

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
11 represent a broad range of expert backgrounds and perspectives. We are very
12 grateful to these authors for giving their time so willingly to this task at such short
13 notice. They are not intended to constitute publication quality pieces – a process
14 that would naturally take 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
18 encompasses all aspects from data recovery through analysis and delivery to end
19 users. They provide a framework for undertaking the breakout and plenary
20 discussions at the workshop. The IOC felt strongly that starting from a blank sheet
21 of paper would not 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 a
27 moderated blog.
- 28 2. At the meeting the approx. 75 experts in attendance will discuss and finesse plans both
29 in breakout groups and in plenary. Stringent efforts will be made to ensure that public
30 comments are taken into account to the extent possible.

31
32

33 White Paper on Dataset algorithm performance assessment based
34 upon all efforts

35
36 Peter Stott¹, Philip Brohan¹, Dick Dee², Alistair Forbes³, Peter Guttorp⁴, Ian
37 Jolliffe⁵, Peter Thorne⁶

38
39 ¹ Met Office Hadley Centre, FitzRoy Road, Exeter, EX1 3PB, UK

40 ² ECMWF, Reading, UK

41 ³ National Physical Laboratory

42 ⁴ University of Washington, Seattle, USA

43 ⁵ University of Exeter, Exeter, EX4 4SB, UK

44 ⁶ National Climatic Data Center, USA.

45

46

47 **Introduction**

48

49 This white paper provides a consideration of the steps that need to be
50 taken to provide suitable performance assessment for algorithms used in
51 the development of climate datasets. It discusses how to make use of the
52 benchmarking results described in White Paper no 9 in addition to other
53 methodological decisions made in order to assess suitability of datasets for
54 a particular purpose. Fundamentally this White Paper is therefore about
55 how well derived datasets represent actual real world behaviour and the
56 attendant uncertainties in that assessment. The technical details of the
57 development and testing of benchmarking homogenization algorithms are
58 provided in White Paper no 9.

59

60 **Considerations for Dataset Algorithm performance assessment**

61

62 In testing and validating dataset algorithm the first step is to determine
63 what the purpose is for which the dataset is being assessed. Only then
64 does it become feasible to determine how good a particular algorithm is for
65 that particular purpose. Such uses could include, for example,
66 determination of long term mean trends, annual maximum values or
67 diurnal variability of a particular climate variable.

68

69 Assessment criteria should be developed independently of algorithm
70 development teams building on suitable expertise gained from a variety of
71 perspectives, ie not just the climate science community. It is strongly
72 desirable that assessment criteria should be determined before any
73 assessments are carried out. Such criteria should be clearly documented
74 and should be approved by the core implementation group who have the
75 task of over-seeing progress on the initiative to develop surface

76 temperature datasets to meet 21st century challenges as set out in White
77 Paper 15.

78

79 It is important that new activities resulting from this initiative are
80 coordinated with on-going national and regional activities to rescue and
81 homogenize data, including for example the European COST-ES0601
82 initiative to produce a benchmark dataset. Further details of some of these
83 initiatives are given in White Paper no 9.

84

85

86 Possible assessment criteria, depending on the use to which the data are
87 being used, could include an ability to identify long term trends as well as
88 to what extent the algorithm has retained the variance of the climate
89 variable and of the ability of the algorithm to identify and adjust for
90 individual breaks in data series as well as a suitable criteria for the number
91 of acceptable false-positive identification of non-existent break points.
92 Further examples of possible assessment criteria include whether
93 simulated data series have non-stationarity in the mean (beyond just
94 seasonal effects), include both short- and longterm "memory" (or temporal
95 correlation) and incorporate realistic continental or global scale spatial
96 correlations. Homogeneity issues are discussed in greater depth in White
97 Paper No 9. They are raised here as examples of assessment criteria that
98 could be appropriate to a particular user and which will therefore need to
99 be assessed.

100

101 The assessed quality of an algorithm for a particular purpose would likely
102 depend on the risks associated with making wrong inferences on the basis
103 of the algorithm outputs. If the risks/uncertainties associated with the
104 outputs of an algorithm are deemed to be too high according to the users,
105 it would be useful to determine what new information would reduce the
106 risks optimally, for example would it be better to replace many inexpensive
107 sensors with limited accuracy with a few highly accurate sensors.

108

109 A key issue is to determine how well uncertainty estimates in datasets
110 provide a measure of the difference between the derived value and the
111 "true" value, recognizing that any uncertainty statement summarises the
112 uncertainty from all factors, including those that operate over different
113 timescales. Establishing traceability is essential as a first step in combining
114 measurements from different measurement systems. Most data analysis
115 algorithms are designed to operate optimally (unbiased, minimum
116 variation) for data generated according to a particular model. For a fully
117 characterised algorithm, the model predicts the behaviour of the algorithm,
118 allowing uncertainties to be associated with the parameter estimates, for

119 example. Therefore assessments should be carried out taking account of
120 the sensitivity of algorithm outputs with respect to underlying model
121 assumptions.

122
123

124 It would be worthwhile to consider the future needs for the development of
125 climate services by identifying an appropriate common set of regions or
126 stations that any assessment should include while recognising that not all
127 dataset creators may include these stations or regions. It is suggested that
128 analyses should be carried out to cover densely, moderately and poorly
129 sampled regions as the performance of any algorithm is likely to depend
130 upon network density. Thought would also be needed as to whether it is
131 appropriate to apply a common set of assessment criteria to all regions or
132 whether different assessment criteria are required for each region.

133

134 Validation of an algorithm should always be carried out on a different
135 dataset from that used to develop and tune the algorithm. As discussed in
136 White Paper No 9, idealised datasets can be created which contain break
137 points, data gaps, etc, for example by sub-sampling and distorting climate
138 model data. Once an algorithm is frozen, the validation dataset can be
139 used to test it by measuring against the pre-determined application-specific
140 assessment criteria.

141

142 White Paper no 9 discussed different approaches for testing
143 homogenization algorithms including in idealized model settings in which
144 potential data gaps, inhomegeneities and other data issues are imposed,
145 reanalyses, and the use of high quality reference datasets. Such model-
146 based tests will always be conditional on the idealized nature of the tests
147 and therefore cannot consider all sources of uncertainty. These
148 uncertainties can be amenable to exploration through the use of
149 reanalyses, which combine observational data with the dynamical
150 constraints of a forecasting model in which multiple climate variables are
151 constrained to be physically consistent. Purely statistical approaches and
152 dynamical approaches (via reanalyses) can be used to validate each other.
153 In addition, a testbed complementary to testing algorithms in models in
154 which deliberate errors have been inserted, is to insert similar such
155 deliberate errors into reference very high quality datasets such as
156 CRN/CRN-HCN dataset which has realistic, irregular grid spacing and
157 variability. Such a test would be good for the smaller scale (country to
158 regional) spatial testing of algorithms as well as for testing spatial-
159 interpolation errors.

160

161 Based on these tests, users should be able to determine which are

162 suitable datasets and approaches for their application. The aim of any
163 assessment should be to determine and flag clearly any datasets that are
164 not fit for a particular purpose and to aid users in determining which are the
165 most suitable datasets for their needs, recognizing importantly that some
166 approaches which may be beneficial for one purpose, eg large-scale trend
167 retrieval, may be detrimental for other features, eg station by station
168 quality. Documentation should aim to provide a clear decision tree for the
169 user, supported by the underlying formal assessments. It should also
170 indicate the extent to which uncertainty assessments provided by the
171 dataset creators provide reliable and robust estimates of the underlying
172 uncertainties for the quantities required by the user for their particular
173 purpose.

174

175

176

177 **Recommendations**

178

179

- 180 • Assessment criteria should be developed entirely independently of
181 the dataset developers and should be pre-determined and
182 documented in advance of any tests.
183
- 184 • It is crucial that the purpose to which a dataset could be put be
185 identified and that a corresponding set of assessment criteria are
186 derived that are suitable for that purpose.
187
- 188 • The output of an assessment should be to determine whether a
189 dataset is fit for a particular purpose and to enable users to
190 determine which are most suitable datasets for their needs. Outputs
191 should be clearly documented in such a form as to enable a clear
192 decision tree for users.
193
- 194 • Validation of an algorithm should always be carried out on a different
195 dataset from that used to develop and tune the algorithm.
196
- 197 • A key issue is to determine how well uncertainty estimates in
198 datasets represent a measure of the difference between the derived
199 value and the “true” real world value.
200
- 201 • It would be worthwhile to consider the future needs for the
202 development of climate services by indentifying an appropriate set of
203 regions or stations that any assessment should include.

204
205
206
207
208
209
210
211
212
213

- New efforts resulting from this initiative should be coordinated with on-going regional and national activities to rescue and homogenize data.