

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 experts
10 who were pre-defined by the International Organising Committee to represent a broad
11 range of expert backgrounds and perspectives. We are very grateful to these authors for
12 giving their time so willingly to this task at such short notice. They are not intended to
13 constitute publication quality pieces – a process that would naturally take somewhat
14 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 all
18 aspects from data recovery through analysis and delivery to end users. They provide a
19 framework for undertaking the breakout and plenary discussions at the workshop. The
20 IOC felt strongly that starting from a blank sheet of paper would not be conducive to
21 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 inform
25 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 **Data provenance, version control, configuration management**

34

35 John R. Christy (UAHuntsville)

36 Nick Barnes (Clear Climate Code)

37 Amy Luers (Google)

38 Shawn Smith (FSU/COAPS)

39 Steve Worley (ICOADS/NCAR)

40 Karl Taylor (CMIP)

41 Jay Lawrimore (NCDC)

42

43 White paper authors are requested to consider: The metadata requirements to ensure
44 traceability back to the original data record; the version controlling requirements (what
45 represents Business As Usual and what constitutes a fundamental update to the
46 databank); how to best retain previous versions of the databank; how to configure the
47 archives to maximize usefulness (data format, file indexing, appending the homogenized
48 data estimates etc.); how to version control and archive the datasets produced from the
49 initial databank (see Section 8).

50

51

Introduction

52

53 Through the years it has become apparent that the development of an internationally-
54 organized, transparent, and comprehensive data management system of surface air
55 temperature has been sorely needed. Nations have stepped forward with extensive efforts
56 to perform these functions, e.g. the U.S. National Climatic Data Center Global Historical
57 Climatology Network (NCDC), the UK Met Office HadCRUT3 surface temperature data
58 set, and the International Comprehensive Ocean-Atmosphere Data Set (ICOADS). These
59 efforts, we believe, can be built upon to achieve an even greater level of effectiveness
60 with international support and oversight for global land surface air temperature
61 measurements managed in a centralized system.

62

63 The emphasis on this white paper is to address the requirements for establishing a data
64 system whereby (a) original observations, (b) information about the observing system
65 which made those observations (metadata), and (c) products generated from those
66 observations, have the potential of meeting the reliability requirements desired not only
67 of the scientific community but of many other communities (e.g. policy, legal, etc.) who
68 now rely on climate data. Such a system will require robust methods of data provenance,
69 version control and configuration management as described below and defined in
70 Appendix A, which in this project relate primarily to Level 0 and 1 data categories (see
71 wp3.) If such a system is established, the scientific requirements regarding climate data
72 activities, including accessibility, traceability, reproducibility and reliability, will be met,
73 and, we will have provided to other communities the information about climate data they
74 need.

75

76 To aid in the production of derived products (e.g. Level 5 such as HadCRUTv3, GISS,
77 ERSST) which may exclude and/or modify the archived, primary-source data through
78 quality control, homogenization, and interpolation algorithms, we shall provide various

79 standardized testing datasets which may be used for assessing the skill of the algorithm.
80 However, it is not the purview of this project to assess derived products as to whether
81 they meet higher standards required by specific communities (i.e. scientific, legal, etc.) It
82 is one aim of this project to archive and disseminate derived products as long as the
83 homogenization algorithm is documented by the peer review process.

84 **Date Provenance (Traceability to primary source)**

85 Primary sources, referred to in this project as Level 0 data (see wp3), of instrumental
86 temperature data fall into two broad categories (a) paper documents and (b) digital
87 computer files from electronic sensors. Establishing an archive that provides a pathway
88 from the heavily-used, uniformly formatted data records (Level 2 and 3 data) in the
89 project’s archive back to their primary source (or relevant information if no primary
90 source exists) is a critical function of the project.

91 **Paper Documents**

92 The earliest primary-source climate records, and even many today, consist of hand-
93 written or machine-printed values, or pin-traces on paper documents. Scientific (and
94 other) research is performed using digital files created from these data after they have
95 been electronically keyed, becoming Level 1 data.

96 Associated with these primary-source paper documents are many metadata documents
97 that describe the instrument location (including maps), type of instrument, condition of
98 site and instrument, directions for taking observations, calibration of instruments, etc.
99 These metadata will also be archived according to “Levels” as described in wp3. This
100 metadata information will require archiving as we anticipate that one of the major uses of
101 our primary-source archive and its associated data records is the construction of
102 homogenized long-time series for which metadata are vital.

103 One goal of this project is to make all images of the primary-source (Level 0) paper
104 documents available to investigators to address traceability and authenticity.
105 Unfortunately, in many cases the primary source documents have been lost, destroyed or
106 for some reason have become unreadable. In their place quite often are secondary-
107 source documents (i.e. official monthly summaries, newspaper reports, etc.) or digital
108 files that may have been derived from a primary-source document before its demise.
109 These data are known as Level 1 data but which may not have traceability to an archived
110 Level 0 document or file. Traceability and authenticity in these cases are more difficult.
111 Thus, there are two types of Level 1 data – that which is traceable to an archived Level 0
112 primary source and that which is not. As indicated below, there will likely be multiple
113 versions of Level 1 data due to the loss of Level 0 for a particular station.

114 **Electronically measured and transmitted data**

115 There has been a relatively rapid conversion from printed data collection methods to
116 electronic measuring, reporting, and quality control so that the human eye never

117 witnesses the observation nor its transmission and archival. In some of these cases, the
118 electronically measured observation is produced originally as a geophysical parameter
119 and reported on an electronic network, usually in an obscure, digitally-packed file
120 structure. In these cases it is important to archive the original transmission as well as the
121 unpacking algorithm so that traceability to and from the primary-source evidence may be
122 achieved.

123 The output of some electronic sensors is recorded in raw data files that require
124 specialized unpacking, conversion and calibration algorithms to generate temperatures. In
125 the purest sense, the Level 0 data are machine-readable files of, for example, the voltages,
126 digital counts, refractivities, etc. In such cases, both the fundamental data stream and the
127 associated conversion algorithm would be considered together as Level 0, primary-source
128 information.

129 As with the instrumental data records recorded by hand on paper, these electronic
130 measurements and algorithms will require metadata documentation that defines the
131 process. While these activities appear especially onerous to climate researchers, the
132 reliability, reproducibility, and traceability requirements insist that such burdens be
133 accommodated wherever possible. We recognize that in many electronic systems, these
134 primary-source data and algorithms may be impossible to recover. In such circumstances
135 where fundamental source data are unavailable, the project should provide in the
136 associated Level 0 metadata archive of information, i.e. technical manuals, to describe the
137 instrumentation and conversion techniques which generated the archived geophysical
138 value at Level 1.

139 Because of the differing methods needed to discover and archive primary and secondary
140 sources based on their original form (i.e. documents vs. digital files) and the time frame
141 covered by these time series, it is anticipated that multiple, parallel (or perhaps
142 sequential) efforts will be required by teams of experts. In other words, if funding is
143 limited, the project may begin with data records deemed most vulnerable to loss, e.g. pre-
144 1950 paper documents in developing countries.

145 • *As an outcome of the workshop, there should be a clear definition of primary*
146 *(Level 0) and secondary (Level 1) source database across the spectrum of observing*
147 *systems which may contribute data to the land surface temperature database.*

148 • *We should establish a coordinated international search and rescue of Level 0,*
149 *primary-source climate data and metadata both documentary and electronic (see wp3.)*
150 *This effort would recognize and support similar on-going national projects. Once*
151 *located, the project should (a) provide, if necessary, a secure storage facility for these*
152 *documents or hard-copies of same, (b) create, where appropriate, digital images of the*
153 *documents for the archive for traceability and authenticity requirements, (c) key*
154 *documentary information into digital files (native format in Level 1 and uniform format*
155 *in Level 2), (d) archive, test and quality-assure raw data files, technical manuals and*
156 *conversion algorithms which are necessary to understand how the geophysical variable*

157 *may be unpacked and generated from electronic instrumentation, and (e) securely*
158 *archive the files for public access and use.*

159 • *A certification panel will be selected to rate the authenticity of source material as*
160 *to its relation to the “primary-source”, i.e. to certify a level of confidence that the Level 1*
161 *data, as archived, represents the original values from the Level 0 primary source. The*
162 *process will often be dynamic, since we anticipate that new information will always*
163 *become available to confirm or cast doubt on the current authenticity rating.*

164 Version Control

165 Real-world experience indicates that archives necessarily become dynamic libraries that
166 must adapt to the unforeseen appearances of new, competing or corrected data sets. Time
167 series which proceed through the Levels to the integrated databank (Level 3) must be
168 carefully indexed for traceability. Indeed, given the wide variety of surface temperature
169 observational methods and institutions that have been involved at some point in their
170 collection and use, this project understands that version control is a vital and complex
171 requirement.

172 Archives of historical observations can and should be updated as new information
173 becomes available. This, at the outset, presents a very difficult problem for the project at
174 hand, namely, at what scale (spatial and temporal) should version control be applied
175 when new or revised data are discovered and accepted as authentic? Should the time
176 series of each individual station be assigned a version number which may be up-
177 versioned when new information is gathered? Should temporal blocks be considered for
178 versioning (i.e. decade by decade or pre-1950 vs. post-1950, etc.)? Should the block of
179 documents pertaining to metadata for a station be up-versioned when a new document is
180 discovered? Should time series of stations be bundled into regional, national, continental
181 or global extents for version designations and updated only occasionally when a
182 significant number of changes are warranted? In cases with unformatted, raw digital data
183 files, what should be done when a change is warranted in applying the unpacking and
184 conversion algorithm of the original file?

185 A practical guideline to follow is that new versions should be adopted when there is a
186 significant addition or change to be instituted. This change would be considered
187 significant if products generated from the archive may be altered in some noticeable way
188 with the new information. In the meantime, new data may be placed in a pre-
189 authentication database, with appropriate caveats and indexing, for usage until an
190 authenticity investigation is completed and the data deemed suitable for the permanent
191 archive.

192 Flexibility will of course be required as versions of the most recent decade or two will be
193 subjected to a greater frequency of added observations, and thus a versioning system will
194 need to accommodate newer versus older portions of the archive. In addition, the
195 constant inclusion of current observations as time goes on will be accommodated.

196 The project will also supply to the community datasets for testing algorithms for accuracy
197 in the various aspects of spatial and temporal homogenization. These specialized datasets
198 will also require versioning so that for publications (or litigation) an investigator will
199 have a clear pathway defined to replicate the findings. It is possible that version control
200 of electronic data and software can be handled through commercial off-the-shelf systems
201 or open source systems such as Subversion (<http://subversion.apache.org/>).

- 202 • *Given the extent of this project and the unpredictable nature of the evolution of*
203 *the archive, the reliance on an active panel to address version-control issues as they*
204 *arise will be necessary. The panel will investigate the possibility of utilizing commercial*
205 *off-the-shelf or open-source version control software for electronic files and software*
206 *code (e.g. Subversion (<http://subversion.apache.org/>)).*

- 207 • *Since one requirement of this project is to preserve older versions of the archive,*
208 *and that a considerable amount of tedious research will be performed on any one*
209 *version, it is generally assumed that up-versioning will be performed of the basic, Level 2*
210 *digital archive as sparingly as possible.*

- 211 • *The algorithms that produce the datasets used for testing and the datasets*
212 *themselves must be documented and version-controlled.*

213 Configuration Management

214 The data system being proposed will be a dynamic system that satisfies a number of
215 requirements in order to establish and maintain a consistent set of products. An example
216 of a system designed to store and provide access to document images of the original
217 paper records is EDADS operated by NOAA/NCDC.
218 [<http://www.ncdc.noaa.gov/oa/climate/cdmp/edads.html>]. Configuration management
219 includes the designation of format selection of accessible files for convenient public use.
220 These are the Level 2 products that provide a uniform format structure for ready analysis
221 and when integrated with other stations become the Level 3 databank. In most cases a
222 fundamental station record format that includes the core information (station identifier,
223 location, date, time (i.e. 0900) or category (i.e. TMax), temperature value, version, etc.)
224 will be the most commonly accessed files and provide the greatest utility for the users.
225 Certain relatively vital pieces of information may also be contained in this fundamental
226 data record such as type of instrument, type of shelter, and length of record with
227 consistent parameters (i.e. location, instrument, etc.) Associated with the Level 2 data
228 will also be pathways back to the primary source data (Levels 0 and 1) and to the
229 available metadata documents which explain the observation and establish traceability.

230 Though most surface stations are assigned a World Meteorological Organization
231 identifying number, there are many for which this has not been done. For these stations
232 the CM team will work with the WMO to assign new WMO-qualified identifiers.
233 However, if it is deemed necessary, the team may construct a new system of
234 identification that is highly expandable and contains within it an apperency feature which

235 allows an individual to easily recognize the region wherein the station resides. This
236 could include the use of a FIPS country code and network identifier.

237 Many builders of large datasets have encountered the existence of duplicate datasets of a
238 particular station whose values differ or data records whose values seem at odds with the
239 climate of the identified station. If the Level 0 primary source evidence is available, the
240 time series can be reconstituted and authenticated. If the primary source document or
241 datafile has been lost, a pathway to a decision is required. In a parallel to “textual
242 criticism” in literature, it is expected that investigators, in many cases external to this
243 project, will need to sift through the information when multiple copies of a station record
244 are discovered, determine which files or documents are most closely related to the
245 missing primary source, and express (and document) an expert judgment as to the
246 decision which is ultimately made. It is likely that in many cases, it will be necessary to
247 archive multiple versions of a station record at Level 1. In some cases there may be
248 greater confidence in some records than others. Confidence ratings and discussions may
249 help in documenting the unique character of each record. This process will likely also
250 occur in situations where a primary source paper form (Level 0) of daily observations
251 may disagree with an official report of a monthly summary supposedly derived from the
252 primary source. The CM functions will obviously be closely linked with the version
253 control functions. Defining other requirements associated with hardware, security,
254 support, and financial resources will be necessary as part of this project.

255 • *A configuration management board will be selected to initially define the*
256 *necessary infrastructure, formats and other aspects of archive practices. A permanent*
257 *board will then be selected to oversee the operation. This board and the version-control*
258 *panel may be coincident or at least overlapping in membership.*

259 Summary

260 To provide primary-source surface temperature data for the research community through
261 a suitable and convenient archive, and to meet new requirements now being demanded of
262 climate observations (i.e. “admissible evidence” in the legal sense) a significant and on-
263 going investment will be required. To the extent resources are limited, the utility of this
264 project will also be limited in meeting the various requirements placed on data by the
265 differing communities who now utilize climate observations. Now is the time to initiate
266 this project (and support similar on-going projects) as institutional memory, documents
267 and critical material are being lost with each passing day. To whatever extent this project
268 is successful with data provenance, version control and configuration management, a
269 century from now, our descendants will either thank us or criticize us.

270

271

Appendix A

272 Definitions

273 **Data Provenance (DP)** refers to the confirmation or gathering of evidence
274 as to the time, place, and -- when appropriate -- the person responsible for
275 the creation, production, or discovery of the data.

276 **Version Control (VC,** also known as revision control, source control, or
277 software configuration management) is the management of changes to
278 documents, programs, and other information stored as computer files.

279 **Configuration Management (CM)** is a field of management that focuses
280 on establishing and maintaining consistency of a system's or product's
281 performance and its functional and physical attributes with its
282 requirements, design, and operational information throughout its life.

283

284