the uncertainties in SST data sets is the

1403     impossibility of tracing individual observations back via an unbroken chain to

1404     international measurement standards. The creation of a global array of reference stations

1405     each making simultaneous redundant measurements of a variety of marine variables

1406     could solve some of the problems of SST analysis that have bedeviled the understanding

1407     of historical SST change and would provide a gold standard against which the *future*

1408     wider observing system – incorporating observations from ships, buoys, profiling floats

1409     and satellites – can be assessed. Even without such traceability a climate record could be

1410     more easily maintained by stricter adherence to the Global Climate Observing System

1411     [*GCOS* 2003] climate monitoring principles.

1412

1413     In the absence of such a network the estimation of uncertainties has depended heavily on

1414     redundancies in measurement systems and in analysis techniques. Full use of the

1415     redundancies is now being made in the modern period via comparisons of the many

1416    available satellite sources with each other and with in situ sources [*O'Carroll et al.*, 2008;

1417    *Merchant et al.*, 2012] and sub-surface data [*Gille*, 2012]. Analyses that ingest a variety

1418    of data sources can produce bias statistics for each of the inputs [*Brasnett*, 2008; *Xu and*

1419    *Ignatov*, 2010]. Such information can be exploited to assess their relative quality and, as

1420    the analyses are pushed further back in time [*Roberts-Jones et al.*, 2012], they will help

1421    assess uncertainties through a larger part of the record.

1422

1423    SSTs are physically related to other measurements including surface pressures and winds,

1424    salinity, air temperatures, sub-surface temperatures and ocean biology amongst others.

1425    Information from SST can be supplemented by analyses based on physical understanding

1426    of the climate system. It has already been shown that by combining information from

1427    night marine air temperatures with SST it was possible to greatly reduce uncertainties in

1428    early 20$^{th}$ and late 19$^{th}$ century SST. *Yu et al.* [2004] used a joint estimation method to

1429    minimize uncertainties in flux estimates based on a range of different variables mostly

1430    based on satellite data. Other studies [*Tung and Zhou*, 2010; *Deser et al.*, 2010] have used

1431    physical reasoning based on a host of variables to explore uncertainties in the long-term

1432    trends of tropical Pacific SSTs first raised by *Vecchi et al.* [2008]. It has even been

1433    suggested that proxy records such as isotope ratios from corals and ice cores could be

1434    used, with appropriate care, to understand uncertainties in the longest-term changes in

1435    SST [*Anderson et al.*, 2013]. The most advanced exemplars of physical and statistical

1436    synthesis are ocean and coupled reanalyses which will play an increasingly important role

1437    in understanding observational uncertainty and long-term climate change.

1438

1439      A key barrier to understanding SST uncertainty is a lack of appropriate metadata. Better

1440      information is needed concerning how measurements were made, which method was

1441      used to make a particular observation, calibration information, the depths at which

1442      observations were made, and even basic information such as the call sign or name of the

1443      ship that made a particular observation.

1444

1445      Some of this information can be inferred from data already contained in marine reports.

1446      Where reports in ICOADS cannot be associated with a particular ship, either because

1447      they have a missing ID, or a generic ID, there is much to be gained by grouping

1448      observations to give plausible ship tracks, or voyages. By using data association

1449      techniques to infer such metadata from the location information and other clues such as

1450      how frequently observations were made and which variables were observed, it should be

1451      possible to assess systematic and random errors on a ship-by-ship basis going back to the

1452      start of the record and even infer likely measurement methods based on characteristic

1453      variations of the measurements with the meteorological conditions.

1454

1455      A more systematic approach to the assessment of analysis techniques is needed to

1456      elucidate the reasons for the differences between analyses and to assess the verisimilitude

1457      of analysis uncertainty estimates. Approaches could include theoretical inter-comparisons

1458      of statistical methods, comparisons based on well-defined sets of common input

1459      observations, and benchmarks built from datasets (such as model output) where the truth

1460      is known a piori. Benchmark tests like those planned by the International Surface

1461      Temperature Initiative [*Thorne et al.* 2011b] provide an objective measure against which

1462 analysis techniques can be evaluated. Both analysis techniques and benchmarks will have

1463 to be tailored appropriately for the particular problems affecting SST measurements and

1464 the latest understanding of measurement uncertainties.

1465

1466 A key weakness of historical SST data sets is the lack of attention paid to evaluating the

1467 effects of data biases particularly in the post-1941 records. Further independent estimates

1468 of the biases produced need to be undertaken using as diverse a range of means as

1469 possible and the robust critique of existing methods must continue. Ideally, these would

1470 be complemented by carefully-designed field tests of buckets and other measurement

1471 methods.

1472

1473 Combining new analysis techniques that have been appropriately benchmarked with

1474 novel approaches to assessing uncertainty arising from systematic errors, pervasive

1475 systematic errors and their adjustments will give new end-to-end analyses that will help

1476 to explore the uncertainties in historical SSTs in a more systematic manner.

1477

1478 For long-term historical analyses, there is no substitute for actual observations and

1479 relevant metadata. Efforts to identify archives of marine observations and digitize them

1480 are ongoing [*Brohan et al.*, 2009; *Wilkinson et al.*, 2011]. Such programs are labor

1481 intensive, first in identifying and cataloguing the holdings in archives around the world,

1482 then in creating and storing digital images of the paper books and finally in keying the

1483 observations. The difficulty of decoding hand written entries in a variety of languages,

1484 formats and scripts means that optical character recognition technologies are of limited

1485 use. A number of popular crowd-sourcing projects have been started to key information

1486 from ships logs that have historical as well meteorological interest. OldWeather.org has

1487 keyed data from Royal Navy logs from the First World War [*Brohan et al.*, 2009] and is

1488 now working on logs from polar expeditions. Digitization of data also holds the

1489 possibility of extending instrumental records further back in time [*Brohan et al.*, 2010].

1490 New observations, with reliable metadata, can be used not only to reduce uncertainty in

1491 SST analyses, but also to test the reliability of existing interpolated products and their

1492 uncertainties.

1493

1494 The ultimate destination of newly digitized observations is the International

1495 Comprehensive Ocean Atmosphere Data Set (ICOADS) [*Woodruff et al.*, 2011]. The

1496 ICOADS repository of marine meteorological data has long been the focus of advances in

1497 the understanding of marine climatology. It provides a consistent baseline for a wide

1498 range of studies, providing a solid basis for traceability and reproducibility. The

1499 continued existence, maintenance and improvement of ICOADS are essential to the

1500 future understanding of the global climate.

1501

1502 Finally, the work of identifying and quantifying uncertainties will be pointless, if those

1503 uncertainties are not used. Uncertainty estimates provided with data sets have sometimes

1504 been difficult to use or easy to use inappropriately. As pointed out by *Rayner et al.*

1505 [2009], "more reliable and user-friendly representations of uncertainty should be

1506 provided" in order to encourage their widespread and effective use.

1507

1520

1521   **Appendix A**

1522

1523   Figure 1 was calculated in the following way. Observations were separated into three

1524   groups – shallow, deep and unknown – using the metadata assignments of *Kennedy et al.*

1525   [2011c]. Bucket and buoy measurements were considered to be shallow. Engine intake

1526   and hull contact measurements were considered to be deep. Shallow measurements were

1527   assumed to exhibit a diurnal cycle equal to that measured by drifting buoys [*Kennedy et*

1528   *al.*, 2007]. Deep measurements were assumed to have no diurnal cycle. The two groups

1529   were assumed to measure the same temperature just before sunrise. The relative bias

1530   between the two was calculated by subtracting the minimum of the diurnal cycle from the

1531     daily average. This value varies by location and calendar month. The bias in each grid

1532     box was estimated by multiplying the relative bias by the fraction of shallow

1533     measurements. The bias was then normalized relative to the period 1961-1990, the

1534     anomaly period used for HadSST3. Figure 1 shows the global monthly average of the

1535     bias.

1536

1537     **References**

1538

1539     Abraham, J.P. et al. (2013), A review of global ocean temperature observations:

1540     Implications for ocean heat content estimates and climate change. Reviews of

1541     Geophysics. doi: 10.1002/rog.20022

1542

1543     Amot, A. (1954), Measurements of sea surface temperature for meteorological purposes.

1544     Results of observations from ocean weather station M, Meteorologiske Annaler, 4(1), 1-

1545     11.

1546

1547     Atkinson, C.P., N.A. Rayner, J. Roberts-Jones, R.O. Smith (2013), Assessing the quality

1548     of sea surface temperature observations from drifting buoys and ships on a platform-by-

1549     platform basis. Journal of Geophysical Research - Oceans, 118, 3507–3529,

1550     doi:10.1002/jgrc.20257.

1551

1552  Anderson, D. M., E. M. Mauk, E. R. Wahl, C. Morrill, A. J. Wagner, D. Easterling, and

1553  T. Rutishauser (2013), Global warming in an independent record of the past 130 years.

1554  Geophysical Research Letters, 40, 189–193, doi:10.1029/2012GL054271.

1555

1556  Beckers, J. M., M. Rixen (2003), EOF Calculations and Data Filling from Incomplete

1557  Oceanographic Datasets. J. Atmos. Oceanic Technol., 20, 1839–1856. doi: 10.1175/1520-

1558  0426(2003)020<1839:ECADFF>2.0.CO;2

1559

1560  Beggs H., R. Verein, G. Paltoglou, H. Kippo and M. Underwood (2012), Enhancing ship

1561  of opportunity sea surface temperature observations in the Australian region, Journal of

1562  Operational Oceanography, (ISSN: 1755-8778), 5, 59-73.

1563

1564  Bernie, D.J., E. Guilyardi, G. Madec, J.M. Slingo, S.J. Woolnough and J. Cole (2008),

1565  Impact of resolving the diurnal cycle in an ocean–atmosphere GCM. Part 2: A diurnally

1566  coupled CGCM. Climate Dynamics, 31(7-8), 909-925

1567

1568  Bernstein, R., and D. Chelton (1985), Large-Scale Sea Surface Temperature Variability

1569  from Satellite and Shipboard Measurements, Journal of Geophysical Research, 90(C6),

1570  11619-11630.

1571

1572  Berry, D.I. and E.C. Kent (2011), Air-Sea fluxes from ICOADS: the construction of a

1573  new gridded dataset with uncertainty estimates. International Journal of Climatology,

1574  31(7) 987-1001 doi:10.1002/joc.2059

1575

1576    BIPM, I., IFCC, I., & ISO, I. (2008). IUPAP and OIML, Evaluation of measurement

1577    data-Guide to the expression of uncertainty in measurement. International Organization

1578    for Standardization (ISO), Online: http://www. bipm. org/en/publications/guides/gum.

1579    html.

1580

1581    Bottomley, M., C.K. Folland, J. Hsiung, R.E. Newell, and D.E. Parker (1990), Global

1582    Ocean Surface Temperature Atlas. 20 + iv pp. and 313 color plates, Her Majesty's Stn.

1583    Off., Norwich, UK, 1990.

1584

1585    Brasnett, B. (2008), The impact of satellite retrievals in a global sea-surface-temperature

1586    analysis. Quarterly Journal of the Royal Meteorological Society, 134(636), 1745-1760

1587

1588    Breaker, L.C., W.W. Broenkow, M.W. Denny, and L.V. Beatman (2005), Reconstructing

1589    an 83-Year Time Series of Daily Sea Surface Temperature at Pacific Grove, California.

1590    Moss Landing, CA, Moss Landing Marine Laboratories,

1591    http://aquaticcommons.org/id/eprint/3129

1592

1593    Brohan P., J.J. Kennedy, I. Harris, S.F.B. Tett and P.D. Jones (2006), Uncertainty

1594    estimates in regional and global observed temperature changes: a new dataset from 1850.

1595    Journal of Geophysical Research, 111, D12106, doi:10.1029/2005JD006548.

1596

1597    Brohan, P., R. Allan, J.E. Freeman, A.M. Waple, D. Wheeler, C. Wilkinson, S. Woodruff

1598    (2009), Marine Observations of Old Weather. Bulletin of the American Meteorological

1599    Society, 90, 219-230. doi: 10.1175/2008BAMS2522.1

1600

1601    Brohan, P., C. Ward, G. Willetts, C. Wilkinson, R. Allan and D. Wheeler (2010), Arctic

1602    marine climate of the early nineteenth century. Climates of the Past, 6, 315-324,

1603    doi:10.5194/cp-6-315-2010

1604

1605    Brooks, C. (1926), Observing water-surface temperatures at sea. Monthly Weather

1606    Review, 54(6), 241-253, doi:10.1175/1520-0493(1926)54<241:OWTAS>2.0.CO;2.

1607

1608    Brooks, C. (1928), Reliability of different methods of taking sea-surface temperature

1609    measurements. Journal of the Washington Academy of Sciences, 18, 525-545.

1610

1611    Cannaby, H., and Hüsrevoğlu, Y. S. (2009), The influence of low-frequency variability

1612    and long-term trends in North Atlantic sea surface temperature on Irish waters. – ICES

1613    Journal of Marine Science, 66: 1480-1489.

1614

1615    Castro, S. L., G. A. Wick, and W. J. Emery (2012), Evaluation of the relative

1616    performance of sea surface temperature measurements from different types of drifting

1617    and moored buoys using satellite-derived reference products. Journal of Geophysical

1618    Research, 117, C02029, doi:10.1029/2011JC007472.

1619

1620    Chiodi, A.M. and D.E. Harrison (2006), Summertime subtropical sea surface temperature

1621    variability. Geophysical Research Letters, 33, L08601, doi:10.1029/2005GL024524

1622

1623    Collins, C., L. Giovando and K. Abbott-Smith (1975), Comparison of Canadian and

1624    Japanese merchant-ship observations of sea-surface temperature in the vicinity of present

1625    ocean station P 1927-33. Canadian Journal of Fisheries and Aquatic Science, 32(2), 253-

1626    258, doi:10.1139/f75-023.

1627

1628    Compo G.P., P.D. Sardeshmukh, J.S. Whitaker, P. Brohan, P.D. Jones and C. McColl

1629    (2013), Independent confirmation of global land warming without the use of station

1630    temperatures. Geophysical Research Letters doi: 10.1002/grl.50425

1631

1632    Cummings, J.A. (2005), Operational multivariate ocean data assimilation. Quarterly

1633    Journal of the Royal Meteorological Society, 131: 3583–3604. doi: 10.1256/qj.05.105

1634

1635    de Boyer Montégut, C., G. Madec, A.S. Fischer, A. Lazar, and D. Iudicone (2004),

1636    Mixed layer depth over the global ocean: an examination of profile data and a profile-

1637    based climatology. Journal of Geophysical Research, 109, C12003.

1638    doi:10.1029/2004JC002378

1639

1640    Deser C., A.S. Phillips and M.A. Alexander (2010), Twentieth century tropical sea

1641    surface temperature trends revisited. Geophysical Research Letters, 37, L10701, doi:

1642    10.1029/2010GL043321.

1643

1644    Dommenget, D. (2007), Evaluating EOF modes against a stochastic null hypothesis.

1645    Climate Dynamics, 28(5), 517-531.

1646

1647    Donlon, C., et al. (2007), The Global Ocean Data Assimilation Experiment High-

1648    resolution Sea Surface Temperature Pilot Project. Bulletin of the American

1649    Meteorological Society, 88, 1197–1213. doi: 10.1175/BAMS-88-8-1197

1650

1651    Emery, W. J., K. Cherkauer, B. Shannon and R.W. Reynolds (1997), Hull-Mounted Sea

1652    Surface Temperatures from Ships of Opportunity. Journal of Atmospheric and Oceanic

1653    Technology, 14, 1237–1251. doi: http://dx.doi.org/10.1175/1520-

1654    0426(1997)014<1237:HMSSTF>2.0.CO;2

1655

1656    Embury, O., C.J. Merchant and G.K. Corlett (2012), A Reprocessing for Climate of Sea

1657    Surface Temperature from the Along-Track Scanning Radiometers: Initial validation,

1658    accounting for skin and diurnal variability. Remote Sensing of Environment, 116(15) 62-

1659    78. DOI:10.1016/j.rse.2011.02.028

1660

1661    Emery, W., D. Baldwin, P. Schlüssel and R. Reynolds (2001), Accuracy of in situ sea

1662    surface temperatures used to calibrate infrared satellite measurements. Journal of

1663    Geophysical Research, 106(C2), 2387-2405, doi:10.1029/2000JC000246.

1664

1665    Folland, C., R. Reynolds, M. Gordon and D. Parker (1993), A study of six operational sea

1666    surface temperature analyses. Journal of Climate, 6(1), 96-113, doi:10.1175/1520-

1667    0442(1993)006<0096:ASOSOS>2.0.CO;2.

1668

1669    Folland, C.K. and D.E. Parker (1995), Correction of instrumental biases in historical sea

1670    surface temperature data. Quarterly Journal of the Royal Meteorological Society 121:

1671    319-367.

1672

1673    Folland C.K., N.A. Rayner, S. J. Brown, T.M. Smith, S.S.P. Shen, D.E. Parker, I.

1674    Macadam, P.D. Jones, R.N. Jones, N. Nicholls and D.M.H. Sexton (2001), Global

1675    temperature change and its uncertainties since 1861, Geophysical Research Letters,

1676    28(13), 2621–2624.

1677

1678    Folland, C. (2005), Assessing bias corrections in historical sea surface temperature using

1679    a climate model. International Journal of Climatology, 25(7), 895-911,

1680    doi:10.1002/joc.1171

1681

1682    Folland, C. K. and Salinger, M. J. (1995), Surface temperature trends and variations in

1683    New Zealand and the surrounding ocean, 1871–1993. International Journal of

1684    Climatology, 15: 1195–1218. doi: 10.1002/joc.3370151103

1685

1686    GCOS (2003), The second report on the adequacy of the Global Observing Systems for

1687    Climate in support of the UNFCCC. GCOS–82, WMO Tech. Doc. 1143, 85

1688    pp. [Available online at www.wmo.int/pages/prog/gcos/Publications/gcos-82_2AR.pdf.]

1689

1690    Gilhousen, D.B. (1987), A field evaluation of NDBC moored buoy winds. Journal of

1691    Atmospheric and Oceanic Technology, 4(1), 94-104

1692

1693    Gille, S.T. (2012), Diurnal variability of upper ocean temperatures from microwave

1694    satellite measurements and Argo profiles. Journal of Geophysical Research, 117, C11027,

1695    doi:10.1029/2012JC007883.

1696

1697    Good, S.A., G.K. Corlett, J.J. Remedios, E.J. Noyes, and D.T. Llewellyn-Jones (2007),

1698    The global trend in sea surface temperature from 20 years of Advanced Very High

1699    Resolution Radiometer data. Journal of Climate, 20(7), 1255-1264

1700    doi:10.1175/JCLI4049.1

1701

1702    Gouretski, V., J.J. Kennedy, T. Boyer, and A. Köhl (2012), Consistent near-surface ocean

1703    warming since 1900 in two largely independent observing networks. Geophysical

1704    Research Letters, 39, L19606, doi:10.1029/2012GL052975.

1705

1706    Grodsky, S.A., J.A. Carton and H. Liu (2008), Comparison of bulk sea surface and mixed

1707    layer temperatures. Journal of Geophysical Research 113, C10026,

1708    doi:10.1029/2008JC004871.

1709

1710    Hanawa, K., S. Yasunaka, T. Manabe, and N. Iwasaka (2000), Examination of correction

1711    to historical SST data using long-term coastal SST data taken around Japan. Journal of

1712    the Meteorological Society of Japan, 78, 187-195.

1713

1714    Hannachi A., I.T. Jolliffe, and D.B. Stephenson (2007), Empirical orthogonal functions

1715    and related techniques in atmospheric science: A review. International Journal of

1716    Climatology, 27:1119-1152.

1717

1718    Houseago-Stokes, R.E., and P.G. Challenor (2004), Using PPCA to Estimate EOFs in the

1719    Presence of Missing Values. J. Atmos. Oceanic Technol., 21, 1471–1480.

1720    doi: 10.1175/1520-0426(2004)021<1471:UPTEEI>2.0.CO;2

1721

1722    Hurrell, J.W. and K.E. Trenberth (1999), Global Sea Surface Temperature Analyses:

1723    Multiple Problemsand Their Implications for Climate Analysis, Modeling, and

1724    Reanalysis. Bulletin of the American Meteorological Society, 80, 2661–2678. doi:

1725    http://dx.doi.org/10.1175/1520-0477(1999)080<2661:GSSTAM>2.0.CO;2

1726

1727    Ilin A. and A. Kaplan (2009), Bayesian PCA for Reconstruction of Historical Sea Surface

1728    Temperatures. Proceedings of the International Joint Conference on Neural Networks

1729    (IJCNN 2009), pp. 1322-1327, Atlanta, USA, 2009.

1730

1731    Ingleby, B. (2010), Factors Affecting Ship and Buoy Data Quality: A Data Assimilation

1732    Perspective. Journal of Atmospheric and Oceanic Technology 27:9, 1476-1489

1733

1734   Ishii, M., M. Kimoto, and M. Kachi (2003), Historical Ocean Subsurface Temperature

1735   Analysis with Error Estimates. Mon. Wea. Rev., 131, 51–73. doi: 10.1175/1520-

1736   0493(2003)131<0051:HOSTAW>2.0.CO;2

1737

1738   Ishii, M., A. Shouji, S. Sugimoto, and T. Matsumoto (2005), Objective analyses of sea-

1739   surface temperature and marine meteorological variables for the 20th century using

1740   ICOADS and the Kobe collection. International Journal of Climatology, 25(7), 865-879,

1741   doi:10.1002/joc.1169.

1742

1743   James, R., and P. Fox (1972), Comparative sea surface temperature measurements in

1744   WMO reports on marine science affairs, rep 5, Tech. Rep. 336, WMO.

1745

1746   Jones, G.S., and P.A. Stott (2011), Sensitivity of the attribution of near surface

1747   temperature warming to the choice of observational dataset. Geophysical Research

1748   Letters, 38, L21702, doi:10.1029/2011GL049324.

1749

1750   Jones, P.D. (1994), Hemispheric Surface Air Temperature Variations: A Reanalysis and

1751   an Update to 1993. Journal of Climate, 7, 1794–1802. doi: 10.1175/1520-

1752   0442(1994)007<1794:HSATVA>2.0.CO;2

1753

1754    Jones, P.D., T.J. Osborn and K.R. Briffa (1997), Estimating Sampling Errors in Large-

1755    Scale Temperature Averages. Journal of Climate, 10, 2548–2568. doi: 10.1175/1520-

1756    0442(1997)010<2548:ESEILS>2.0.CO;2

1757

1758    Jones, P.D. and T.M.L. Wigley (2010), Estimation of global temperature trends: What's

1759    important and what isn't. Climatic Change, 100 (1). pp. 59-69.

1760

1761    Kaplan, A., M. Cane, Y. Kushnir, A. Clement, M. Blumenthal, and B. Rajagopalan

1762    (1998), Analyses of global sea surface temperature 1856-1991. Journal of Geophysical

1763    Research, 103(C9), 18,567-18,589, doi:10.1029/97JC01736

1764

1765    Karnauskas, K.B. (2013), Can we distinguish canonical El Niño from Modoki?

1766    Geophysical Research Letters, doi: 10.1002/grl.51007

1767

1768    Karspeck, A.R., A. Kaplan, and S.R. Sain, (2012), Bayesian modelling and ensemble

1769    reconstruction of mid-scale spatial variability in North Atlantic sea-surface temperatures

1770    for 1850-2008. Quarterly Journal of the Royal Meteorological Society, 138, 234-248. doi:

1771    10.1002/qj.900

1772

1773    Kawai, Y. and H. Kawamura, (2000), Study on a Platform Effect in the In Situ Sea

1774    Surface Temperature Observations under Weak Wind and Clear Sky Conditions Using

1775    Numerical Models. Journal of Atmospheric and Oceanic Technology, 17, 185–196. doi:

1776    10.1175/1520-0426(2000)017<0185:SOAPEI>2.0.CO;2

1777

1778    Kawai, Y. and A. Wada (2007), Diurnal sea surface temperature variation and its impact

1779    on the atmosphere and ocean: a review. Journal of Oceanography, 2007, 63: 721-744.

1780

1781    Kennedy, J.J., P. Brohan and S.F.B. Tett (2007), A global climatology of the diurnal

1782    variations in sea-surface temperature and implications for MSU temperature trends.

1783    Geophysical Research Letters, 34(5), L05712 doi:10.1029/2006GL028920

1784

1785    Kennedy, J.J., R. Smith, and N. Rayner (2011a), Using AATSR data to assess the quality

1786    of in situ sea surface temperature observations for climate studies. Remote Sensing of

1787    Environment. 116, 79–92 http://dx.doi.org/10.1016/j.rse.2010.11.021

1788

1789    Kennedy, J.J., N.A. Rayner, R.O. Smith, M. Saunby and D.E. Parker (2011b),

1790    Reassessing biases and other uncertainties in sea-surface temperature observations since

1791    1850 part 1: measurement and sampling errors. Journal of Geophysical Research, 116,

1792    D14103, doi:10.1029/2010JD015218

1793

1794    Kennedy, J.J., N.A. Rayner, R.O. Smith, M. Saunby and D.E. Parker (2011c),

1795    Reassessing biases and other uncertainties in sea-surface temperature observations since

1796    1850 part 2: biases and homogenisation. Journal of Geophysical Research, 116, D14104,

1797    doi:10.1029/2010JD015220

1798

1799

1800    Kent, E., P. Taylor, B. Truscott, and J. Hopkins (1993), The accuracy of Voluntary

1801    Observing Ships' meteorological observations - results of the VSOP-NA. Journal of

1802    Atmospheric and Oceanic Technology, 10(4), 591-608, doi:10.1175/1520-

1803    0426(1993)010<0591:TAOVOS>2.0.CO;2.

1804

1805    Kent, E.C., P.G. Challenor and P.K. Taylor (1999), A statistical determination of the

1806    random observational errors present in voluntary observing ships meteorological reports.

1807    Journal of Atmospheric and Oceanic Technology, 16(7), 905-914

1808

1809    Kent, E.C. and D.I. Berry (2005), Quantifying random measurement errors in Voluntary

1810    Observing Ships' meteorological observations. International Journal of Climatology,

1811    25(7), 843-856, doi:10.1002/joc.1167

1812

1813    Kent, E., and P. Challenor (2006), Toward estimating climatic trends in SST. Part II:

1814    Random errors. Journal of Atmospheric and Oceanic Technology, 23(3), 476-486,

1815    doi:10.1175/JTECH1844.1.

1816

1817    Kent, E., and A. Kaplan (2006), Toward estimating climatic trends in SST. Part III:

1818    Systematic biases. Journal of Atmospheric and Oceanic Technology, 23(3), 487-500,

1819    doi:10.1175/JTECH1845.1.

1820

1821    Kent, E.C., S.D. Woodruff and D.I. Berry (2007), Metadata from WMO Publication No.

1822    47 and an Assessment of Voluntary Observing Ship Observation Heights in ICOADS.

1823    Journal of Atmospheric and Oceanic Technology, 24, (2), 214-234.

1824    doi:10.1175/JTECH1949.1

1825

1826    Kent, E., and D. Berry (2008), Assessment of the marine observing system (ASMOS):

1827    Final report, Tech. Rep. 32, Natl. Oceanogr. Cent., Southampton, U. K.

1828

1829    Kent, E.C., J.J. Kennedy, D.I. Berry and R.O. Smith (2010), Effects of instrumentation

1830    changes on sea surface temperature measured in situ. Wiley Interdisciplinary Reviews:

1831    Climate Change.1(5) 718-728 doi:10.1002/wcc.55

1832

1833    Kent, E.C., N.A. Rayner, D.I. Berry, M. Saunby, B.I. Moat, J.J. Kennedy, and D.E.

1834    Parker (2013), Global analysis of night marine air temperature and its uncertainty since

1835    1880: The HadNMAT2 data set, J. Geophys. Res. Atmos., 118, 1281–1298,

1836    doi:10.1002/jgrd.50152.

1837

1838    Kirk, T., and A. Gordon (1952), Comparison of intake and bucket methods for measuring

1839    sea temperature. Marine Observer, 22, 33-39.

1840

1841    Knudsen, J. (1966), An experiment in measuring the sea surface temperature for synoptic

1842    purposes. Tech. Rep. 12, Det. Norske Meteor. Inst.

1843

1844    Kondrashov, D. and M. Ghil (2006), Spatio-temporal filling of missing points in

1845    geophysical data sets, Nonlin. Processes Geophys., 13, 151-159, doi:10.5194/npg-13-

1846    151-2006.

1847

1848    Lindau, R. (2003), Errors of Atlantic air-sea fluxes derived from ship observations.

1849    Journal of Climate, 16, 783–788.

1850

1851    Lumby, J. (1927), The surface sampler, an apparatus for the collection of samples from

1852    the sea surface from ships in motion with a note on surface temperature observations. J.

1853    Cons. Perm. Int. Explor. Mer., 2, 332-342.

1854

1855    Luttinen, J. and A. Ilin. (2009), Variational Gaussian-Process Factor Analysis for

1856    Modeling Spatio-Temporal Data. In Proceedings of the 23rd Annual Conference on

1857    Neural Information Processing Systems (NIPS 2009), Vancouver, Canada, 2009.

1858

1859    Luttinen J. and A. Ilin. (2012), Efficient Gaussian Process Inference for Short-Scale

1860    Spatio-Temporal Modeling. Accepted to the 15th International Conference on Artificial

1861    Intelligence and Statistics (AISTATS 2012).

1862

1863    Lyman, J.M., S.A. Good, V.V. Gouretski, M. Ishii, G.C. Johnson, M.D. Palmer, D.M.

1864    Smith and J.K. Willis (2010), Robust warming of the global upper ocean. Nature 465,

1865    334-337 doi:10.1038/nature09043

1866

1867     MacKenzie B.R., and D. Schiedek (2007), Long-term sea surface temperature

1868     baselines—time series, spatial covariation and implications for biological processes,

1869     Journal of Marine Systems, 68(3–4), pp 405-420 doi: 10.1016/j.jmarsys.2007.01.003.

1870

1871     Matthews, J.B.R. (2013), Comparing historical and modern methods of sea surface

1872     temperature measurement – Part 1: Review of methods, field comparisons and dataset

1873     adjustments, Ocean Sci., 9, 683-694, doi:10.5194/os-9-683-2013.

1874

1875     Matthews, J.B.R., and Matthews, J.B. (2013), Comparing historical and modern methods

1876     of sea surface temperature measurement – Part 2: Field comparison in the central tropical

1877     Pacific, Ocean Sci., 9, 695-711, doi:10.5194/os-9-695-2013.

1878

1879     Maul, G.A., A.M. Davis, and J.W. Simmons (2001), Seawater temperature trends at USA

1880     Tide Gauge sites. Geophysical Research Letters, 28: 3935–3937. doi:

1881     10.1029/2001GL013458

1882

1883     Mears, C. A., F.J. Wentz, P. Thorne and D. Bernie (2011), Assessing uncertainty in

1884     estimates of atmospheric temperature changes from MSU and AMSU using a Monte-

1885     Carlo estimation technique. Journal of Geophysical Research, 116, D08112,

1886     doi:10.1029/2010JD014954.

1887

1888     Merchant, C., D. Llewellyn-Jones, R. Saunders, N. Rayner, E. Kent, C. Old, D. Berry, A.

1889     Birks, T. Blackmore and G. Corlett (2008a), Deriving a sea surface temperature record

1890    suitable for climate change research from the along-track scanning radiometers.

1891    Advances in Space Research 41 (1), 1–11.

1892

1893    Merchant C.J., P. Le Borgne, A. Marsouin and H. Roquet (2008b), Optimal estimation of

1894    sea surface temperature from split-window observations. Remote Sensing of

1895    Environment, 112, 2469–2484

1896

1897    Merchant, C.J., O. Embury, N.A. Rayner, D.I. Berry, G. Corlett, K. Lean, K.L. Veal, E.C.

1898    Kent, D. Llewellyn-Jones, J.J. Remedios, and R. Saunders (2012), A twenty-year

1899    independent record of sea surface temperature for climate from Along Track Scanning

1900    Radiometers . Journal of Geophysical Research, 117, C12013,

1901    doi:10.1029/2012JC008400.

1902

1903    Merchant, C.J., S. Matthiesen, N.A. Rayner, J.J. Remedios, P.D. Jones, F. Olesen, B.

1904    Trewin, P.W. Thorne, R. Auchmann, G.K. Corlett, P.C. Guillevic, and G.C. Hulley

1905    (2013), The surface temperatures of the earth: steps towards integrated understanding of

1906    variability and change, Geosci. Instrum. Method. Data Syst. Discuss., 3, 305-345,

1907    doi:10.5194/gid-3-305-2013.

1908

1909    Morice, C.P., J.J. Kennedy, N.A. Rayner, and P.D. Jones (2012), Quantifying

1910    uncertainties in global and regional temperature change using an ensemble of

1911    observational estimates: The HadCRUT4 dataset. Journal of Geophysical Research, 117,

1912    D08101, doi:10.1029/2011JD017187.

1913

1914   Morrissey, M., and J. Greene (2009), A theoretical framework for the sampling error

1915   variance for three-dimensional climate averages of ICOADS monthly ship data.

1916   Theoretical and Applied Climatology, 96(3-4), 235-248, doi:10.1007/s00704-008-0027-3

1917

1918   Moyer, K.A., and R.A. Weller (1997), Observations of Surface Forcing from the

1919   Subduction Experiment: A Comparison with Global Model Products and Climatological

1920   Datasets. Journal of Climate, 10, 2725–2742. doi: http://dx.doi.org/10.1175/1520-

1921   0442(1997)010<2725:OOSFFT>2.0.CO;2

1922

1923   Nixon, S.W., S. Granger, B.A. Buckley, M. Lamont, B. Rowell (2004), A one hundred

1924   and seventeen year coastal water temperature record from Woods Hole, Massachusetts.

1925   Estuaries 27(3), 397-404, doi: 10.1007/BF02803532

1926

1927   O'Carroll, A.G., J.R. Eyre and R.W. Saunders (2008), Three-way error analysis between

1928   AATSR, AMSR-E and in situ sea surface temperature observations. Journal of

1929   Atmospheric and Oceanic Technology, 25(7), 1197-1207

1930

1931   Palmer, M. D., K. Haines, S.F.B. Tett and T.J. Ansell (2007), Isolating the signal of

1932   ocean global warming, Geophysical Research Letters, 34, L23610,

1933   doi:10.1029/2007GL031712.

1934

1935    Palmer, M.D. and P. Brohan (2011), Estimating sampling uncertainty in fixed-depth and

1936    fixed-isotherm estimates of ocean warming. International Journal of Climatology, 31(7),

1937    980–986, doi: 10.1002/joc.2224

1938

1939    Parker, D.E. (1987), The sensitivity of estimates of trends of global and hemispheric

1940    marine temperatures to limitations in geographical coverage.  Long Range Forecasting

1941    and Climate Research Memo LRFC 12.

1942

1943    Perlroth, I. (1962), Relationship of central pressure of hurricane Esther (1961) and the sea

1944    surface temperature field. Tellus, 14: 403–408. doi: 10.1111/j.2153-3490.1962.tb01353.x

1945

1946    Prytherch, J., J.T. Farrar, and R.A. Weller (2013), Moored surface buoy observations of

1947    the diurnal warm layer. Journal of Geophysical Research (in press) doi:

1948    10.1002/jgrc.20360

1949

1950    Rayner, N.A., D.E. Parker, E.B. Horton, C.K. Folland, L.V. Alexander, D.P. Rowell, E.C.

1951    Kent and A. Kaplan (2003), Global analyses of sea surface temperature, sea ice, and night

1952    marine air temperature since the late nineteenth century. Journal of Geophysical

1953    Research, Vol. 108, No. D14, 4407 10.1029/2002JD002670

1954

1955    Rayner, N., P. Brohan, D. Parker, C. Folland, J. Kennedy, M. Vanicek, T. Ansell, and S.

1956    Tett (2006), Improved analyses of changes and uncertainties in sea surface temperature

1957    measured in situ since the mid-nineteenth century: The HadSST2 data set. Journal of

1958    Climate, 19(3), 446-469, doi:10.1175/JCLI3637.1.

1959

1960    Rayner, N., et al. (2009), Evaluating climate variability and change from modern and

1961    historical SST observations, in Proceedings of OceanObs'09: Sustained Ocean

1962    Observations and Information for Society, vol. 2, Venice, Italy, 21-25 September 2009,

1963    ESA Publ. WPP-306, edited by J. Hall, D. Harrison, and D. Stammer, ESA, Paris,

1964    doi:10.5270/OceanObs09.cwp.71

1965

1966    Reverdin, G., J. Boutin, N. Martin, A. Lourenco, P. Bouruet-Aubertot, A. Lavin, J.

1967    Mader, P. Blouch, J. Rolland, F. Gaillard and P. Lazure (2010), Temperature

1968    Measurements from Surface Drifters. Journal of Atmospheric and Oceanic Technology,

1969    27, 1403-1409. doi: 10.1175/2010JTECHO741.1

1970

1971    Reynolds, R.W., N.A. Rayner, T.M. Smith, D.C. Stokes and W.Q. Wang (2002), An

1972    improved in situ and satellite SST analysis for climate. Journal of Climate, 15(13), 1609-

1973    1625

1974

1975    Reynolds, R.W., C.L. Gentemann and G.K. Corlett, (2010), Evaluation of AATSR and

1976    TMI Satellite SST Data. Journal of Climate, 23, 152–165. doi: 10.1175/2009JCLI3252.1

1977

1978    Roberts-Jones, J., E.K. Fiedler and M.J. Martin (2012), Daily, Global, High-Resolution

1979    SST and Sea Ice Reanalysis for 1985–2007 Using the OSTIA System. Journal of Climate,

1980    25, 6215–6232. doi: 10.1175/JCLI-D-11-00648.1

1981

1982    Roll, H. (1951), Water temperature measurements on deck and in the engine room. Ann.

1983    Meteor., 4, 439-443.

1984

1985    Roweis, S. (1998), EM Algorithms for PCA and SPCA. Neural Information Processing

1986    Systems 10 (NIPS'97) pp.626-632

1987

1988    Rutherford, S., M.E. Mann, T.L. Delworth, and R.J. Stouffer (2003), Climate Field

1989    Reconstruction under Stationary and Nonstationary Forcing. *J. Climate*, **16**, 462–479.

1990    doi: 10.1175/1520-0442(2003)016<0462:CFRUSA>2.0.CO;2

1991

1992    Sarachik, E.S., (1984), Large-scale surface heat fluxes. Large-Scale Oceanographic

1993    Experiments and Satellites, C Gautier and M. Fieux, Eds., Reidel, 147–165.

1994

1995    Saur, J. (1963), A study of the quality of sea water temperatures reported in the logs of

1996    ships' weather observations. Journal of Applied Meteorology, 2(3), 417-425,

1997    doi:10.1175/1520-0450(1963)002<0417:ASOTQO>2.0.CO;2.

1998

1999    Schell, I.I. (1959), On a criterion of representativeness of sea-surface data. Bull. Amer.

2000    Met. Soc. 40(11) pp 571-574.

2001

2002    Schneider, T. (2001), Analysis of Incomplete Climate Data: Estimation of Mean Values

2003    and Covariance Matrices and Imputation of Missing Values. J. Climate, 14, 853–871.

2004    doi:10.1175/1520-0442(2001)014<0853:AOICDE>2.0.CO;2

2005

2006    She, J., J.L. Høyer and  J. Larsen (2007), Assessment of sea surface temperature

2007    observational networks in the Baltic Sea and North Sea. Journal of Marine Systems,

2008    65(1–4), 314–335. doi:10.1016/j.jmarsys.2005.01.004,

2009

2010    Shen, S.S.P., H. Yin and T.M. Smith, (2007), An Estimate of the Sampling Error

2011    Variance of the Gridded GHCN Monthly Surface Air Temperature Data. Journal of

2012    Climate, 20, 2321-2331. doi: 10.1175/JCLI4121.1

2013

2014    Sheppard, C.R.C. and N.A. Rayner (2002), Utility of the Hadley Centre sea-Ice and Sea

2015    Surface Temperature data set (HadISST1) in two widely contrasting coral reef areas.

2016    Marine Pollution Bulletin, 44 303 - 308 (0025-326X)

2017

2018    Simmons, A.J., K.W. Willett, P.D. Jones, P.W. Thorne and D. Dee (2010), Low-

2019    frequency variations in surface atmospheric humidity, temperature and precipitation:

2020    Inferences from reanalyses and monthly gridded observation datasets. Journal of

2021    Geophysical Research - Atmospheres.115, D01110, doi:10.1029/2009JD012442.

2022

2023 Smith, D.M., S. Cusack, A.W. Colman, C.K. Folland, G.R. Harris, and J.M. Murphy

2024 (2007), Improved Surface Temperature Prediction for the Coming Decade from a Global

2025 Climate Model. Science 317 (5839), 796-799. doi:10.1126/science.1139540

2026

2027 Smith, D., C. Mutlow, J. Delderfield, B. Watkins, G. Mason (2012), ATSR infrared

2028 radiometric calibration and in-orbit performance. Remote Sensing of Environment 116,

2029 4–16, doi:10.1016/j.rse.2011.01.027

2030

2031 Smith, T., and R. Reynolds (2002), Bias corrections for historical sea surface

2032 temperatures based on marine air temperatures. Journal of Climate, 15(1), 73.

2033

2034 Smith, T.M. and R.W. Reynolds, (2003), Extended Reconstruction of Global Sea Surface

2035 Temperatures Based on COADS Data (1854-1997). Journal of Climate, 16, 1495-1510.

2036

2037 Smith, T.M., and R.W. Reynolds, (2004), Improved Extended Reconstruction of SST

2038 (1854-1997). Journal of Climate, 17, 2466-2477.

2039

2040 Smith, T.M., R.W. Reynolds, (2005), A Global Merged Land-Air-Sea Surface

2041 Temperature Reconstruction Based on Historical Observations (1880-1997). Journal of

2042 Climate, 18, 2021-2036. doi: 10.1175/JCLI3362.1

2043

2044 Smith, T.M., R.W. Reynolds, T.C. Peterson and J.Lawrimore, (2008), Improvements to

2045 NOAA's Historical Merged Land-Ocean Surface Temperature Analysis (1880-2006).

2046 Journal of Climate, 21, 2283-2296.

2047

2048 Stevenson, R.E. (1964), The Influence of a Ship on the Surrounding Air and Water

2049 Temperatures. Journal of Applied Meteorology, 3, 115–118. doi: 10.1175/1520-

2050 0450(1964)003<0115:TIOASO>2.0.CO;2

2051

2052 Strong, A.E. and E.P. McClain (1984), Improved Ocean Surface Temperatures From

2053 Space-Comparisons With Drifting Buoys. Bulletin of the American Meteorological

2054 Society, 65, 138–142. doi: 10.1175/1520-0477(1984)065<0138:IOSTFS>2.0.CO;2

2055

2056 Stubbs, M.W. (1965), The standard error of a sea surface temperature as measured using

2057 a canvas bucket. The Meteorological magazine 94(1112).

2058

2059 Tabata, S. (1978a), On the accuracy of sea-surface temperatures and salinities observed in

2060 the Northeast Pacific Ocean. Atmosphere Ocean, 16(3), 237-247.

2061

2062 Tabata, S. (1978b), Comparison of observations of sea-surface temperatures at ocean

2063 staion P and NOAA buoy stations and those made by merchant ships travelling in their

2064 vicinities in the northeast Pacific Ocean. Journal of Applied Meteorology, 17(3), 374-

2065 385, doi:10.1175/1520-0450(1978)017<0374:COOOSS>2.0.CO;2.

2066

2067    Tauber, G. (1969), The comparative measurements of sea surface temperature in the

2068    USSR, Tech. Rep. 103, WMO.

2069

2070    Thompson, D.W.J., J.J. Kennedy, J.M. Wallace and P.D. Jones (2008), A large

2071    discontinuity in the mid-twentieth century in observed global-mean surface temperature.

2072    Nature 453, 646-649

2073

2074    Thompson, D.W.J., J.M. Wallace, J.J. Kennedy and P.D. Jones (2010), An abrupt drop in

2075    Northern Hemisphere sea surface temperature around 1970. Nature 467, 444-447,

2076    doi:10.1038/nature09394

2077

2078    Thorne, P.W., D.E. Parker, J.R. Christy and C.A. Mears (2005), Uncertainties in climate

2079    trends: Lessons from upper-air temperature records. Bulletin of the American

2080    Meteorological Society, 86, 1437–1442.

2081

2082    Thorne, P.W., J.R. Lanzante, T.C. Peterson, D.J. Seidel, K.P. Shine (2011), Tropospheric

2083    temperature trends: history of an ongoing controversy. WIREs Climate Change 2(1) 66-

2084    88 doi:10.1002/wcc.80

2085

2086    Thorne, P.W., and Coauthors, (2011b), Guiding the Creation of A Comprehensive

2087    Surface Temperature Resource for Twenty-First-Century Climate Science. Bull. Amer.

2088    Meteor. Soc., 92, ES40–ES47. doi: http://dx.doi.org/10.1175/2011BAMS3124.1

2089

2090    Tokinaga, H., S. Xie, C. Deser, Y. Kosaka and Y.M. Okumura (2012), Slowdown of the

2091    Walker circulation driven by tropical Indo-Pacific warming. Nature 491, 439–443.

2092    doi:10.1038/nature11576

2093

2094    Tung, K. and J. Zhou, (2010), The Pacific's Response to Surface Heating in 130 Yr of

2095    SST: La Niña–like or El Niño–like?. Journal of Atmospheric Science, 67, 2649–2657.

2096    doi: 10.1175/2010JAS3510.1

2097

2098    Vecchi G.A., A. Clement, and B.J. Soden (2008), Examining the Tropical Pacific's

2099    Response to Global Warming. Eos, Vol. 89, No. 9, 26 February 2008

2100

2101    Venema, V.K.C., O. Mestre, E. Aguilar, I. Auer, J.A. Guijarro, P. Domonkos, G.

2102    Vertacnik, T. Szentimrey, P. Stepanek, P. Zahradnicek, J. Viarre, G. Müller-Westermeier,

2103    M. Lakatos, C.N. Williams, M.J. Menne, R. Lindau, D. Rasol, E. Rustemeier, K.

2104    Kolokythas, T. Marinova, L. Andresen, F. Acquaotta, S. Fratianni, S. Cheval, M.

2105    Klancar, M. Brunetti, C. Gruber, M. Prohom Duran, T. Likso, P. Esteban and T.

2106    Brandsma (2012), Benchmarking homogenization algorithms for monthly data. Climates

2107    of the Past, 8, 89-115, doi:10.5194/cp-8-89-2012

2108

2109    Wahl, E. (1948), Water temperature measurements on deck and in the engine room. Ann.

2110    Meteor., 1(7).

2111

2112    Walden, H. (1966), On water temperature measurements aboard merchant vessels (in

2113    German). Ocean Dynamics, 19, 21-28, doi:10.1007/BF02321345.

2114

2115    Weare, B.C. (1989), Uncertainties in estimates of surface heat fluxes derived from marine

2116    reports over the tropical and subtropical oceans. Tellus A, 41A: 357–370. doi:

2117    10.1111/j.1600-0870.1989.tb00388.x

2118

2119    Weare, B.C. and Strub, P.T. (1981), The significance of sampling biases on calculated

2120    monthly mean oceanic surface heat fluxes. Tellus, 33: 211–224. doi: 10.1111/j.2153-

2121    3490.1981.tb01745.x

2122

2123    Wilkerson, J.C., and M. D. Earle (1990), A study of differences between environmental

2124    reports by ships in the voluntary observing program and measurements from NOAA

2125    buoys. Journal of Geophysical Research, 95(C3), 3373–3385,

2126    doi:10.1029/JC095iC03p03373.

2127

2128    Wilkinson, C., S.D. Woodruff, P. Brohan, S. Claesson, E. Freeman, F. Koek, S.J. Lubker,

2129    C. Marzin and D. Wheeler (2011), Recovery of logbooks and international marine data:

2130    the RECLAIM project. International Journal of Climatology, 31(7) 968-979

2131    doi:10.1002/joc.2102

2132

2133    Woodruff, S.D., S.J. Worley, S.J. Lubker, Z. Ji, E. Freeman, D.I. Berry, P. Brohan, E.C.

2134    Kent, R.W. Reynolds, S.R. Smith and C. Wilkinson (2011), ICOADS Release 2.5:

2135  extensions and enhancements to the surface marine meteorological archive. International

2136  Journal of Climatology. 31(7) 951-967 doi:10.1002/joc.2103

2137

2138  Worley, S.J., S.D. Woodruff, R.W. Reynolds, S.J. Lubker, and N. Lott, 2005: ICOADS

2139  Release 2.1 data and products. International Journal of Climatology. (CLIMAR-II Special

2140  Issue), 25, 823-842 doi:10.1002/joc.1166

2141

2142  Xu, F., and A. Ignatov (2010), Evaluation of in situ sea surface temperatures for use in

2143  the calibration and validation of satellite retrievals. Journal of Geophysical Research,

2144  115, C09022, doi:10.1029/2010JC006129.

2145

2146  Yasunaka, S. and K. Hanawa (2002), Regime shifts found in the Northern Hemisphere

2147  SST field. Journal of the Meteorological Society of Japan 80: 119-135.

2148

2149  Yasunaka, S. and K. Hanawa (2011), Intercomparison of historical sea surface

2150  temperature datasets. International Journal of Climatology, 31(7) 1056-1073

2151  doi:10.1002/joc.2104

2152

2153  Yu, L., R.A. Weller, and B. Sun, 2004: Improving latent and sensible heat flux estimates

2154  for the Atlantic Ocean (1988-1999) by a synthesis approach. J. Climate, 17, 373–393

2155

2156

2157

| References | Estimated measurement uncertainty for ship measurements |
|---|---|
| *Stubbs* [1965] | 0.11±0.01K for canvas bucket measurements from an Ocean Weather Ship |
| *Strong and McLean* [1984] | 1.8K RMS difference between ship and AVHRR data |
| *Bernstein and Chelton* [1985] pg 11620 | 1.1 K |
| *Sarachik* [1984], *Weare* [1989] pg 359 | 1 K |
| *Wilkerson and Earle* [1990] pg 3381 | 3.5 K |
| *Cummings* [2005] Table 1, pg 3592 | 1.3 K (ERI) 0.6 K (Hull sensor) 1.2 K (bucket) |
| *Kent and Challenor* [2006] pg 484 | 1.2±0.4 K or 1.3±0.3 K depending on how measurements were weighted |
| *Kent et al.* [1999] abstract | 1.5±0.1 K |
| *Kent and Berry* [2005] Table 2 pg 853 | 1.3±0.1 K and 1.2±0.1 K |
| *Reynolds et al.* [2002] pg 1613 | 1.3 K |
| *Kennedy et al.* [2011a] pg 83 | 1.0 K |
| *Ingleby* [2010] Table 10 pg 1487 | 0.9 K for automatic systems 1.2 K for manual measurements |
| *Kent and Berry* [2008] Table 5a pg 11 | 1.1 K |
| *Xu and Ignatov* [2010] pg 16 of 18 | 1.16 K |

2158    **Table 1**: List of estimates of measurement error uncertainties for ships where random and

2159    systematic errors were not dealt with separately.

2160

| References | Estimated measurement uncertainty for drifting buoy measurements |
|---|---|
| *Strong and Mclean* [1984] | 0.6K RMS difference between drifter and AVHRR |
| *Reynolds et al.* [2002] pg 1613 | 0.5 K |
| *Emery et al.* [2001] pg 2393 | 0.3 K |
| *Cummings* [2005] Table 1, pg 3592 | 0.12 K |
| *O'Carroll et al.* [2008] abstract | 0.23 K |
| *Kent and Berry* [2008] Table 5c pg 12 | 0.67 K |
| *Ingleby* [2010] Table 10 pg 1487 | 0.33 K |
| *Kennedy et al.* [2011a] pg 83 | 0.2-0.4 K |
| *Xu and Ignatov* [2010] pg 16 of 18 | 0.26 K |
| *Merchant et al.* [2012] Table 2 pg 8 of 18 | 0.15-0.19 K |

2161    **Table 2**: List of estimates of measurement error uncertainties for drifting buoys where

2162    random and systematic errors were not dealt with separately.

2163

| Reference | Estimated measurement uncertainty for moored buoy measurements |
|---|---|
| *Cummings* [2005] Table 1, pg 3592 | 0.05 K |
| *Kent and Berry* [2008] Table 5b pg 11 | 0.4 K |

| | |
|---|---|
| *Kennedy et al.* [2011a] pg 83 | tropical moorings, 0.12 K; all moorings, 0.21 K |
| *Xu and Ignatov* [2010] pg 16 of 18 | tropical moorings, 0.30 K; coastal moorings, 0.39 K |
| *Gilhousen* [1987] Table 6 pg 104 | 0.22 K |

2164  **Table 3**: List of estimates of measurement error uncertainties for moored buoys where

2165  random and systematic errors were not dealt with separately.

2166

| Reference | Platform type | Random | Systematic | Notes |
|---|---|---|---|---|
| *Kent and Berry* [2008] pg 11 Table 5a | Ship | 0.7 K | 0.8 K | From comparison with Numerical Weather Prediction fields provided with VOSClim data |
| Pg 12 Table 5c | Drifter | 0.6 K | 0.3 K | |
| Pg 11 Table 5b | Mooring | 0.3 K | 0.2 K | |
| *Kennedy et al.* [2011a, 2011b] pg 86 | Ship | 0.74 K | 0.71 K | From comparison with Along Track Scanning Radiometer SST retrievals |
| Pg 86 | Drifter | 0.26 K | 0.29 K | |
| *Brasnett* [2008] values estimated for present study by | Ship | 1.16 K | 0.69 K | From comparison with interpolated fields |

| | | | | |
|---|---|---|---|---|
| author | | | | |
| *Xu and Ignatov* [2010] values estimated for present study by author | Ship | 0.81 K | 0.53 K | From comparison with multisensor satellite SST fields |
| *Kennedy et al.* [2011a, 2011b] method using *Atkinson et al.* [2013] whitelist | Ship | 0.56 K | 0.37 K | From comparison with multisensor satellite SST fields |
| *Gilhousen* [1987] Table 6 pg 104 | Mooring | 0.22 K | 0.13 K | Comparison of moored buoys |

2167 **Table 4**: List of estimates of measurement error uncertainties for all platforms for studies

2168 where the measurement error uncertainty is decomposed into random and systematic

2169 components.

2170

| Data set | Input data set | Interpolation method | Resolution |
|---|---|---|---|
| ICOADS summaries [*Woodruff et al.,* 2011] | ICOADS 2.5 | None | 2°x2° monthly |
| HadSST2 [*Rayner et al., 2006*] | ICOADS 2.1 | None | 5°x5° monthly |

| | | | |
|---|---|---|---|
| HadSST3 [*Kennedy et al.,* 2011b; *Kennedy et al.,* 2011c] | ICOADS 2.5 | None | 5°x5° monthly |
| TOHOKU [*Yasunaka and Hanawa*, 2002] | ICOADS 2.1 | None | 5°x5° monthly |
| HadISST1.1 [*Rayner et al.*, 2003] | Met Office Marine Databank and COADS, AVHRR satellite retrievals | Reduced Space Optimal Interpolation | 1°x1° monthly |
| ERSSTv3b [*Smith et al.*, 2008] | ICOADS 2.1 | Separate low and high frequency reconstructions. High frequency component based on EOTs | 2°x2° monthly |
| COBE [*Ishii et al.*, 2005] | ICOADS 2.1 and Kobe collection | Optimal interpolation | 1°x1° monthly |
| COBE-2 [*Hirahara et al.,* 2013] | ICOADS 2.5 and Kobe collection, AVHRR satellite retrievals | Multi scale analysis based on EOFs | 1°x1° daily and monthly |

| | | | |
|---|---|---|---|
| Kaplan [*Kaplan et al.,* 1998] | Met Office Marine Databank | Reduced Space Optimal Smoothing | 5°x5° monthly |
| NOCS [*Berry and Kent*, 2011] | ICOADS 2.5 | Optimal Interpolation | 1°x1° daily and monthly |
| VBPCA [*Ilin and Kaplan*, 2009] | ICOADS 2.5 | Variational Bayesian Principal Component Analysis | 5°x5° monthly |
| GPFA [*Luttinen and Ilin,* 2009] | ICOADS 2.5 | Gaussian Process Factor Analysis | 5°x5° monthly |
| GP [*Luttinen and Ilin,*2012] | ICOADS 2.5 | Gaussian Process | 5°x5° monthly |

2171 **Table 5**: List of datasets used and referred to in the review.

2172

2173 **Figure Captions**

2174

2175 **Figure 1:** (a) Estimated bias (with respect to the 1961-1990 average) on global average

2176 SST anomalies associated with measurement depth as a function of time (upper panel).

2177 (b) Global average SST anomaly from the HadSST3 [Kennedy et al. 2011b, 2011c]

2178 median before (black) and after (red) the measurement-depth bias has been subtracted.

2179 The two red lines reflect different assumptions concerning data that could not be

2180 definitively assigned to any particular measurement type. The large dip during World

2181 War 2 arises because the majority of observations were ERI measurements.

2182

2183  **Figure 2**: Time series of upper ocean temperatures (0-30 m) from nine moorings in the

2184  Tropical Ocean Atmosphere (TAO) array and the Subduction Array. The mooring and its

2185  location are given above each plot. The different coloured lines represent different depths

2186  and these are indicated by the legends in each panel. The Subduction Array data are

2187  described in *Moyer and Weller* [1997].

2188

2189  **Figure 3:** Distributions of estimated measurement errors and uncertainties from ships. (a)

2190  distributions of systematic measurement errors for all entries (2003-2007) in *Kennedy et*

2191  *al.* [2011a], *Brasnett* [2008], *Berry and Kent* [2008] and *Xu and Ignatov* [2010]. (b)

2192  distributions of random measurement error uncertainties (expressed as variances) from

2193  the same analyses as in the top left panel and *Atkinson et al.* [2013]. (c) as for top left

2194  except each ship now has only a single entry so the analyses are directly comparable. (d)

2195  scatter plot showing systematic measurement errors estimated by *Brasnett* [2008] and

2196  *Berry and Kent* [2008] showing the good correlation between the estimates.

2197

2198  **Figure 4:** (a) Estimated global average SST anomaly from HadSST3 [Kennedy et al.

2199  2011b, 2011c] (red) and for subsamples of the HadSST3 dataset reduced to 19th century

2200  coverage. The black line is the median of the samples and the blue area gives the range.

2201  (b) difference, on an expanded temperature scale, between the global average SST

2202  anomaly from the full HadSST3 data set and global averages calculated from the

2203  subsamples.

2204

2205    **Figure 5:** Global average sea-surface temperature anomalies and night marine air

2206    temperature anomalies from a range of data sets. (a) Simple gridded SST data sets

2207    including ICOADS v2.1 (red), 200 realizations of HadSST3 (pale grey), HadSST2 (dark

2208    green), TOHOKU (darker grey), ARC (*Merchant et al.* [2012] lime green) and the

2209    COBE-2 dataset sub-sampled to observational coverage (pale blue). (b) 8 Interpolated

2210    SST analyses including the COBE-2 dataset (pale blue), HadISST1.1 (gold), ERSSTv3b

2211    (orange), VBPCA, GPFA and GP (deep magenta), Kaplan (pink), NOCS (black). (c)

2212    shows the series in (a) and (b) combined. (d) NMAT: Ishii et al. (2005, red and blue),

2213    MOHMAT4N3 and HadMAT (*Rayner et al.* [2003], pink and orange), *Berry and Kent*

2214    [2009] (green), HadNMAT2 (*Kent et al.* [2013], gold).

2215

2216    **Figure 6**: Comparison between COBE-2 (black) and HadSST3 (red) metadata and bias

2217    estimates for the period 1920 to 2010. (a) Fraction of buckets assessed as being

2218    uninsulated. The two red lines indicate the earliest and latest switchover dates allowed in

2219    the generation of the HadSST3 ensemble. (b) Fractional contribution to the global

2220    average from buckets, buoys and engine room measurements. The total is less than unity;

2221    the remainder are either unknown (in the HadSST3 analysis) or uncategorized (COBE-2).

2222    (c) Estimated bias. There are 100 versions of HadSST3 and a single estimate from

2223    COBE-2.

2224

2225    **Figure 7**: Maps showing climatological standard deviation of SST (a, g, m), Structural

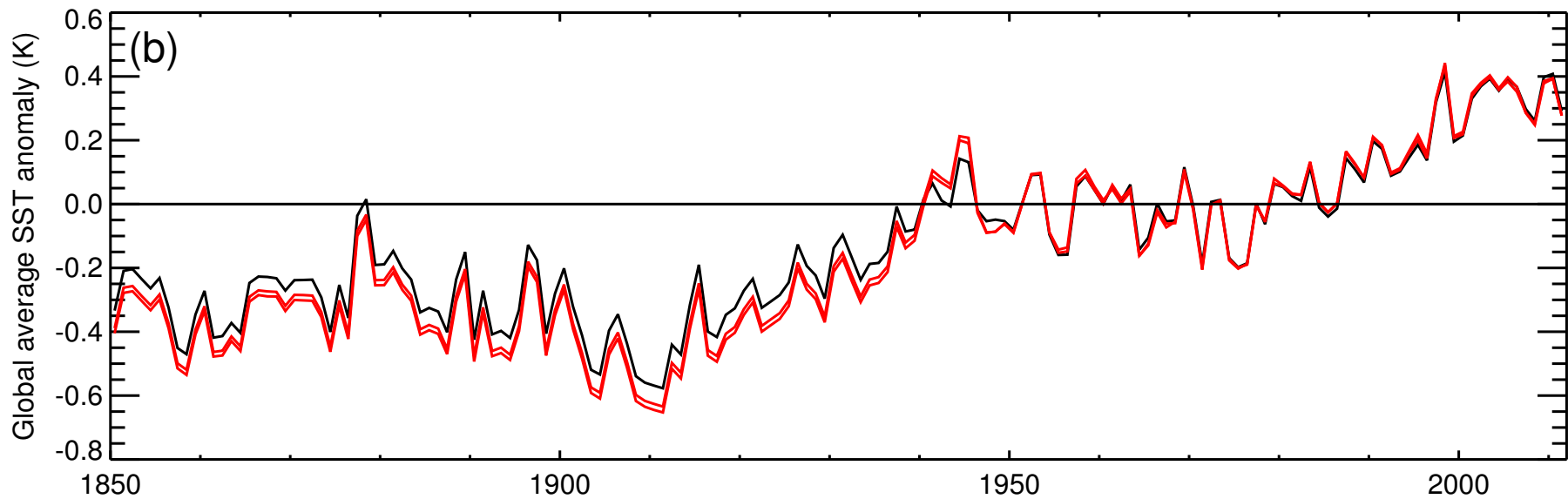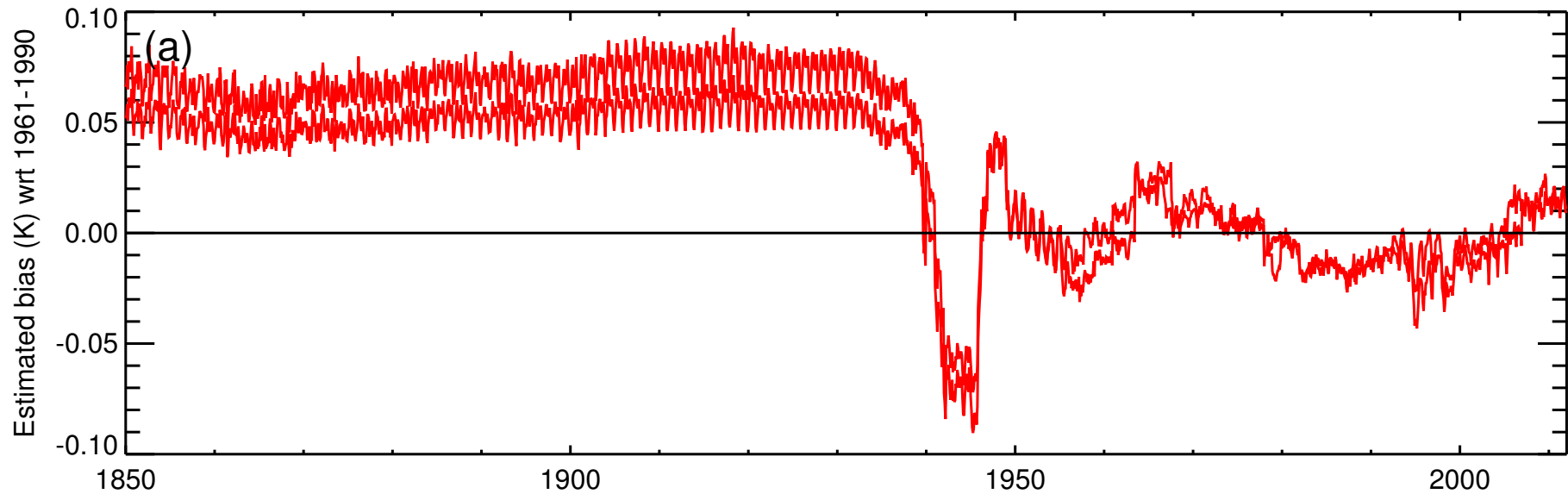2226    uncertainty (b, h, n), Sampling uncertainty (c, i, o), measurement uncertainty (d, j, p), bias

2227   uncertainty (e,k,q) and analysis uncertainty from ERSST (f, l, r). Three months are

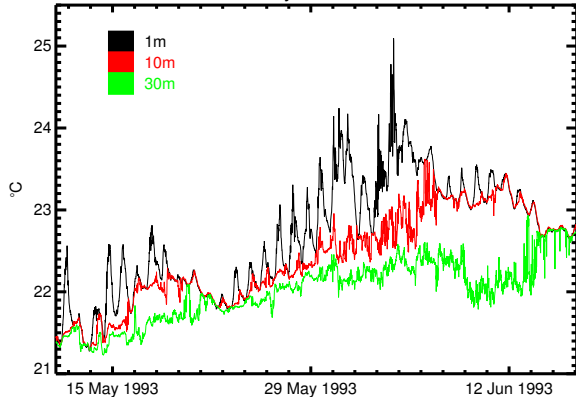2228   shown: (a-f) June 1891, (g-l) April 1944 and (m-r) August 2003.

2229

2230   **Figure 8**: Time series of estimated uncertainties arising from different sources in area-

2231   averages: (a) Global annual, (b) Northern hemisphere annual, (c) North Pacific annual,

2232   (d) North Atlantic annual and (e) a 5-degree grid box centered on 42.5°W, 27.5°N

2233   monthly. Uncertainty components shown are: (pale blue) grid-box sampling uncertainty,

2234   (green) uncorrelated measurement uncertainty, (red) correlated measurement uncertainty,

2235   (dark blue) parametric bias uncertainty from a 200-member ensemble based on HadSST3,

2236   (black) large-scale sampling uncertainty, and (magenta) structural uncertainty estimated

2237   by taking the range of the area-average calculated from seven near-globally-complete
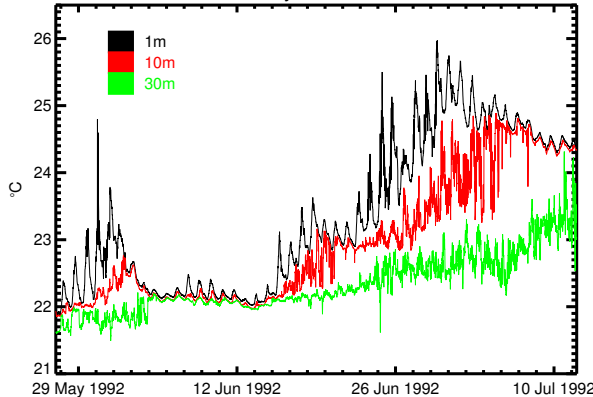
2238   analyses.

2239

2240   **Figure 9**: (a) Global, (b) Northern Hemisphere, (c) Southern Hemisphere and (d)

2241   Tropical average sea-surface temperature anomalies with estimated 95% confidence

2242   range for ERSSTv3b (1880-2012 dark blue line and pale blue shading) and for the

2243   HadSST3 based analysis described in section 3.5 (1850-2011 red line and orange and

2244   yellow shading). The yellow shading indicates an estimate of the additional structural
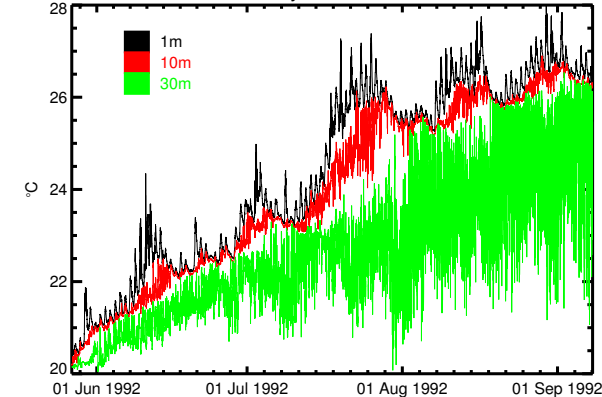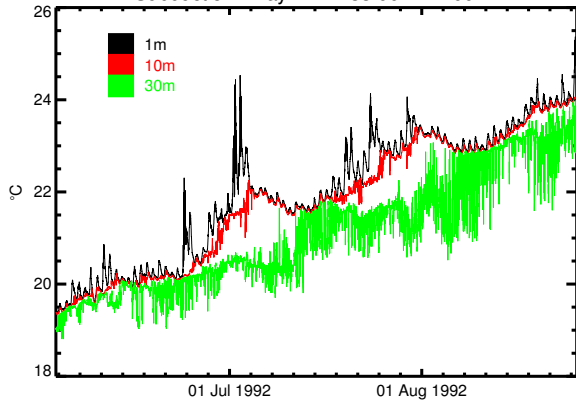
2245   uncertainty in the HadSST3 series.

(a)

(b)

**(a) Systematic Error (all values)**

Probability vs K
- K11 0.67°C
- B08 0.93°C
- BK08 0.89°C
- XI10 0.65°C

**(b) Random Error Variance**

Probability vs $K^2$
- K11 0.91°C
- B08 1.16°C
- BK08 1.04°C
- XI10 0.81°C
- A12 1.05°C

**(c) Systematic Error (ships common to all analyses)**

Probability vs K
- K11 0.67°C
- B08 0.69°C
- BK08 0.67°C
- XI10 0.53°C
- A12 0.43°C

**(d) Systematic Error vs Systematic Error**

Berry and Kent (2008) vs Brasnett (2008)

(a) Sea-surface Temperature: 5+200 datasets

(b) Sea-surface Temperature Analyses: 8 datasets

(c) Sea-surface Temperature: 13+200 datasets

(d) Marine Air Temperature: 6 Datasets

(a)

Fraction of buckets that are uninsulated

(b)

Fractional contribution

ERI

Bucket

Buoy

(c)

Estimated bias (K)

(a) Standard deviation June 1891    (b) Structural Uncertainty    (c) Grid-box Sampling Uncertainty

(d) Measurement Uncertainty    (e) Bias Uncertainty    (f) ERSST analysis uncertainty

(g) Standard deviation April 1944    (h) Structural Uncertainty    (i) Grid-box Sampling Uncertainty

(j) Measurement Uncertainty    (k) Bias Uncertainty    (l) ERSST analysis uncertainty

(m) Standard deviation August 2003    (n) Structural Uncertainty    (o) Grid-box Sampling Uncertainty

(p) Measurement Uncertainty    (q) Bias Uncertainty    (r) ERSST analysis uncertainty

(a) Globe
Grid-box Sampling
Measurement (assumed uncorrelated)
Measurement (assumed correlated)
Bias
Large-scale Sampling
Structural

(b) Northern Hemisphere

(c) North Pacific

(d) North Atlantic

(e) Grid box

(a) Globe

(b) Northern Hemisphere

(c) Southern Hemisphere

(d) Tropics (20°S-20°N)