

1 **Creating surface temperature datasets to meet 21st Century challenges**

2
3 **Met Office Hadley Centre, Exeter, UK**

4
5 **7th-9th September 2010**

6
7 **White papers background**

8
9 Each white paper has been prepared in a matter of a few weeks by a small set of
10 experts who were pre-defined by the International Organising Committee to represent
11 a broad range of expert backgrounds and perspectives. We are very grateful to these
12 authors for giving their time so willingly to this task at such short notice. They are not
13 intended to constitute publication quality pieces – a process that would naturally take
14 somewhat longer to achieve.

15
16 The white papers have been written to raise the big ticket items that require further
17 consideration for the successful implementation of a holistic project that encompasses
18 all aspects from data recovery through analysis and delivery to end users. They
19 provide a framework for undertaking the breakout and plenary discussions at the
20 workshop. The IOC felt strongly that starting from a blank sheet of paper would not
21 be conducive to agreement in a relatively short meeting.

22
23 It is important to stress that the white papers are very definitely not meant to be
24 interpreted as providing a definitive plan. There are two stages of review that will
25 inform the finally agreed meeting outcome:

- 26 1. The white papers have been made publicly available for a comment period
27 through a moderated blog.
- 28 2. At the meeting the approx. 75 experts in attendance will discuss and finesse
29 plans both in breakout groups and in plenary. Stringent efforts will be made
30 to ensure that public comments are taken into account to the extent
31 possible.

32

33 **Creation of quality controlled homogenised datasets from the databank**

34
35 Blair Trewin, Xiaolan L. Wang, Phil Jones, Matthew Menne, Robert Lund, Graham
36 Machin

37 Remit

- 38
- 39
- 40 • Whether different methods will be required for different timescales (monthly,
- 41 daily, sub-daily).
- 42 • Whether some common metrics would be appropriate.
- 43 • Whether efforts should be focussed on a particular area or timescale first, and
- 44 if so, why.
- 45 • The role of national and regional assessments.
- 46 • How to engender multiple efforts.
- 47 • Potential for novel approaches (e.g. an open source community effort, the use
- 48 of reanalyses output instead of neighbours as the expectation field, Bayesian
- 49 approaches).
- 50 • The importance of a consistent approach to assessing uncertainty.

51 Introduction and terminology

52
53
54 Inhomogeneities in temperature records can arise for a wide variety of reasons (for a
55 broad review of these see Trewin (2010)). The process of creating a homogenised
56 dataset from the databank involves two principal stages: the **detection** of
57 inhomogeneities (also known as **changepoints or shifts**) in the data, and making
58 **adjustments** to remove those inhomogeneities and create a homogeneous dataset. For
59 some applications only detection is required, with users making their own judgements
60 about adjustments (if any). The extent to which users wish to remove inhomogeneities
61 may also vary depending on the application: for example, in detection of global
62 climate change signals, it is desirable to remove any inhomogeneities arising from
63 urbanisation, but some users may be explicitly interested in anomalous local trends
64 arising from urban growth.

65
66 There are multiple approaches in the literature to both the detection and adjustment
67 problem. In order to preserve the real climate signal (trend and low-frequency
68 variations) in the series under homogeneity study (the candidate series), most
69 approaches compare the **candidate series** to one or more **reference series**. The
70 reference series should contain the same regional climate signal as the candidate
71 series and be highly positively correlated with the candidate series. The reference
72 series most commonly used in currently published methods for surface temperature
73 data are based on data from a neighbouring station, or a combination of neighbouring
74 stations (more in the section on Reference series)).

75
76 **Quality control** refers to the identification of errors in an individual observation, or
77 small number of observations (such as might arise from, for example, a clerical or
78 observation error, or a short-term instrument fault). These are sometimes
79 characterised as random (as opposed to systematic) errors. This has often been a
80 neglected area in the development of homogenised data sets. Whilst most national
81 databanks now have quality control procedures implemented, this has not necessarily
82 been the case historically, and an important part of the development of a homogenised

83 dataset is applying, as far as possible, a level of quality control to the existing data
84 which is comparable to that applied to current data. Particularly important in the
85 context of extremes is not automatically rejecting observations purely because they
86 are outliers in a statistical distribution. Large random data errors may also complicate
87 detection and adjustment of inhomogeneities. The boundary between data quality
88 issues and inhomogeneities may sometimes be blurred (e.g. in a case where an
89 instrument is out of calibration over a period of several months) and such cases
90 require decisions as to whether the data can be adjusted or should be rejected
91 altogether.

92
93 Detection and adjustment of inhomogeneities are complex problems. Metadata are
94 often incomplete or non-existent, and the majority of long-term time series used in
95 climate change analyses will have multiple inhomogeneities, which substantially
96 increases the statistical challenges of detection. Creating a ‘perfect’ data set from
97 historical data is a practical impossibility, as there is a lower bound below which
98 inhomogeneities, and errors in individual observations, are effectively undetectable.
99 Quantifying this lower bound gives important information as to one source of
100 systematic (type B) uncertainty in the data. Where inhomogeneities are identified,
101 there is also an uncertainty in determining their size and hence any necessary
102 adjustments.

103 104 Methods of detection of inhomogeneities

105
106 Many changepoint detection methods have been developed and applied to time series
107 of climate data. Most commonly used are likelihood ratio tests, such as the Standard
108 Normal Homogeneity (SNH) test and its variants, two-phase regression (TPR) based
109 tests, Potter’s method, etc. There are also tests that use penalized likelihood criteria
110 (e.g., Akaike’s information criteria), as well as CUSUM and nonparametric tests, and
111 Bayesian approaches. Reeves et al. (2007) and Peterson et al. (1998) give relatively
112 comprehensive reviews on this topic although there are more recent developments.
113 The regression based tests are most powerful when the assumed mean structure (e.g.,
114 no trend in the SNH test, a constant trend in a TPR test...) and normality of errors
115 hold. One should not expect a single method to perform optimally over all forms of
116 mean structures. This stresses the need for using multiple methods. The chosen
117 method(s) should be able to address the most common types of inhomogeneities
118 expected to exist in climate data sets.

119
120 Many of these methods are developed for time series containing “at most one
121 changepoint (AMOC)”. However, a long term climate data time series often contains
122 multiple changepoints. A commonly used method to deal with multiple changepoints
123 is a sort of stepwise testing algorithm, in which the time series being tested is divided
124 into segments by a number of most probable changepoints and then an AMOC sub-
125 series containing two neighbouring segments is tested to determine the most probable
126 position and significance of the changepoint within this sub-series (e.g., Wang 2008).
127 Segmentation methods usually perform reasonably well, especially when the mean
128 shifts in a series all are in the same direction (either increasing or decreasing), but
129 they could become more delicate when the shifts take opposite signs or occur close
130 together. A direct approach to the problem of multiple changepoints is a penalized
131 log-likelihood procedure developed by Caussinus and Mestre (2004).

133 There exist two types of changepoints, depending on whether the time of change is
134 known or unknown. As pointed out by Lund and Reeves (2002), test statistics for
135 detecting an unknown changepoint are different from those for assessing statistical
136 significance of a known changepoint. This is because detection of an unknown
137 changepoint involves a search for the most probable time of change by maximizing
138 the related statistic and thus the test statistic has an extreme type of distribution (much
139 higher percentiles), while such a search is not needed when the time of change is
140 known (documented in metadata). However, apart from the RHtestsV3 package
141 (Wang and Feng 2010), few other packages include tests for both known and
142 unknown changepoints. Hierarchical Bayesian methods with a prior chosen to depend
143 on a metadata record have the potential to consolidate both known and unknown
144 changepoints, but have not been extensively developed to date.

145

146 Such techniques are generally designed to detect changepoints at a specific point in
147 time. Most are less suited to detecting cases where an anomalous local trend may
148 develop over a period of time, for example as a result of urbanisation (although the
149 extent to which urbanisation would be manifested as an anomalous trend, as opposed
150 to one or more step changes as conditions change in the vicinity of the observation
151 site, is not fully resolved).

152

153 Reference series and network-wide inhomogeneities

154

155 Most published homogenisation methods for temperature have used one or more
156 reference series based on a number of neighbouring stations. This may involve, for
157 example, the estimation of a background field at the location of the candidate station
158 from a distance-weighted mean of neighbouring stations, or the pairwise comparison
159 of the candidate station with neighbours. The best choice of methods is an area of
160 active discussion.

161

162 A good reference series should contain the same regional climate signal as the
163 candidate series and be homogeneous, or at least be homogeneous in a sub-period in
164 which the candidate series is likely inhomogeneous when a pair-wise comparison
165 method is used. The uncertainty in the extent to which the reference series represents
166 the regional climate signal in the candidate series is a source of uncertainty in the
167 homogenization process. The accuracy of individual values in the reference series will
168 also affect the accuracy of the homogenized data when the reference series is involved
169 in deriving the adjustments. In most cases (especially for temperature), it is reasonable
170 to assume that the candidate and reference series have the same regional climate
171 signal; homogeneity can also be assumed, especially over sub-periods of the record.

172

173 Station-based reference series are, however, not appropriate for a situation where a
174 change is implemented across a network at the same time (e.g. a change in instrument
175 type, or a change in observation time). Furthermore, a network-wide change may
176 induce an inhomogeneity which is undetectable at an individual station but may be
177 highly significant in a larger dataset (for example, a hypothetical change in instrument
178 type that caused an inhomogeneity of $+0.2^{\circ}\text{C}$ would have a major impact on a global
179 mean). For known network-wide changes, an option is, where possible, to establish an
180 experimental comparison (e.g. between an old and new instrument type), or a physics-
181 based model linking old and new methods. The former approach has been used in a
182 number of studies to assess the impact of changing from manual thermometers to

183 automated temperature sensors. Comparing instruments with a traceable standard is
184 significant in this context (although in theory all instruments should be traceably
185 calibrated to reliable national standards for the appropriate quantities, in practice this
186 is unlikely to have occurred in all countries or throughout the period of historical
187 record). There have also been attempts to recreate historical thermometer exposures to
188 compare with modern standards (e.g. Böhm et al., 2010). Physics-based modelling has
189 been used to assess the impact of changes in sea-surface temperature measurement
190 methodology.

191

192 While reanalyses have inhomogeneities of their own, some of them are independent
193 of surface data (especially surface air temperature data) and thus may have potential
194 as a reference series for network-wide changes where no other suitable reference
195 series exist. Because upper-air temperatures (on which reanalyses are based) generally
196 have longer decorrelation length scales than surface temperatures, reanalyses may
197 also be of value as a reference series where no good surface neighbours exist (e.g.
198 sparsely populated regions and remote islands). For coastal and island stations, sea
199 surface temperatures may be another reference series possibility. As yet no major
200 published data set has used either method.

201

202 Methods of adjustment to remove inhomogeneities

203

204 Many existing homogenised datasets at the national and international level take one of
205 two approaches: either the exclusion of stations found to be inhomogeneous, or the
206 application of a single adjustment for each inhomogeneity, applied uniformly across
207 the year. More recently, a wider range of adjustment techniques have been put
208 forward. These include:

209

- 210 (a) Calculation of adjustments separately for each month, or season.
- 211 (b) Calculation of adjustments calculated from monthly data but smoothed across
212 the annual cycle with a different adjustment for each date (widely described in the
213 literature (e.g. Brunet et al, 2007) as ‘daily adjustment’). A number of national
214 and regional datasets have used this method.
- 215 (c) Calculation of different adjustments for different parts of the daily temperature
216 frequency distribution, or for different weather types. Such methods have been
217 applied to a number of test sets of stations, but Australia and Canada are the only
218 countries known to have produced a national dataset adjusted in this way.

219

220 Historically, adjustment methods have received less attention in the literature than
221 detection methods, but are currently an active area of research, in particular through
222 the European HOME project (www.homogenisation.org).

223

224 Appropriate timescales for homogenisation

225

226 Temperature data exist on a number of timescales, the most commonly-used being
227 annual, monthly, daily and sub-daily. This raises issues such as:

228

- 229 • At which timescale can inhomogeneities most effectively be detected, and
230 adjusted for? (this will not necessarily be the same for detection and
231 adjustment).

- 232 • Should homogenisation be carried out on a ‘base’ set of data from which
233 further variables are derived (e.g., a daily maximum/minimum temperature
234 dataset from which daily, monthly and annual means can be derived), or
235 should different variables be homogenised separately?
236

237 Most published work has involved detection at an annual or monthly timescale. Initial
238 results from the HOME project suggest that this is likely to be a more effective
239 approach than detection on shorter timescales, with signal-to-noise being an issue.
240 Note that detection power is inversely related to sample size (series length); and larger
241 uncertainty is usually associated with statistical estimates from smaller samples.
242 Annual data series could be too short to detect changepoints with acceptable
243 uncertainty.
244

245 While a number of existing datasets involve different variables being homogenised
246 separately, a major issue with such an approach is that they may no longer be
247 internally consistent (e.g. a monthly mean may not be the mean of the daily values, or
248 a daily mean may no longer be consistent with the maximum and minimum).
249

250 The homogenisation of fixed-hour sub-daily observations is a largely unexplored
251 scientific question, although some attention has been given to homogenising daily
252 means which are based on fixed-hour observations, in those countries where daily
253 means are calculated using that method, and there has been some limited work done
254 on assessing sub-daily inhomogeneities using weather types.
255

256 Scale of homogenisation and quality control efforts

257

258 It is likely that homogenisation work will be most effectively carried out at the
259 national/regional level, for a number of reasons:
260

- 261 • Researchers working within their own country or region are likely to have
262 access to a greater range of data (e.g. additional stations, or more daily/sub-
263 daily observations) than are available internationally. They are also likely to
264 have access to a wider range of metadata (including pictures, results of
265 calibration checks etc.), and language difficulties will be less of a factor in
266 interpreting that metadata.
- 267 • National-level researchers are more likely to be familiar with the local
268 geography and climate around their observing locations. This is particularly
269 important in quality control in assessing what level of spatial variation
270 between sites is reasonable.
- 271 • Homogenisation and quality control can be a very labour-intensive process,
272 and carrying out the work on a national/regional basis limits the resources that
273 need to be committed by any one institution.
- 274 • Carrying out work at a national/regional level will give national institutions
275 greater ownership of the project, and may help smooth the path in resolving
276 data policy issues.
277

278 The development of a consistent framework between nations/regions will be an
279 important part of this project. The use of inconsistent methods between regions raises
280 the possibility of inconsistencies between data sets once national/regional data sets are
281 consolidated into a global set.

282

283 A significant question which will need to be addressed is that of the relative merits of
284 fully automated homogenisation and quality control methods, and those with some
285 level of manual intervention. Automated methods have the advantages of, for
286 example, being much less resource-intensive than manual methods, being fully
287 reproducible, and of being much more amenable to regular updating. It remains to be
288 determined whether such methods are capable of matching, or at least approaching to
289 within an acceptable level, the accuracy of more manually-intensive methods.

290

291 The potential for carrying out the work on a more distributed basis still (e.g. through
292 volunteer individuals) is an interesting one which is worth exploring, but would
293 require the development of suitable tools and training. Such approaches may be better
294 deployed to areas such as digitisation in the first instance.

295

296 Uncertainty assessment

297

298 Very limited attention has been given to uncertainty assessment in existing
299 homogenised data sets, although some, mostly theoretical, attention has been given to
300 the question of what the minimum detectable inhomogeneity is in any given time
301 series. This is unfortunate because no measurement is properly complete without a
302 rigorous assessment of its associated uncertainty.

303

304 A proper assessment of uncertainty, including uncertainties arising from the data itself
305 and those uncertainties associated with any homogenisation procedure, will be an
306 important aspect of the development of homogenised datasets in this project. The
307 uncertainty, properly quantified, will give significantly enhanced confidence in the
308 data to researchers who use the data products and, probably more importantly, to
309 those outside of the community. Potential approaches include:

310

- 311 • The use of multiple reference series (e.g. through using different combinations
312 of neighbouring stations).
- 313 • The use of different detection and/or adjustment methods on the same dataset,
314 or alternatively the comparison of two different methods used in neighbouring
315 regions (e.g. cross-border comparisons of two different national data sets in
316 adjoining countries).

317

318 Possible metrics include the proportion of techniques which detect a known (from
319 metadata) inhomogeneity, or the spread of magnitude of inhomogeneities detected
320 using different methods or reference series.

321

322 To retrofit a fully assessed uncertainty analysis to past data is very difficult, if not
323 impossible, however a rigorous approach to quantifying the uncertainty will go a long
324 way to improving confidence in the data. To do this the following will need to be
325 performed:

326

- 327 • Understanding quantitatively the influence factors and their relative
importance on the measured quantities
- 328 • Establishment of a uniform approach to quantifying uncertainty in data
329 analysis following internationally accepted guidelines (e.g. the ISO Guide to
330 Uncertainty in Measurement) (for example developing a transparent way of

331 identifying and addressing discontinuities in (non-contiguous) data sets and
332 outliers)

333

334 Recommendations

335

336 • To use daily maximum/minimum temperature as the ‘base’ data set to which
337 adjustments are made, with data at monthly and longer timescales derived
338 from the daily data (adjusted where appropriate) rather than adjusted
339 separately.

340 • To ensure that all detection and adjustment of inhomogeneities is fully
341 documented, allowing reassessments to be made in the future (e.g. if new
342 techniques are developed or previously unknown data or metadata become
343 available).

344 • To carry out an objective evaluation of known methods for
345 homogenisation/adjustment, in collaboration with the COST action;

346 • To establish a testbed of data for this purpose (see white paper 9);

347 • To seek to ensure that all sources of uncertainty are well quantified and
348 defined.

349

350 References

351

352 Böhm, R., Jones, P.D., Hiebl, J., Frank, D., Brunetti, M. and Maugeri, M. 2010. The
353 early instrumental warm bias: a solution for long central European temperature
354 series 1760-2007. *Climatic Change*, 101, 41-67.

355 Brunet M, Saladié O, Jones P, Sigró J, Aguilar E, Moberg A, Lister D, Walther A,
356 Almarza C, 2007: A case-study/guidance on the development of long-term
357 daily adjusted temperature datasets. WMO/TD No. 1425, World
358 Meteorological Organization, Geneva. Available online at
359 [http://www.wmo.int/pages/prog/wcp/wcdmp/wcdmp_series/documents/WCD](http://www.wmo.int/pages/prog/wcp/wcdmp/wcdmp_series/documents/WCDMP_Spain_case_study-cor_ver6March.pdf)
360 [MP_Spain_case_study-cor_ver6March.pdf](http://www.wmo.int/pages/prog/wcp/wcdmp/wcdmp_series/documents/WCDMP_Spain_case_study-cor_ver6March.pdf).

361 Caussinus, H., and O. Mestre, 2004: Detection and correction of artificial shifts in
362 climate. *Journal of the Royal Statistical Society, Series C*, 53, 405-425.

363 Lund, RB, and Reeves J, 2002: Detection of undocumented changepoints --- a
364 revision of the two-phase regression model. *Journal of Climate*, 17, 2547-
365 2554.

366 Peterson, T. C. and co-authors, 1998: Homogeneity adjustments of in situ atmospheric
367 climate data: A review. *Int. J. Climatol.* 18: 1493-1517.

368 Reeves J, Chen J, Wang XL, Lund R, Lu Q, 2007: A review and comparison of
369 changepoint detection techniques for climate data. *J. Appl. Met. Climatol.*, 46:
370 900-915, doi: 10.1175/JAM2493.1.

371 Trewin, B.C. 2010: Exposure, instrumentation and observing practice effects on land
372 temperature measurements. *Wiley Interdisciplinary Reviews: Climate Change*,
373 published online at
374 <http://www3.interscience.wiley.com/journal/123356866/abstract>.

375 Wang, XL, Y. Feng, 2010: RHtestsV3 User Manual. Climate Research Division,
376 Environment Canada. 28pp. Available at
377 <http://cccma.seos.uvic.ca/ETCCDMI/software.shtml>

378 Wang, XL, 2008: Accounting for autocorrelation in detecting mean-shifts in climate
379 data series using the penalized maximal *t* or *F* test. *J. App. Meteor. Climatol.*,
380 47, 2423–2444. DOI: 10.1175/2008JAMC1741.1.