

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 through
27 a moderated blog.
28 2. At the meeting the approx. 75 experts in attendance will discuss and finesse plans
29 both in breakout groups and in plenary. Stringent efforts will be made to ensure that
30 public comments are taken into account to the extent possible.

31

32 **Benchmarking homogenisation algorithm performance against test cases**

33

34 ¹Kate Willett, ²Matt Menne, ²Peter Thorne, ³Stefan Brönnimann, ⁴Ian Jolliffe, ⁵Lucie
35 Vincent & ⁵Xiaolan Wang

36

37 ¹ Met Office Hadley Centre, FitzRoy Road, Exeter, UK

38 ² National Climatic Data Center, Ashville, NC, USA

39 ³ University of Bern, Switzerland

40 ⁴ University of Exeter, Exeter, UK

41 ⁵ Climate Research Division, Science and Technology Branch, Environment Canada,
42 Toronto, Canada

43

44

45 **Discussion Areas:**

- 46 i. The need for pseudo-data to have full space and time sampling of the
47 observational network
- 48 ii. Source data for pseudo-data construction
- 49 iii. Design of orthogonal discontinuity test cases that can optimally answer
50 multiple questions
- 51 iv. Should discontinuity test case construction be truly blind to the dataset
52 creators to avoid over-tuning of algorithms to the discontinuity test cases?

53

54 **Introduction:**

55 For climate change research it is essential to detect, attribute (if possible) and, adjust
56 for discontinuities (where desirable depending on end-user) originating from changes
57 in station location, instrumentation, recording practice or surrounding environment.
58 Reconciliation of discontinuities using station metadata is preferable but these data
59 are mostly undigitised or in many cases undocumented even on paper. The ability of
60 chosen algorithms to achieve station homogeneity in the absence of complete and
61 verified metadata (methodological uncertainty) has rarely been assessed
62 comprehensively so as to enable useful dataset intercomparison. Algorithms vary in
63 their skill at identifying discontinuities that are: multiple and geographically or
64 temporally clustered; close to end points; gradual; variance based; in the presence of a
65 background trend; and seasonally or diurnally variant. In addition, quantifying
66 discontinuity magnitude (and hence the presence of discontinuities and adjustments
67 required) against background noise (natural variability and white noise of erroneous
68 observations remaining even after quality control) is very difficult, especially if
69 discontinuity magnitude varies diurnally and seasonally. Even after detection a series
70 of decisions are required as to whether and how to adjust. These decisions can have a
71 further non-negligible impact upon the resulting dataset estimates. While decisions
72 are as evidence based as is physically possible, some are unavoidably arbitrary. This
73 is especially problematic for large datasets where the whole process by necessity is
74 automated.

75

76 Benchmarking (measuring performance against a standard) against homogenous real
77 or synthetic data series with known discontinuity test cases applied is increasingly
78 common practice (Easterling & Peterson 1995; Vincent 1998; Ducre-Robitaille et al.
79 2003; Wang et al 2007; Begert et al. 2008; Wang 2008a; Wang 2008b). However,
80 these pseudo-datasets are usually very simplistic e.g. a single timeseries lacking real
81 world spatio-temporal sampling, climatology, variability and noise. Furthermore,

82 reconstructing realistic discontinuity test cases that capture all eventualities to the best
83 of our knowledge is a very large task indeed, especially for datasets covering multiple
84 regions. Recent studies (e.g. Menne & Williams 2005; DeGaetano 2006; Wang et al.
85 2007 and Wang et al. 2008a; 2008b) have generated synthetic pseudo-data with
86 varying degrees of real-world characteristics (e.g. variance, autocorrelation and
87 randomly assigned sign or position of discontinuity). Pseudo-data have also been
88 generated from homogeneous series of real data using a re-sampling technique (e.g.
89 bootstrap; Wang et al. 2010). Recent European efforts from COST-ESO601⁶ have
90 produced a monthly mean benchmark dataset compiling synthetic and real data with a
91 number of discontinuity test cases. Titchner et al. 2009 created a monthly pseudo-
92 dataset by sub-sampling the atmosphere-only general circulation model (GCM)
93 HadAM3 and adjusting the variance, mean and white noise to match real world
94 characteristics. Various discontinuity test cases were explored. Recent work between
95 the Met Office Hadley Centre, University of New South Wales and the National
96 Climatic Data Center has generated a number of daily pseudo-datasets and
97 discontinuity test cases from the GCM HadCM3 under control, natural and all
98 forcings scenarios. Sampling density, variance, climatology and autocorrelation are
99 adjusted to match the HadGHCND daily temperature record (maximum and
100 minimum), and then various scenarios of discontinuities added.

101

102 At present, no agreed global standard exists against which to benchmark multiple
103 datasets. Worse still, existing benchmarks are usually created by the dataset creators
104 themselves leaving potential for unintentional tuning of algorithms or test cases. The
105 issue is becoming increasingly critical. At present no homogenised sub-daily product
106 exists, yet there is an ever growing need for such data. We therefore conclude the
107 need for a global collaborative effort to: create globally applicable pseudo-data that
108 sufficiently represent real world characteristics; comprehensively capture real world
109 characteristics of inhomogeneity for creation of realistic discontinuity test cases; and
110 objectively test all candidate homogenisation procedures produced against this
111 pseudo-data with a realistic suite of discontinuity test cases applied. This effort should
112 be independent from any single group of climate dataset creators (but may include
113 some people with expertise from various institutions).

114

115 **Discussion of the main issues:**

116

117 *Reconstructing full space and time sampling of the observational network* 118 *characteristics*

119 The majority of existing algorithms have already been skill tested using simple
120 synthetic pseudo-data with simple discontinuity test cases applied which runs the risk
121 of artificially awarding high performance to algorithms that cannot cope with real
122 world natural variability and noise and the variation in sampling density and data
123 completeness/continuity. Climatology (including annual and diurnal cycles), variance
124 (white noise and natural variability), autocorrelation, sampling density (network
125 coverage and missing data) and presence of a trend are real world characteristics of
126 high importance that can affect algorithm skill and generally are not all included in
127 such cases. Hence, to truly test algorithm skill on real data we need global
128 reconstruction of real world characteristics including space and time sampling of the
129 observational network as far as possible. Pseudo-data should be created at a range of
130 resolutions (sub-daily to monthly) where the underlying characteristics and
131 discontinuity test cases applied are identical. This will be particularly relevant to those

132 groups considering algorithm temporal transferability e.g. of a monthly algorithm to
133 the daily timescale.

134

135 *Source data for creation of homogeneous pseudo-data*

136 As discussed, pseudo-data should be physically based on real world characteristics
137 with spatio-temporal correlation structures (e.g. ENSO variability and
138 teleconnections) that exist in the real world. There are three realistic options:

139

- 140 • **High quality homogenous observational datasets**

141 These have the benefit of being real observations but may differ from the
142 majority of stations which due to poorer quality will likely retain more white
143 noise even after quality control. Furthermore, whether these constitute truly
144 homogenous stations where instruments have been regularly calibrated to
145 metrological standards is questionable. Critically, coverage of such stations is
146 likely very small and so of limited use for benchmarking data from across the
147 globe.

148

- 149 • **4th generation reanalyses products**

150 These are globally complete, high resolution (gridded), and represent real
151 world characteristics. The 20th Century Reanalysis⁷ continues back to the early
152 20th century and as it only assimilates sea level pressure does not suffer from
153 some of the major discontinuities apparent at major transitions in the global
154 observing system in some other products although time-dependent biases are
155 still possible. The most recent ECMWF reanalyses appear to adequately
156 reflect surface temperature changes over land at least on the monthly mean
157 timescale (Simmons et al., 2010). However, spurious discontinuities are
158 present from major changes in the observing network (e.g. at the beginning of
159 the satellite era) and post-1997/98 in ERA Interim likely due to a source
160 change to the NCEP operational SST product.

161

- 162 • **CMIP5 GCM output⁸**

163 These are globally complete, high spatial resolution (gridded) and 3 hourly for
164 surface variables although verification compared to real-world data is essential
165 first. They are homogeneous (insofar as they are at least consistent with the
166 basic physics underlying each model and the applied forcings). However, care
167 should be taken to choose a GCM, or preferably multiple GCMs, that
168 represent low-resolution natural variability (e.g. ENSO, PDO etc.) and
169 teleconnections satisfactorily as not all do. There are a range of forcing
170 scenarios including control, high emissions (e.g. A1B) and natural only which
171 can be used to assess some range of dataset algorithm performance with and
172 without underlying trends.

173

174 Spatially and temporally complete gridded fields from reanalyses or GCMs can be
175 sub-sampled and nudged with real world climatology (allowing for algorithms that
176 can work in either absolute or anomaly space), variance, white noise and adjustment
177 for autocorrelation from the databank (see White Papers 3 to 6). Where possible,
178 efforts should collaborate with and build on existing work.

179

180 *Exploring all eventualities with optimum discontinuity test case design*

181 Test cases should be designed in parallel with homogenisation assessment
182 requirements and with the broader assessment discussed in White Paper 10. They
183 should encapsulate physically plausible effects of changes in: station location;
184 instrument; recording practice; or surrounding environment and analyses should lead
185 to clear and useful results. Some discontinuities are well documented and pinned to a
186 period and region (i.e. mid-1990s automation in the USA) and these should be
187 included. However, far more are undocumented and unknown and could be of any
188 magnitude, frequency, clustering or sign bias and are likely a combination of all and a
189 mix of abrupt and more graduated discontinuities. While we can characterise the main
190 features of real world inhomogeneities relatively comprehensively (Elliott 1995,
191 Peterson et al. 1998; Vincent et al. 2007) metadata is vastly incomplete and
192 undigitised. A thorough review is necessary prior to creation of the discontinuity test
193 cases for benchmarking and also to algorithms design (see White Paper 8). This could
194 be facilitated via a relevant conference such as RMS, EGU or AGU. Digitising all
195 available metadata and adding it to the databank is also strongly recommended (see
196 white paper 3).

197
198 Approximately 10 global discontinuity test cases should ideally be constructed for
199 benchmarking, full assessment (see White Paper 10), publication and collation of data
200 products (see White Paper 13). These should be physically plausible scenarios based
201 on our understanding of real world issues that likely pertain and include the control
202 case of a homogeneous world to assess the effect of algorithms giving false positive
203 detections and adjustments. These should incorporate a mix of abrupt, gradual and
204 seasonally/diurnally varying discontinuities. They should methodically address key
205 questions by testing skill under situations of: discontinuity clustering versus sparsity;
206 proximity to endpoints versus midpoints; large versus small discontinuities; a
207 combination of both; and the presence of strong versus no background trend (i.e.
208 taken from control, A1B, c20c and natural climate runs).

209
210 Homogenisation algorithms will produce a new estimate of large scale features
211 (climatology, variance and trends) that will fall somewhere on a benchmark spectrum
212 populated by the inhomogeneous world (pseudo-data with the suite of discontinuity
213 test cases applied) and the target truth (the homogeneous pseudo-data) for the
214 region/regions in question. A perfect algorithm would recreate the target truth across a
215 range of space and time scales. Performance in terms of percentage of successful,
216 missed and incorrect adjustments should also be noted. This assessment will feed into
217 the broader assessment described in White Paper 10. Although successful
218 homogenisation of sub-daily or even daily resolution data has not yet been
219 demonstrated the Surface Temperature initiative will likely spur attempts at this and
220 so the benchmarking dataset needs to be available at high resolution.

221
222 *Avoiding prior knowledge of test cases and resulting over-tuning of algorithms to*
223 *discontinuity test cases*

224 The benchmark dataset (pseudo-data with discontinuity test cases) will be freely
225 available although it is advisable that assessment is done by an independent group
226 (see White Paper 10). In the interests of transparency we argue that all pseudo-data
227 creation, discontinuity test case creation and benchmarking is done independently
228 from any single group of derived dataset creators. A multi-institution group or groups
229 including expertise from non-traditional participants could be set up via the
230 September workshop of those already involved in or with an interest in pursuing such

231 work. Applying to funding bodies will be essential. The methodology underlying this
232 work should be fully documented and published via peer review including the pseudo-
233 data with discontinuity test cases but withholding the ‘solutions’ (the original
234 homogenous pseudo-data). Having a single suite of test cases does still pose the
235 problem of potential to over-tune through multiple iterations of dataset creation in
236 order to reduce uncertainty although it is unlikely that any knowing malpractice of
237 tuning to the test data will occur, especially as full audit trails and code for dataset
238 creation will also be published (White paper 13). Furthermore, there is low likelihood
239 of over-tuning as long as we use a wide range of discontinuity types across the test-
240 cases that are physically based and fully represent real world eventualities.

241

242 **Recommendations:**

- 243 • Global pseudo-data with real world characteristics
- 244 • GCM or Reanalyses data should be used as source base with real spatial,
245 temporal and climatological characteristics applied to recreate to a reasonable
246 approximation the observational record statistics
- 247 • Review of inhomogeneity across the globe finalised via a session at an
248 international conference (link with White Paper 8) to ensure plausibility of
249 discontinuity test cases
- 250 • Suite of ~10 discontinuity test cases that are physically based on real world
251 inhomogeneities and orthogonally designed to maximise the number of
252 objective science questions that can be answered
- 253 • Benchmarking to rank homogenisation algorithm skill in terms of performance
254 using climatology, variance and trends calculated from homogenous pseudo-
255 data and inhomogenous data (discontinuity test cases applied) (skill
256 assessment to be synchronised with broader efforts discussed in White Paper
257 10)
- 258 • Independent (of any single group of dataset creators) pseudo-data creation, test
259 case creation and benchmarking
- 260 • Peer-reviewed publication of benchmarking methodology and pseudo-data
261 with discontinuity test cases but ‘solutions’ (original homogenous pseudo-
262 data) to be withheld

263

264

265 **References:**

266 ⁶ <http://www.homogenisation.org>

267 ⁷ http://www.esrl.noaa.gov/psd/data/gridded/data.20thC_Rean.html

268 ⁸ http://cmip-pcmdi.llnl.gov/cmip5/data_description.html?submenuheader=1

269

270 Begert, M., Zenklusen, E., Haberli, C., et al., 2008: An automated procedure to detect
271 discontinuities; performance assessment and application to a large European climate
272 data set. *Meteorologische Zeitschrift*. **17**, (5), 663-672.

273

274 DeGaetano, A. T., 2006: Attributes of several methods for detecting discontinuities in
275 mean temperature series. *Journal of Climate*. **19** (5), 838-853.

276

277 Ducré-Robitaille, J.-F., Vincent, L. A. & Boulet, G., , 2003: Comparison of
278 techniques for detection of discontinuities in temperature series. *International Journal*
279 *of Climatology*, **23**, 1087-1101.

280

281 Easterling, D. R. & Peterson, T. C., 1995: The effect of artificial discontinuities on
282 recent trends in minimum and maximum temperatures. International Minimax
283 Workshop on Asymmetric Change of Daily Temperature Range, SEP 27-30, 1993
284 COLLEGE PK, MD. *Atmospheric Research*. **37**, 19-26.
285

286 Elliott, W. P., 1995: On detecting long-term changes in atmospheric moisture.
287 International Meeting of Experts on Long-Term Climate Monitoring by the Global
288 Climate Observing System, JAN 09-11, 1995 ASHEVILLE, NC. *Climatic Change*.
289 **31**, 349-367.
290

291 Menne, M. J. & Williams, C. N., 2005: Detection of undocumented changepoints
292 using multiple test statistics and composite reference series. *Journal Of Climate*. **18**,
293 4271-4286.
294

295 Peterson, T. C., Easterling, D. R., Karl, T. R., et al., 1998: Homogeneity adjustments
296 of in situ atmospheric climate data: A review. *International Journal Of Climatology*.
297 **18**, 1493-1517.
298

299 Simmons, A. J., Willett, K. M., Jones, P. D., Thorne, P. W. & Dee, D., 2010: Low-
300 frequency variations in surface atmospheric humidity, temperature, and precipitation:
301 Inferences from reanalyses and monthly gridded observational data sets. *Journal Of*
302 *Geophysical Research-Atmospheres*. **115**, D01110, doi:10.1029/2009JD012442.
303

304 Titchner, H. A., Thorne, P. W., McCarthy, M. P. et al. 2009: Critically Reassessing
305 Tropospheric Temperature Trends from Radiosondes Using Realistic Validation
306 Experiments. *Journal Of Climate*. **22**, 465-485.
307

308 Vincent, L.A., 1998: A technique for the identification of inhomogeneities in
309 Canadian temperature series. *Journal of Climate*, **11**, 1094-1104.
310

311 Vincent, L. A., van Wijngaarden, W. A. & Hopkinson, R., 2007: Surface temperature
312 and humidity trends in Canada for 1953-2005. *Journal of Climate*, **20**, 5100-5113.
313 DOI: 10.1175/JCLI4293.1.
314

315 Wang, X. L., Wen, Q. H., and Wu, Y., 2007: Penalized Maximal t Test for Detecting
316 Undocumented Mean Change in Climate Data Series. *Journal of Applied Meteorology*
317 *and Climatology*. **46**, 916-931. DOI:10.1175/JAM2504.
318

319 Wang, X. L., 2008a: Accounting for autocorrelation in detecting mean-shifts in
320 climate data series using the penalized maximal t or F test. *Journal of Applied*
321 *Meteorology and Climatology*. **47**, 2423–2444. DOI: 10.1175/2008JAMC1741.1
322

323 Wang, X. L., 2008b: Penalized maximal F test for detecting undocumented mean-
324 shift without trend change. *Journal of Atmospheric and Oceanic Technology*, **25**, 368-
325 384. DOI:10.1175/2007/JTECHA982.1.
326

327 Wang, X. L., Chen, H., Wu, Y. et al., 2010: New techniques for detection and
328 adjustment of shifts in daily precipitation data series. *Journal of Applied Meteorology*
329 *and Climatology*. (accepted)
330