

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 **Creating surface temperature datasets to meet 21st Century** 33 **challenges**

34 35 **- 14: Solicitation of input from the community at large including** 36 **non-climate fields and discussion of web presence**

37
38 Authors:

39 Nigel Fox	NPL
40 Ian Joliffe	Univ of Exeter
41 Nick Barnes	Ravenbrook
42 Amy Luers	Google inc
43 Rob Allan	UK met office
44 Philip Brohan	UK met office
45 Michael De Podesta	NPL
46 Richard Chandler	UCL

47
48 Scope of document/task:

- 49
50 • Mechanisms to facilitate communication with all stakeholders including non-
51 climate science community.
 - 52 ○ Out going - Awareness
 - 53 ○ Incoming – input/views particularly from those with relevant expertise
 - 54 ○ Leading to “buy-in”
- 55 • Use of internet tools and web presence
 - 56 ○ Test-bed for new methodologies: blogs, social networking sites, data
 - 57 visualisation ...
 - 58 ○ Provision of data, tools and products that can be easily utilised by all
 - 59 stakeholders
- 60 • Maximising productive input and debate, not fuelling controversy and
61 confusion.
 - 62 ○ Providing adequate background / education
 - 63 ○ Controlled but flexibly moderated debate
 - 64 ○ Responsive feedback
 - 65 ○ Access to data
- 66 • Need to identify desired stakeholder input

67
68 Introduction

69
70 Today there is an unprecedented demand for climate information by all sectors of society. Local,
71 national and international organizations grappling with development and security issues are
72 increasingly seeking information about how the climate is changing and what the impacts of these
73 changes are likely to be, to inform how to respond to these risks. Adequately meeting society’s
74 growing need for trusted climate information will require a more open, transparent and
75 participatory approach to climate change science than we have seen in the past— one that engages
76 scientists, governments and the public.

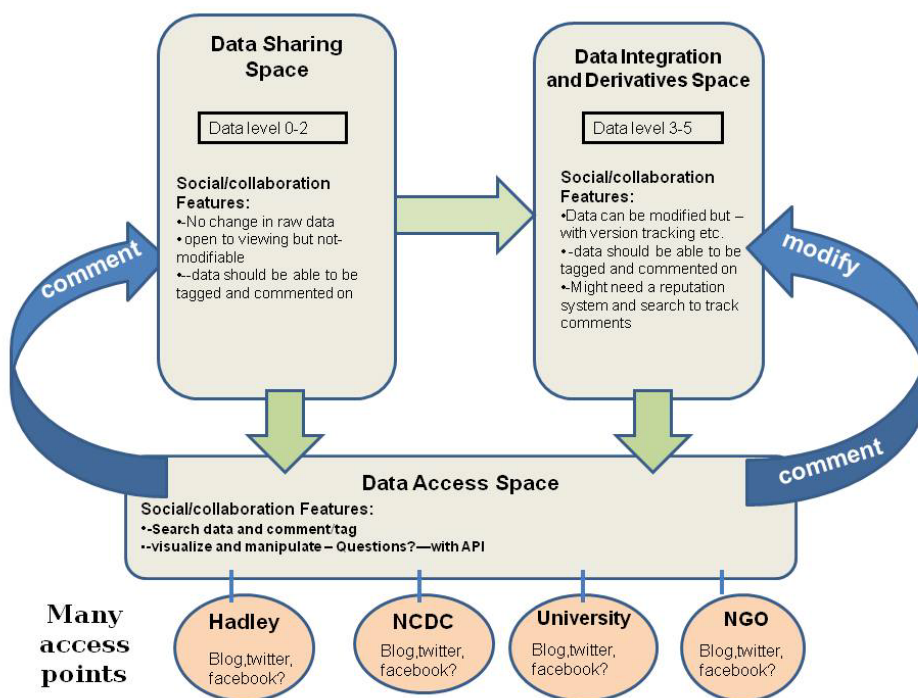
77
78
79 In this age of widespread internet access, information is available and shared almost
80 instantaneously on so many topics—through websites, facebook and twitter among other tools –
81 yet, the scientific community has not yet adapted. If science is to effectively support society’s

82 needs, the research and publication process must find its position in this new transparent
83 information culture.

84
85 This global surface temperatures project is a critical high profile opportunity to define a
86 new model of open, transparent and collaborative science and establish a precedent for
87 this new approach.

88
89 The challenge – is how to do this in a manner that maintains the rigor of the scientific
90 process? In this paper we explore how to integrate ICT tools to enable openness,
91 transparency and collaboration in the three key components of the global temperature
92 project: 1) Data Sharing; 2) Data integration and derivative products; 3) Data access and
93 useability.

94
95
96
97



98 Figure1 highlights how we might integrate social and collaboration tools in each of these
99 three areas. See White Paper 3 Figure 1 for elucidation of the data levels.

100

101 In considering the above scope it is essential that we take note of recent media and public
102 interest/sensitivity to the topic of climate change and temperature data in particular. We
103 need to be sure that we make every effort to provide full transparency of the activities in
104 the project and the methodologies being used or associated with it. However, to minimise
105 risk of further confusion from this process we will need to be sure that there is adequate
106 public awareness of basic scientific principles, statistics and uncertainty statements, the
107 latter both in terms of data and prediction.

108

109 The nature of this project also provides the opportunity, by making this a highly visible
110 activity, to rebuild public confidence in climate science (and science as a whole) and
111 facilitate “buy-in” to the need for a lower carbon future. Of course such visibility brings
112 with it obvious risks.

113

114 The global nature of this initiative together with its high profile can take the advantage to
115 experiment with new media/communication approaches which may also attract resources
116 in their own right (e.g. open source dataset creation algorithms, novel applications and
117 visualisation tools). The use of such tools particularly to visualise data or in
118 communication, may also help to engage with the younger generation.
119 One of the most important issues that needs to be addressed early on is “what type of
120 active input is being sought”?. Is it – ideas/agreement on methodologies for data
121 collection, analysis etc, or something else. For the purpose of this document it is assumed
122 that from a technical perspective the answer is largely as indicated i.e. data collection,
123 analysis and visualisation. However, there is also the presumption that we seek input or at
124 least feedback that the public as a whole understand the issues and importance of the
125 activity and its consequences.

126
127 The other key aspect is the need for openness and transparency. Without the entire process
128 being as open and transparent as feasible with exceptions to such an ethos clearly flagged
129 and justified (e.g. employer asserts IPR to some aspect of a given effort) it will be hard to
130 justify this project meeting these aims.

131
132 Stakeholder engagement for what and how?

133
134 There are three distinct groups that need to be considered when seeking an answer to this
135 question:

136
137 *i/ those involved in the community already or with technical expertise from another*
138 *discipline or community e.g. statisticians, data processing etc but also others engaged in*
139 *activities involving multiple-sampling, thermometry, metrology and probably many*
140 *others. Input would be solicited on all aspects of the process - methodologies being*
141 *proposed for data analysis, including any data screening criteria and visualisation as well*
142 *as data collection and its quality assurance.*

143
144 *If all data were accessible then this would also allow any user to apply independent tools*
145 *and analysis and report back to the project for subsequent discussion.*

146
147 *ii/ those likely to provide little technical input but who will be an active user of derivative*
148 *products from these core data sets to make decisions of environmental, societal and*
149 *economic importance.*

150
151 *iii/ those likely to provide little technical input e.g. General public, policy makers, but*
152 *whom we need to be sure understand and have full buy-in to what is being done*

153
154 Technically proficient/aware
155 In the first case, “technically proficient/aware”, it will probably be adequate to make use
156 of normal established communication methods, e.g. scientific journals, open
157 meetings/workshops and a website. The latter would of course require some level of
158 interactability and feedback probably with the ability for wiki/blog discussion to seeded
159 questions and or unsolicited input. We should explore new collaborative dataspace models
160 to support tagging and commenting whilst recognising that not all scientists will be
161 comfortable with such technologies (e.g. the blog for discussion has seen minimal use
162 despite being flagged to numerous expert communities through multiple listservs).
163 However, as a starting point a few questions phrased to allow simple YES/NO responses

164 wherever possible, but with options for free flow comment may be simpler and less
165 intimidating. For this to work effectively, in terms of data useage, it will be important to
166 give clear examples and clarity on any analysis approach or methodology that is being
167 proposed for discussion with real or simulated data to scope that expected in practise.
168 Because it is important that this project not be seen to be owned by a single entity but
169 rather be a truly international multi-partner effort we should actively consider a cloud
170 based model taht is not centered on one website but can be accessed from several sites.
171 This may yield an additional organisational overhead to achieve in total but would ensure
172 that the overhead did not fall on a single institution which may be more sensible. For all
173 practical purposes the website would become the principle focal point for interaction.
174

175 It would need to have sufficient resources to ensure it was maintained and also responsive
176 to input including feedback to correspondents. Here, again, a distributed model may work
177 better whereby the comments relevant to a given component of the project are managed
178 by those with relevant expertise. This will of course also require some level of
179 appropriate and consistent moderation. Any answers to questions would need to be
180 automatically compiled into a summary table and any free-flow answers/comments copied
181 into a blog-like format to allow further community discussion.
182

183 It would be appropriate to have some form of registration before any inputs are
184 submittable to the site, however full visibility should be available without registration.
185 Registration should not normally preclude anyone, but some reasonable filtering system
186 may need to be employed to reduce poor quality, non-attributable, input e.g.
187 correspondents should be happy to provide real contact details and use regular email
188 addresses, and of course in many countries this would probably be subject to the data
189 protection act. More important is that any input will have to meet certain norms: be on
190 topic; non-inflammatory; avoid accusations of mal-practice etc. if we are to ensure the
191 project retains focus.
192

193 Whilst dedicated workshops and meetings are likely to bring the most rapid and
194 comprehensive interaction they are very costly in terms of organisation and therefore
195 probably best arranged as a session as part of other already organised events/conferences
196 which a critical mass of project participants are already planning to attend rather than as
197 stand alone events. However, after some pre-defined timescale or if some significant
198 controversy were to emerge then there may be justification to have an open event to ratify
199 any decision. There may also be a need for some meetings of governance structures
200 (which should be covered under the governance white paper and its discussions).
201

202 The scientific media, particularly those of a more general nature and wide circulation such
203 as “new scientist” but also Nature should be used as “pointers” to the web-site. General
204 awareness news articles rather than costly adverts being the preferred approach using
205 appropriate press-release mechanisms to support this. If written correctly these would
206 likely trigger broad media coverage and so the web presence would need to be ready for
207 this and also appropriate speakers available to answer questions.
208

209 In terms of data access the most transparent but also challenging approach that could be
210 followed is to provide full access to all data sets and support information for anyone to
211 analyse in their own way. Whilst this is attractive it does bring with it some significant
212 risks and issues not least the ability to obtain a uniform data access policy from all data
213 suppliers and a uniform data format to permit easy intercomparisons. There may also be

214 issues regarding misunderstandings of aspects of data quality which could lead to
215 confusing analysis being carried out and/or the need for time consuming discussions with
216 individual data providers to clarify issues. However, this would be the ideal scenario. The
217 application of any restrictions would be difficult. It is relatively easy to justify a limit on
218 access to data providers or funders to the project but as soon as this is extended to try to
219 include any one else, and not everyone it is hard to have what would be seen as fair
220 criteria.

221

222 **Non-technical specialist**

223 Whilst it is unlikely that there will be detailed comment and input on specific
224 methodologies being used or data quality from the non-technical specialists we should not
225 ignore their importance. This project is likely to be of wide interest and provides an
226 opportunity to increase awareness in climate change issues and further the public
227 understanding of science in general. It is thus important to give them the same access to
228 the web presence and opportunity to comment. Again, it would be beneficial to have a
229 cloud-based presence with multiple “front ends” to the portal and redundant serving
230 capacity. There may be benefit to having a presence aimed at different levels for the
231 technical expert, the non-technical end-user and the lay person (and possibly others). One
232 such entry point should be set-up with an educational objective and have sufficient
233 resources attached to take the technical data and distill this in an unbiased way into
234 information that is usable to a lay-audience. This could, for example, provide
235 information and background on the key scientific practices that are likely to cause
236 confusion, e.g. uncertainty and the need/methods for data screening/weighting etc. The
237 use of good graphics and ideally interactive visualisation tools to support this would be
238 expected by the public and organisations like Google are probably well placed to develop
239 an appropriate learning and visualisation environment. Many related organisations;
240 NOAA, NASA, ESA etc have similar educational sections of their websites and
241 encourage the same for individual projects as well.

242

243 It may be possible to propose small experimental projects to encourage engagement with
244 schools and help them gain an understanding of some of the sampling issues, e.g.
245 establishing their own small local networks of temperature monitoring stations and
246 observing variations due to location, sampling strategies etc. They could also be provided
247 with access to typical or real data for their own analysis, (sub-setted and stored
248 independent of the primary data sets) and perhaps comparing results and approaches with
249 other schools, again as a learning and awareness activity. Such educational activities are
250 likely to be relatively low cost but could have major benefit in terms of public awareness
251 and buy-in. These could also include efforts to digitise local data to augment the databank
252 (see White Paper 3).

253

254 **Data and metadata accessibility**

255 In considering data accessibility, it is important to consider two things,

- 256 • Ease of access to the user
- 257 • Responsibility for data quality and access

258

259 *Databank* (see also white papers 3-6)

260

261 In considering data accessibility from a user perspective, particularly if the users are to be
262 global and of varying levels of expertise, it is important that all data from whatever source
263 is provided in a similar (ideally identical) format and with consistent and transparent

264 metadata and quality descriptors. This will ensure users can make reasonable decisions on
265 data selection and useage. Since such datasets are expected to be used for climate studies
266 it is important that they are stored in an invariant form for all future generations. Any
267 modification of the data (even if for good reason) should be done and described as part of
268 any processing activity to avoid potential confusion through use of different source data
269 and version controlled. From a simplistic perspective a single (and of course mirrored)
270 archive of data, accessible from a single portal would appear to be the easiest solution as
271 the cost of data storage is relatively low. If the archive was truly an archive, where data
272 was only added or accessed and not subject to any modification, then this may be the best
273 solution and would ensure that any and all users would always have access to the same
274 source data. Once submitted it would be fixed for all time. If a distributed data system
275 were established, this would place severe requirements on data suppliers to freeze and
276 maintain links, and guarantee not to update any data sets. There are opportunities to
277 engage the community in the creation of this databank through crowd-sourcing of
278 digitisation with suitable quality control (double keying from unique IPs with quality
279 checks by full-time experts).

280
281 In both cases there is an assumption that data suppliers will allow free and open access to
282 their data and take responsibility for ensuring it is correct at the time of submission. If
283 any subsequent error is found or an update required the data-base will need to be able to
284 tag the data set with some indicator that a change is required and that this should be stated
285 in any future use of the data. It will of course be critical to have a common set of
286 guidelines for data format, and metadata with an appropriate QA process to be
287 demonstrated by data suppliers going forwards. For existing data, an appropriate
288 screening and subsequent quality indicators should be created for each data set to ensure
289 consistent analyses.

290
291 It seems reasonable that any user of the data archive should be required to register and
292 agree to reasonable acknowledgements in any published useage. Data should be
293 searchable and downloadable for off-line processing, ideally making use of graphical
294 interfaces. However, there should also be on-line tools to visualise and manipulate data
295 available. The nature of the data is not highly complex and so such tools should be
296 readily available from many commercial suppliers who might be encouraged to volunteer
297 to provide them as part of a promotional activity. This can be solicited in the early stages
298 of the projects development and facilitated in the seeded questions section of the web
299 front end.

300
301 *Datasets (see also white paper 13)*

302
303 Datasets should be made available through the cloud-based web presence along with their
304 attendant metadata upon acceptance for publication. The metadata should include all
305 intermediate processing steps, a description of the algorithm and its code. As has been
306 proven by e.g. www.clearclimatecode.org there is substantial value to release of the code
307 enabling a greater understanding of the processing by involvement of interested members
308 of the public on a pro bono basis as well as scientific peers. However, a distinction needs
309 to be made between the scientifically useful reproduction activity undertaken by e.g.
310 clearclimatecode.org and the scientifically meaningless replication process of running the
311 same code on my PC as you ran on yours which only proves I can replicate your mistakes
312 on my PC. The scientific value in providing the process metadata accrues through
313 reproduction and assessment of sensitivity to methodologically uncertain choices.

314 Datasets and metadata, including their quality assessment against the benchmarks should
315 be made available in a common format to aid users. In the first instance it would make
316 sense for the dataset, its metadata and associated web presence to be hosted by the dataset
317 creators or a group nominated by them. They will be best placed to answer questions
318 related to their particular product.

319

320 *Data products*

321

322 Once multiple groups have produced datasets for a given station / region / timescale then
323 tools will be required that enable this data to be collated along with attendant metadata
324 pertaining to dataset quality and additional flags. Then additional capabilities will be
325 required to aid in interpretation. This may require the user to input what is of most interest
326 to them: characterisation of extremes; seasonal variability; trend characterisation etc. as
327 each dataset will have different capabilities and strengths and weaknesses. A suite of
328 graphical capabilities will be required as well as a capability to tabulate the data in the
329 way the user desires.

330

331 New web based communications

332 If we are to maximise impact and exposure of the project we should take advantage of all
333 media/communication opportunities. This should include, at least for promotional
334 purposes, tools such as Twitter, Facebook etc and also virtual worlds such as Second Life
335 for those aspects deemed appropriate whilst recognising that not all project aspects will
336 necessarily be suited to these tools. All of the above are being used by many
337 organisations to promote awareness but also in some cases to solicit opinions and views.
338 Whilst it is arguable that the use of these tools will not provide any new information or
339 key stakeholders they do offer different opportunities to promote the project and its
340 outputs through a secondary marketing approach, i.e. there will be communication
341 opportunities from those interested primarily in the tool being used and where our project
342 may be of secondary interest.

343

344 Organisations like Google might be able to help promote the project and provide
345 significant support in visualisation. For example through adding additional layers or
346 markers on for example “Google Earth”, identifying the location of all sample sites,
347 together with the associated high resolution satellite imagery this could help identify
348 sampling errors due to environmental effects as well as historical images for the location.
349 These pointers could also be linked to the source data records and with some relatively
350 simply additional tools could allow time series visualisation for any site and subsequently
351 global or grouped visualisation of any data set in a simple intuitive manner.