

# Has poor station quality biased U.S. temperature estimates?

Ronan Connolly <sup>\*1</sup>, Michael Connolly <sup>1</sup>

<sup>1</sup> Connolly Scientific Research Group. Dublin, Ireland.

## Abstract

Two independent surveys have found that about 70% of the thermometer stations in the U.S. Historical Climatology Network (USHCN) dataset are currently poorly or badly sited. Previous investigations into how this poor siting has affected estimates of U.S. temperature trends have led to apparently contradictory conclusions. However, in this study, these contradictions are resolved, and it is shown that poor station quality has introduced a noticeable warming bias into temperature trend estimates for the U.S.

For the unadjusted station records, this poor siting increased the mean temperature trends by about 32%. When time-of-observation adjustments were applied to the records, this increased temperature trends by about 39%, and so the relative fraction of the trends due to the siting bias decreased. However, the siting biases were still substantial, and increased trends by about 18%.

The step-change homogenization algorithm which had been developed to remove non-climatic biases such as siting biases was shown to be seriously problematic. Instead of correcting the poorly- and badly-sited station records to match the trends of the well-sited stations, it appears to have blended the temperature records of all stations to match the trends of the poorly-sited stations.

It seems likely that similar poor siting biases also exist in global thermometer datasets, and this has probably led to an overestimation of the amount of “global warming” since the 19th century.

### Citation:

R. Connolly, and M. Connolly (2014). *Has poor station quality biased U.S. temperature estimates?*, Open Peer Rev. J., 11 (*Clim. Sci.*), ver. 0.1 (non peer reviewed draft). URL: <http://oprj.net/articles/climate-science/11>

**Version:** 0.1 (non peer-reviewed)

**First submitted:** January 8, 2014.

**This version submitted:** January 31, 2014.

This work is licensed under a [Creative Commons Attribution-ShareAlike 4.0 International License](https://creativecommons.org/licenses/by-sa/4.0/).



wind[3, 4] and thermometers kept above asphalt concrete can report temperatures considerably warmer than over soil or grass[5, 6]. As a result, the temperature records of poorly exposed stations are likely to contain non-climatic biases from localised changes in the micro-climate immediately surrounding the thermometer housing.

These “siting biases” or “inadequate station exposure biases” are different from the more widely-studied “urbanization biases” which we discuss in Refs. [11–13]. There are some similarities between both types of bias, e.g., they can both arise as a result of modernization and/or urban development in the area. However, while urbanization bias can physically alter the local climate of a large area, siting biases are strictly confined to the localized micro-climate in the immediate vicinity of the thermometer station. Because of this, the two biases can occur independently of each other, e.g., an urbanized station with a strong urbanization bias may have a very good station exposure, while a rural station with no urbanization bias may have a strong siting bias due to inadequate station exposure. In this study, we will

## 1 Introduction

Two recent surveys[1, 2] of the local surroundings of the thermometer stations in the [U.S. Historical Climatology Network](https://www.ncep.noaa.gov/historical/) (often abbreviated USHCN), have revealed that about 70% of the stations are currently sited in poorly or badly exposed locations. It is well-known that the local environment within a few hundred metres of a thermometer station can lead to unusual “micro-climates”, which are unrepresentative of the climate in the surrounding area[3–10]. For instance, nearby trees can reduce sunlight and

\*Corresponding author: [ronanconnolly@yahoo.ie](mailto:ronanconnolly@yahoo.ie). Website: <http://globalwarmingsolved.com>

35 focus on siting biases.

36 The 1218 stations used for the U.S. Historical Cli- 86  
37 matology Network were selected according to their 87  
38 spatial coverage, record length, data completeness, 88  
39 and historical stability. For this reason, the U.S. 89  
40 Historical Climatology Network is the main dataset 90  
41 used for calculating monthly or annual temperature 91  
42 trends for the contiguous United States, i.e., all of 92  
43 the United States except Hawai'i and Alaska. It is 93  
44 compiled and maintained by the National Climatic 94  
45 Data Center of the U.S.-based National Oceanic and 95  
46 Atmospheric Administration (NOAA) - referred to 96  
47 henceforth as the National Climatic Data Center. 97

48 The U.S. Historical Climatology Network is also in- 91  
49 cluded in the National Climatic Data Center's **Global** 92  
50 **Historical Climatology Network** (often abbreviated 93  
51 GHCN), which as we discuss elsewhere[13] is the main 94  
52 weather record dataset used for estimating *global* 95  
53 temperature trends. The stations in the U.S. network 96  
54 account for nearly 17% of the stations in the global 97  
55 network. But, more importantly, they account for 98  
56 the vast majority of the stations in the global net- 99  
57 work which are both rural *and* have relatively long, 100  
58 complete station records. For instance, 219 of the 101  
59 225 stations (i.e., 97.3%) in the Global Historical 102  
60 Climatology Network that are identified as rural in 103  
61 terms of both night-light brightness and associated 104  
62 populations, and have data for at least 95 of the last 105  
63 100 years are in the U.S. network. These long, rural 106  
64 records are the ones least likely to be affected by 107  
65 urbanization bias - a systematic bias which we argue 108  
66 in Refs. [11–13] has introduced an artificial warming 109  
67 trend into weather station-based global temperature 110  
68 trend estimates. 111

69 For these reasons, the U.S. network is an important 112  
70 part of any analysis of global temperature trends. 113  
71 Hence, problems in the reliability of the U.S. His- 114  
72 torical Climatology Network have implications not 115  
73 just for regional U.S. temperature trend estimates, 116  
74 but also for global temperature trend estimates. In 117  
75 this study we attempt to estimate the net sign and 118  
76 magnitude of the non-climatic biases introduced into 119  
77 the U.S. temperature records by inadequate station 120  
78 exposure. 121

79 Our analysis is based on the results of the Surface 122  
80 Stations project carried out by Watts et al. [1], and 123  
81 suggests that inadequate station exposures have in- 124  
82 deed introduced a noticeable warming bias into U.S. 125  
83 temperature trends. It is probable that similar bi- 126  
84 ases exist in global temperature trend estimates. Sev- 127  
85 eral studies have previously attempted to construct 128

estimates of the bias from the Surface Stations re- 86  
sults, but have each reached different conclusions[2, 87  
14–17]. We will attempt to rationalise the apparent 88  
contradictions between the different analyses (includ- 89  
ing ours). 90

The layout of this article is as follows. In Section 91  
2, we will summarise the results of the Surface Sta- 92  
tions project, and review the current literature on 93  
how station quality can influence the station temper- 94  
ature trends. In Section 3, we will present our analy- 95  
sis of the poor station quality problem using the Sur- 96  
face Stations results. In Section 4, we will compare 97  
our analysis with the previous studies of the Surface 98  
Stations results, and discuss the reasons for the differ- 99  
ent conclusions between these studies. In Section 5, 100  
we will make some practical suggestions and recom- 101  
mendations that we believe could substantially im- 102  
prove the quality of the available temperature record 103  
datasets. Finally, we will offer some concluding re- 104  
marks in Section 6. 105

## 2 Literature review 106

### 2.1 Motivation for the Surface 107 Stations project 108

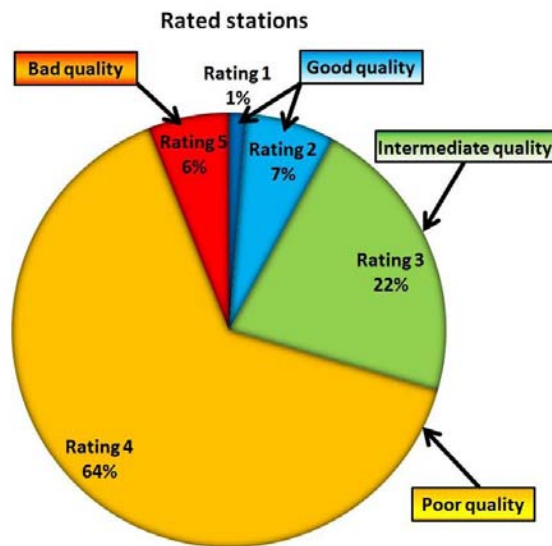
The stations in the U.S. Historical Climatology 109  
Network were taken from a larger network called 110  
the **NOAA Cooperative Observer Program Network**, 111  
(henceforth, the “COOP Network”), which is a 112  
volunteer-run weather observation program for the 113  
U.S. Since the stations are mostly volunteer-run, 114  
sometimes the official recommendations provided to 115  
the observers by NOAA National Weather Service on 116  
how to maintain the station are overlooked. This has 117  
led some researchers to speculate that the exposure of 118  
some of the thermometer shelters may be inadequate. 119  
For instance, Robinson, 1990[18] noticed that some 120  
weather observers had dramatically altered the expo- 121  
sure of their thermometer shelter when they switched 122  
to using electronic thermometers. He cautioned that 123  
this may have biased the temperature records. 124

In 2002, Davey & Pielke, 2005[8] carried out on- 125  
site inspections of 57 COOP stations (including 10 126  
Historical Climatology Network stations) in eastern 127  
Colorado. They found that many of the stations 128  
were poorly exposed. Some stations were located be- 129  
side air conditioner exhausts, some were surrounded 130  
by trees and/or buildings and some were set up 131  
over a gravel surface instead of grass. Vose et al., 132  
2005[19] and Peterson, 2006[20] argued that the Na- 133

134 tional Climatic Data Center's homogenization adjust- 160  
135 ments which were applied to some of the releases of 161  
136 the U.S. Historical Climatology Network had already 162  
137 removed (or at least reduced) any such biases. But, 163  
138 Pielke et al., 2007b[10] questioned the reliability of 164  
139 the adjustments. 165

140 Vose et al., 2005[19] had also claimed that Davey & 166  
141 Pielke's study was over too small a part of the U.S. 167  
142 to assume it was a widespread problem. However, 168  
143 Mahmood et al., 2006 showed that inadequate station 169  
144 siting was also a systematic problem for stations 170  
145 in Kentucky[6]. So, in 2007, Watts decided to extend 171  
146 Davey & Pielke's study. Together with a group of 172  
147 more than 650 volunteers, he began visual and photo- 173  
148 graphic on-site inspections of all of the U.S. Historical 174  
149 Climatology Network stations, with his [Surface 175  
150 Stations](#) project. In Section 2.2, we will summarise 176  
151 their findings and in the following sections we will 177  
152 discuss the implications of these findings.

## 153 2.2 Summary of the Surface Stations 177 154 findings 178



155 **Figure 1:** The relative percentage of different expo- 179  
156 sure ratings in the U.S. Historical Climatology Network 180  
157 stations that have been rated by the [Surface Stations 181  
158 project](#). 182

155 The [Surface Stations](#) project used the rating 177  
156 scheme followed by the National Climatic Data 178  
157 Center when they were establishing a high quality 179  
158 weather station network called the [United States Cli- 180  
159 mate Reference Network \(USCRN\)](#) in 2002. This 181

scheme was described in a NOAA technical docu- 160  
ment[21] and was based on a scheme proposed by 161  
Leroy, 1999[22]. 162

The main parameter by which the stations were 163  
classified was the distance of the thermometer sensor 164  
from artificial heating sources (or reflecting surfaces), 165  
such as buildings, concrete surfaces or parking lots. 166

**Rating 1** There are no artificial heating sources 167  
within at least 100 m of the sensor. 168

**Rating 2** There are artificial heating sources within 169  
30 to 100 m of the sensor. 170

**Rating 3** There are artificial heating sources within 171  
10 to 30 m of the sensor. 172

**Rating 4** There are artificial heating sources less 173  
than 10 m from the sensor. 174

**Rating 5** There are artificial heating sources located 175  
next to, or below, the sensor. 176

The higher the rating number, the more likely it is 177  
that temperature readings have been biased by un- 178  
representative micro-climate conditions. Ratings 1 179  
and 2 have excellent or good site exposure, and so 180  
are unlikely to be biased by micro-climate conditions, 181  
while measurements at Rating 5 sites are likely to be 182  
dominated by micro-climate biases. If a station is 183  
strongly influenced by micro-climate changes, then 184  
this reduces the reliability of its temperature records 185  
for considering climatological trends. Hence, for this 186  
analysis we define stations with ratings of either 1 or 187  
2 to be "good quality", those with a rating of 3 as "in- 188  
termediate quality", those with a rating of 4 as "poor 189  
quality" and those with a rating of 5 as "bad quality". 190

Leroy, 1999 also recommended that all stations 191  
should be located on flat and horizontal ground, sur- 192  
rounded by a clear surface with a slope of less than 193  
19°. He recommended that the stations should also 194  
be far from large bodies of water, unless they are 195  
representative of the area, in which case the station 196  
should still be located at least 100 m away[22]. How- 197  
ever, he did not indicate how to modify the ratings 198  
if stations did not meet those requirements, so these 199  
factors do not appear to have been included in the 200  
[Surface Stations](#) rating system. 201

Additional requirements for a station to be clas- 202  
sified with Rating 1 were that the surrounding 203  
grass/low vegetation ground cover is less than 10 cm 204  
high, and that the sensor stops being shaded once the 205  
sun reaches an elevation of 3° or lower. If either of 206  
these requirements are broken, then the best rating 207

208 the station can receive is Rating 2. However, in this  
209 analysis, since only about 1% of the stations have  
210 Rating 1, we group both Ratings 1 and 2 together,  
211 and so do not consider the vegetation and shading  
212 factors.

213 The Surface Stations group have archived ratings  
214 for 1007 out of the 1218 stations (82.7%) which we  
215 downloaded on 8th March 2012 from their website  
216 at <http://www.surfacestations.org>. We show the  
217 relative percentages of these ratings in Figure 1, and  
218 their locations in Figure 2. Examples of stations from  
219 each of the four subsets are shown in Figure 3. It can  
220 be seen that only about 8% of the stations have a  
221 good quality rating (1 or 2), and the majority of the  
222 rated stations have poor or bad quality ratings of 4  
223 or 5 (70%).

224 For a few of the unrated stations, the Surface Sta-  
225 tions team were able to determine that the rating  
226 changed substantially in recent years, and so, they  
227 were not assigned a rating. For instance, Malad City,  
228 ID (105559) would currently have a rating of 4, but  
229 only since 2008, while the Houma, LA (164407) sta-  
230 tion would have had a rating of 3 during the pe-  
231 riod 2004-2007. This confirms that station qualities  
232 change over time, and so it should be recognised that  
233 the current Surface Stations ratings are only based  
234 on the station quality they had at the time of the  
235 surveys. Obviously, it would be preferable if simi-  
236 lar assessments were also available for early periods.  
237 Nonetheless, we will see that it is still possible to  
238 make useful retrospective assessments of the poor sta-  
239 tion quality problem from just the current ratings.

240 In response to the Surface Stations findings, NOAA  
241 National Weather Service carried out independent re-  
242 assessments of 276 of the stations. Their reassess-  
243 ments confirmed that the Surface Stations ratings  
244 were reasonably accurate, and that poor siting is in-  
245 deed a systemic problem in the U.S. Historical Cli-  
246 matology Network[2].

247 Watts, 2009 provided photographs of many of the  
248 stations and it can be seen that, in many cases, sta-  
249 tion thermometers were situated near (or over) as-  
250 phalt roads or parking lots, and beside buildings,  
251 sometimes beside the exhaust fans of air condition-  
252 ing units, amongst many other problems[1]. This ex-  
253 plains why so many of the stations received a bad  
254 rating. All of these problems could easily have bi-  
255 ased the station records, and so it is important to  
256 reliably account for any such biases.

## 2.3 Previous assessments of the Surface Stations findings

257  
258  
259 When the Surface Stations project had rated most  
260 of the stations, Watts published a photographic re-  
261 port illustrating the remarkably high occurrence of  
262 poor quality siting in the U.S. Historical Climatol-  
263 ogy Network[1]. He discussed how it was plausible  
264 that this poor quality siting had introduced artificial  
265 warming trends into many of the station records, and  
266 that this would have introduced warming biases into  
267 the regional U.S. temperature trends which had been  
268 calculated from the Historical Climatology Network,  
269 e.g., Ref. [23]. He also suggested that similar prob-  
270 lems could exist for the rest of the Global Historical  
271 Climatology Network. He did not attempt to quan-  
272 tify this proposed bias, however.

273 At the time of writing, at least five studies (aside  
274 from our own) have attempted to quantify the extent  
275 of this bias, by using the Surface Stations results[2,  
276 14–17]. However, before we discuss these studies, it  
277 is important to briefly consider the different temper-  
278 ature datasets available for the stations.

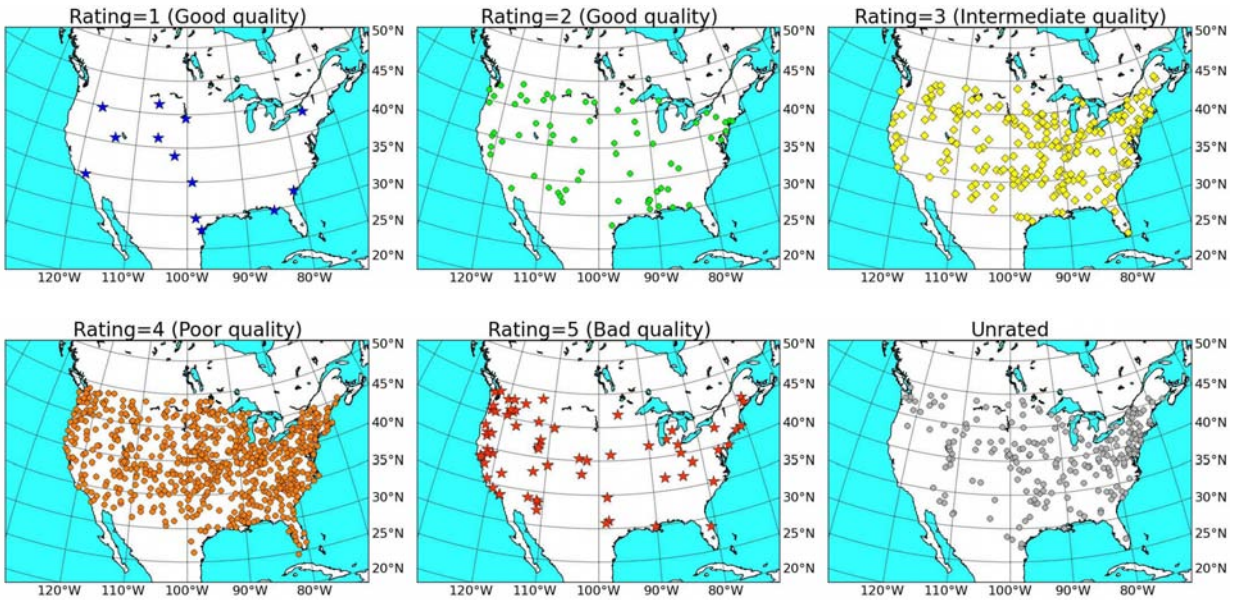
279 The National Climatic Data Center provide three  
280 different releases of their U.S. Historical Climatol-  
281 ogy Network datasets[24], which differ in the degree  
282 to which they have been adjusted for potential non-  
283 climatic biases.

284 One of their releases is essentially unadjusted, but  
285 has undergone a set of quality control checks to re-  
286 move individual monthly values which appear erro-  
287 neous, e.g., monthly temperatures that are far greater  
288 than (or far less than) the seasonal average for that  
289 station. We will refer to this release as the “*Unad-*  
290 *justed*” dataset.

291 In a second release, they have also applied specific  
292 adjustments to the temperature records for each sta-  
293 tion to account for documented changes in the time of  
294 day that the observers made their measurements[25,  
295 26]. We will refer to this release as the “*Time of*  
296 *observation adjusted*” dataset.

297 For their third release, they also carry out a series  
298 of station-station inter-comparisons to identify and  
299 remove station-specific non-climatic step-changes[27].  
300 We will refer to this release as the “*Time of observa-*  
301 *tion and step-change adjusted*” dataset. The National  
302 Climatic Data Center also apply inter-station inter-  
303 polation to “fill in” any missing gaps in the station  
304 records for the third release.

305 Although there is now a general acceptance that  
306 the findings of the Surface Stations project are ac-  
307 curate, and that the majority of the stations in the



**Figure 2:** Location of stations in the U.S. Historical Climatology Network with different siting quality ratings, as well as the remaining unrated stations.

308 U.S. Historical Climatology Network are of a poor or  
 309 bad quality, there has been considerable debate over  
 310 exactly what biases, if any, this poor siting has actu-  
 311 ally introduced into the reported U.S. temperature  
 312 trends.

313 Menne et al., 2010 suggested that the National Cli-  
 314 matic Data Center's homogenization procedure[24]  
 315 had already accounted for any biases which poor sit-  
 316 ing may have introduced. Moreover, they suggested  
 317 that if there was any residual bias, it was probably  
 318 a slight *cooling* bias[2], rather than the warming bias  
 319 Watts, 2009 had suggested.

320 Muller et al., 2013[15] claimed that the linear  
 321 trends of the *Unadjusted* records for stations with  
 322 Ratings 1, 2 or 3 were comparable to those of stations  
 323 with Ratings 4 or 5, and that there was not much  
 324 difference between estimates constructed from the  
 325 Ratings 1-3 and Ratings 4-5 subsets of the USHCN.  
 326 Therefore, they concluded that poor siting did not  
 327 have much effect on the temperature trends.

328 Martinez et al., 2012 used the Surface Stations rat-  
 329 ings in their analysis of temperature trends for the  
 330 state of Florida (USA)[16]. For their study, they  
 331 used the *Time-of-observation and step-change ad-*  
 332 *justed* dataset. As they were only studying the trends  
 333 for Florida, their study only involved 22 Historical  
 334 Climatology Network stations. So, they were cau-

335 tious about drawing definitive conclusions on the ef-  
 336 fects of poor station quality on temperature trends.  
 337 Nonetheless, they found the linear trends were dif-  
 338 ferent for the subsets of the worst rated (4 & 5) and  
 339 best rated (1 & 2) over the two periods they con-  
 340 sidered - 1895-2009 and 1970-2009. From this they  
 341 concluded that station quality *does* influence temper-  
 342 ature trends. However, they were unclear as to the  
 343 sign of this influence. For the 1895-2009 period, the  
 344 poor quality stations showed a greater warming trend  
 345 in the mean monthly temperatures than the good  
 346 quality stations, while for the 1970-2009 period, the  
 347 reverse applied.

348 Fall et al., 2011[14] agreed with Menne et al., 2010  
 349 that the National Climatic Data Center's homoge-  
 350 nization adjustments reduced much (although not all)  
 351 of the difference between the good quality and poor  
 352 quality subsets. But, they argued that the *Unad-*  
 353 *justed* and *Time-of-observation adjusted* trends of the  
 354 poor quality stations showed a substantial warming  
 355 bias relative to the good quality stations. They also  
 356 argued that poor exposure led to biases in Diurnal  
 357 Temperature Range (DTR) trends.

358 The Diurnal Temperature Range is the difference  
 359 between the mean maximum daily temperatures and  
 360 the mean minimum daily temperatures. Although  
 361 it is not as widely studied as the mean temper-



(a) Fallon, Nevada (Good quality)

(b) Boulder, Colorado (Intermediate quality)



(c) Napa State Hospital, California (Poor quality)

(d) Santa Rosa, California (Bad quality)

**Figure 3:** Examples of stations from each of the four subsets - Good quality (Fallon; ID = 262780, Rating 2); Intermediate quality (Boulder; ID = 050848, Rating 3); Poor quality (Napa State Hospital; ID = 046175, Rating 4); Bad quality (Santa Rosa; ID = 047965, Rating 5). Photographs were downloaded from <http://surfacestations.org/>. The photographer for these four stations was Anthony Watts.

362 ature trends, there has been considerable interest  
 363 in Diurnal Temperature Range trends, partly be-  
 364 cause it is thought they can provide insight into the  
 365 cause of mean temperature trends, e.g., Refs. [28-  
 366 30]. Fall et al. calculated that the worst-sited sta-  
 367 tions implied a Diurnal Temperature Range trend of  
 368  $-0.4^{\circ}\text{C}/\text{century}$ , whereas the best-sited stations had  
 369 essentially no long-term trend.

370 Watts et al. (in preparation, 2012)[17] argued that  
 371 the original Leroy et al., 1999[21, 22] rating system  
 372 used by Watts, 2009 was not rigorous enough, and  
 373 they re-evaluated the station exposures using the rec-  
 374 ommendations of Leroy et al., 2010[31]. When they  
 375 applied this new rating system to the stations, they

found a greater difference between the poor quality  
 and good quality stations than before (for the *Unad-*  
*justed* records), with the poorly-sited stations show-  
 ing more warming. They questioned the reliability  
 of the National Climatic Data Center's homogeniza-  
 tion adjustments, and suggested that a combination  
 of poor station exposure, urbanization bias and unre-  
 liable homogenization adjustments had led to a spu-  
 rious doubling of U.S. mean temperature trends over  
 the period 1979-2008.

376  
 377  
 378  
 379  
 380  
 381  
 382  
 383  
 384  
 385

## 2.4 Exposure problems for the early instrumental period

Before we discuss our re-analysis of the Surface Stations findings, it is worth discussing a related station exposure problem which has also received some discussion in the literature lately. Namely, there have been several attempts to estimate the extent to which changes in the types of thermometer shelters used by weather observers have biased 18th, 19th and early 20th century temperature records, relative to modern records[7, 32–44].

During the 20th century, a lot of the thermometers used by weather observers were housed in outdoor shelters like the Stevenson screen<sup>1</sup>. In recent decades, many stations have switched to using automated thermometer systems. Although there may have been instrumental biases from the switch in instrumentation during the recent move to automation[45–47] and some of the station observers also seem to have reduced the quality of the station exposure[1, 2, 18], the new automated stations are also outdoor shelters. However, before the introduction of the Stevenson-type screens, thermometers had quite different station exposures. For instance, in the early 1700s, English weather observers were encouraged to keep their thermometers unshielded and *indoors* in well-ventilated, north-facing, fireless rooms[7].

Chenoweth, 1992[32] and 1993[33] noted that many U.S. observers in the late 19th century and early 20th century were using similar unshielded, north-facing, thermometers, and even when screens were introduced, there was a wide range of different types of screens. In many cases, the siting of the screen was often inappropriate, e.g., attached to a wall[32], and would have received a bad or poor rating under the Surface Stations project. In some cases, thermometers located at railway stations, frequently registered artificially high temperatures when trains arrived at the station, but also sometimes gave minimum temperatures that were too low if the arriving train shook the thermometer[32].

Parker, 1994[34] compiled some information on the thermometer exposures *recommended* by different countries between the mid-19th and mid-20th centuries, and found that the recommended thermometer exposures often varied dramatically over time. Unfortunately, most of the available information is rather limited, and establishing the *actual thermome-*

*ter exposures used* has been contentious. For example, there has been considerable debate over when the use of Stevenson screens first became widespread in Australia[35–39].

It is quite likely that these changes in thermometer shelters introduced biases of some sort (sometimes called “shelter biases”, “screen biases” or “early instrumental period biases”). However, establishing what those biases would have been is difficult. The biases introduced into daily averages by different thermometer exposures could depend on a number of factors, e.g., the observation time and averaging method used; the degree of urbanization of the site; the materials the nearby buildings were constructed from; the station siting of both the old and new exposures.

Some studies have attempted to estimate these biases by comparing temperature measurements recorded simultaneously at the same site by thermometers with different shelters[32–34, 40, 43, 44]. In some cases, these studies were carried out during the actual change-over, i.e., the late-19th/early 20th century[33, 34, 40]. However, in other cases, they are modern attempts to recreate the transition[32, 33, 40, 43, 44].

We note that these experiments are not as straightforward as commonly assumed. For instance, Chenoweth, 1992 calculated different estimates of the screen bias of a north-facing window thermometer relative to a thermometer in a Cotton Region Shelter (a similar shelter to the Stevenson screen), depending on which Cotton Region Shelter he used. For his estimates, he had two different shelters, one located in his backyard, and the other in a nearby field. From the photographs in Chenoweth, 1992, it appears to us that one of the shelters would have had a good quality Surface Stations rating, while the other would have had a poor quality rating[32]. This suggests that any screen biases would have depended on the siting of the old and new thermometers.

We also note that many of the stations with long records are currently more urbanized than they would have been in the 18th/19th centuries. Since urbanization bias as well as other land use changes are known to affect the diurnal temperature range at a station[29], it is possible that estimates of different screen biases carried out using modern field tests may over/underestimate the actual biases. For instance, we used Google Earth to analyse aerial photographs of the study sites used by Brunet et al., 2011[44] (see the Supplementary Information for the Google Earth station location files and some aerial photographs).

<sup>1</sup>Invented in the 19th century by Sir Thomas Stevenson, the father of the well-known author, Robert Louis Stevenson.

485 Both of the sites Brunet et al. used appear to  
486 be highly urbanized, and the local environment has  
487 clearly undergone dramatic changes since the 19th  
488 century. For this reason, it is likely that the cur-  
489 rent “screen bias” is different than it would have  
490 been in the 19th century/early 20th century when the  
491 Montsouri shelters (which Brunet et al., 2011 were as-  
492 sessed) were actually in use. We note that the bias  
493 appears to be greater at the more urbanized of the  
494 two stations (Murcia). From the aerial photographs,  
495 it also appears that both sites would receive a poor  
496 or bad Surface Stations rating, however, it is possible  
497 that this was also the case for the historic sites.

498 In contrast, the site used for the Böhm et al., 2010  
499 study[43] is located in a relatively rural location, at  
500 a Benedictine monastery in Kremsmünster, Austria.  
501 However, while it is a useful comparative study, their  
502 estimates of the screen bias are not necessarily repre-  
503 sentative.

504 Their old window thermometer is located in a tall  
505 astronomical tower with panoramic views of the sur-  
506 rounding town. Meanwhile, the new shelter is in a  
507 heavily shaded garden in front of the building, sur-  
508 rounded on all four sides by either tall buildings or  
509 trees.

510 Shelters exposed in the shade will tend to register  
511 lower maximum temperatures than properly exposed  
512 shelters[32]. So, if their new shelter was too heavily  
513 shaded, then this would have exaggerated Böhm et  
514 al.’s estimate of the screen bias. Analysis of the lo-  
515 cation of the new shelter suggests it does suffer from  
516 shading problems. For instance in the 8th August  
517 2012 aerial Google Earth photograph for the location  
518 (see Supplementary Information), the shelter was in  
519 shadow, and it was also in shadow in the photograph  
520 shown in Böhm et al., 2010, which was taken on 21st  
521 March 2007[43].

522 Another potential problem in estimating the shel-  
523 ter biases is that the observation times and averaging  
524 methods used by observers have also changed over  
525 time. As we will discuss in Section 4.4, different ob-  
526 servation times and averaging methods can lead to  
527 different estimates of the daily mean temperature at  
528 a station. So, the biases introduced by changes in  
529 the thermometer shelter used would have depended  
530 on which averaging method the observers were using.

531 In particular, we note that the apparent biases  
532 often seem to be greatest for the minimum and  
533 maximum daily temperatures, and could have been  
534 smaller for some of the observation times which might  
535 have been used. Many modern weather observers use

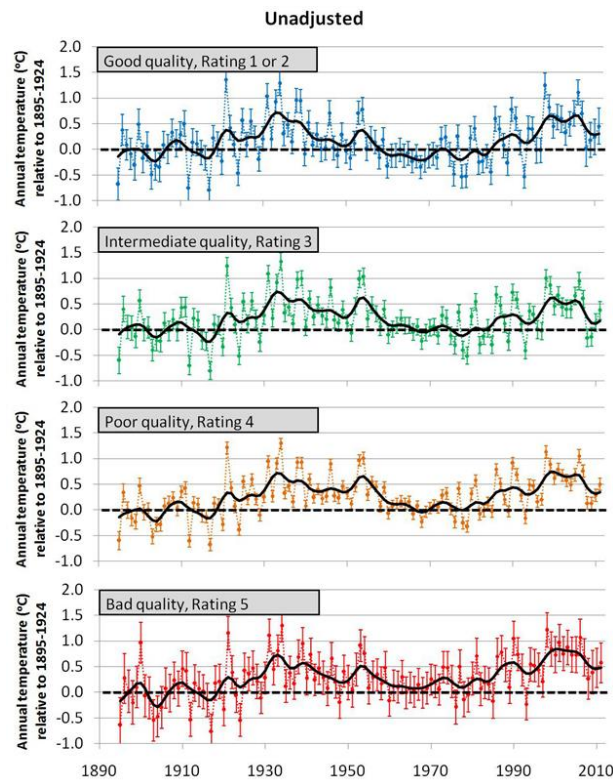
536 minimum-maximum thermometers and approximate  
537 the daily mean temperature by calculating the mean  
538 of the maximum and minimum temperatures reached  
539 during the previous 24 hours. But, especially in the  
540 18th and 19th centuries, many observers would have  
541 measured the temperatures at specific times in the  
542 day, and estimated the daily average using those mea-  
543 surements. When the different shelters were being de-  
544 veloped in the 19th century, there was considerable  
545 awareness of the biases that different thermometer  
546 exposures introduced[33, 34, 40]. So, it is plausible  
547 that, in some cases, the averaging formulae used were  
548 partially chosen in an attempt to minimise the screen  
549 biases.

550 For all these reasons, we find the current estimates  
551 of these screen biases are still incomplete, and require  
552 more careful studies. Despite this, several researchers  
553 believe that screen biases have led to a significant  
554 overestimation of pre-20th century summer tempera-  
555 tures[34, 37, 39–44]. This is considered a particular  
556 problem for the longest thermometer records, which  
557 are mostly European. One of the main reasons for  
558 this belief appears to be because temperature proxies  
559 for the same areas suggest colder 18th century tem-  
560 peratures than the thermometer records[41–43]. We  
561 review this “convergence problem” issue in a separate  
562 paper[48], and a detailed comparison of temperature  
563 proxies and thermometer records is beyond the scope  
564 of this article. However, we will note that it seems  
565 ironic that, while some researchers are arguing that  
566 early thermometer records are unreliable because the  
567 temperature proxies show colder temperatures, other  
568 researchers are arguing that temperature *proxies* are  
569 unreliable for the late 20th century, because they do  
570 not show the warm temperatures of the thermometer  
571 records[49, 50], i.e., the so-called “divergence prob-  
572 lem”.

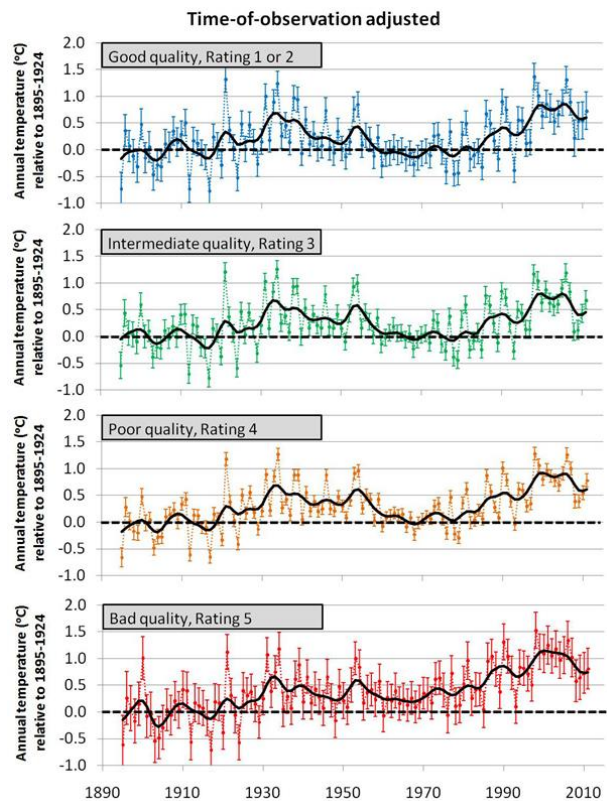
### 3 Our reanalysis 573

574 For our analysis of the siting biases in the U.S. His-  
575 torical Climatology Network, we first downloaded the  
576 station ratings from the Surface Stations website (on  
577 8th March 2012). Ratings were available for 1007  
578 out of the 1218 stations (82.7%). We then divided  
579 the stations with ratings into four subsets, i.e., good  
580 quality, intermediate quality, poor quality and bad  
581 quality, as described in Section 2.2. For each of these  
582 subsets we calculated the gridded mean temperature  
583 trends using the three different U.S. Historical Clima-  
584 tology Network datasets, i.e., the *Unadjusted, Time-*





**Figure 4:** Mean gridded trends for each of the four subsets using the Unadjusted dataset. Solid lines correspond to the 11 point binomial smoothed versions of the annual values. Confidence errors correspond to twice the standard error of the annual means.



**Figure 5:** Mean gridded trends for each of the four subsets using the Time-of-observation adjusted dataset. Solid lines correspond to the 11 point binomial smoothed versions of the annual values. Confidence errors correspond to twice the standard error of the annual means.

585 of-observation adjusted and Time-of-observation and  
586 step-change adjusted datasets.

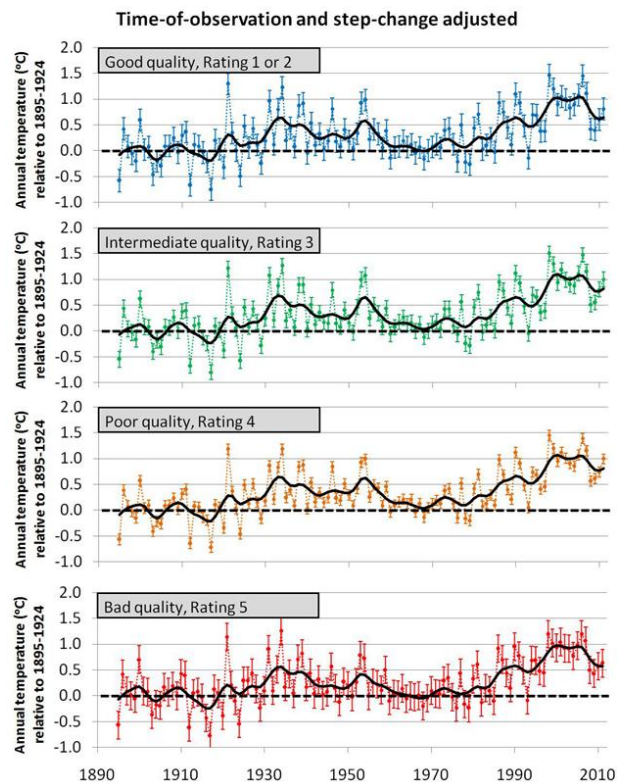
587 To calculate the gridded mean trends we adopted  
588 a similar procedure to the one we used in Refs. [11–  
589 13]. Namely, we first converted all of the station  
590 annual temperature records for a given subset into  
591 temperature anomalies relative to the mean station  
592 temperature during a 30 year baseline period, 1895–  
593 1924. As we will discuss in Section 4.1, this differs  
594 from the 1961–1990 baseline period we used in Refs.  
595 [11–13], because we wanted to better study the di-  
596 vergence between subsets over time. Stations which  
597 did not have at least 5 years of data during this pe-  
598 riod were dropped from our analysis. This also differs  
599 from the 15 year minimum we used in Refs. [11–13].  
600 As we will discuss in Section 4.1, this was because  
601 we wanted to reduce the number of stations dropped  
602 from our analysis. 87 of the 1218 stations ( $\sim 7\%$ )  
603 were dropped from our analysis for this reason.

604 Stations were then assigned to  $5^\circ \times 5^\circ$  grid boxes.  
605 For each year, the mean temperature anomalies for  
606 each of the grid boxes were calculated by determining  
607 the simple mean of the temperature anomalies of all  
608 the stations in that box with data for that year.

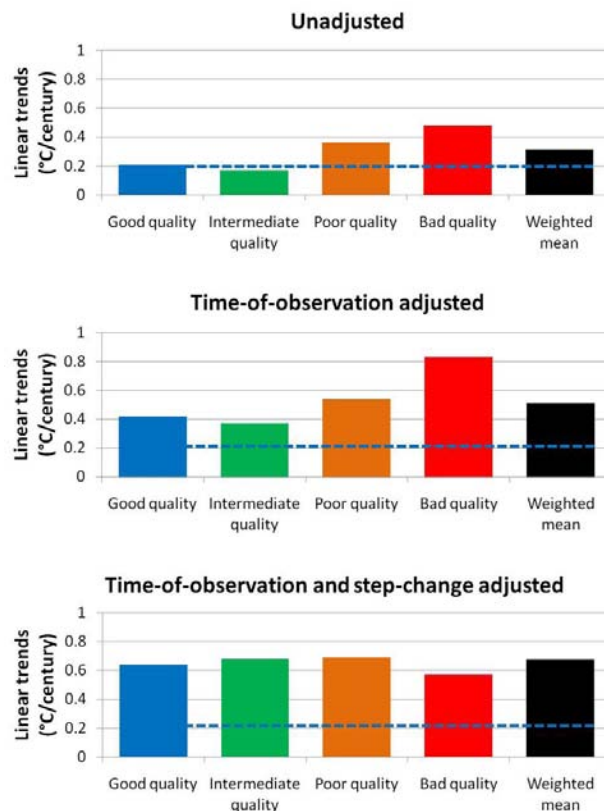
609 The mean U.S. temperature anomaly for each year  
610 was then calculated as the area-weighted mean of all  
611 of the grid box means. Standard errors of the mean  
612 were also calculated using the method described in  
613 the Supplementary Information.

614 The mean U.S. temperature trends for each of the  
615 four subsets, using each of the three datasets are  
616 shown in Figures 4, 5 and 6.

617 The different estimates of U.S. temperature trends  
618 all have a lot of similarities, but there are also sub-  
619 stantial differences between them. In terms of the  
620 similarities, one striking feature is the pronounced al-  
621 ternation between “warming” periods and “cooling”



**Figure 6:** Mean gridded trends for each of the four subsets using the Time-of-observation and step-change adjusted dataset. Solid lines correspond to the 11 point binomial smoothed versions of the annual values. Confidence errors correspond to twice the standard error of the annual means.



**Figure 7:** Bar charts showing the 1895-2011 linear trends of each of the subsets, and the weighted mean of the subsets, for each of the three datasets. The dashed blue lines correspond to the linear trend of the Unadjusted good quality subset.

622 periods, each lasting several decades. Since the start of  
 623 the estimates in 1895, there seem to have been two  
 624 warm periods (1920s-1930s and 1990s-2000s) and two  
 625 cool periods (1900s-1910s and 1960s-1970s).

626 One of the main differences between the estimates  
 627 is in how the warm periods and cool periods compare  
 628 to each other. For the good quality and intermediate  
 629 quality, *Unadjusted* subsets, the early warm period  
 630 appears comparable to the recent warm period, and  
 631 the early cool period appears comparable to the re-  
 632 cent cool period. However, as the quality of the sub-  
 633 sets decreases, or as the datasets become more heavily  
 634 adjusted, the relative warmth of the early warm pe-  
 635 riod appears to decrease, and the recent cool period  
 636 appears warmer.

637 It is worth noting that, if the less heavily-adjusted,  
 638 good quality subsets are reliable, then the recent  
 639 warming in the U.S. does not appear unusual in the

640 context of the overall record. This would contradict  
 641 the popular perception that, because of man-made  
 642 global warming, the recent warm period is the hottest  
 643 on record in the U.S.[51-53]. The 1930s coincided  
 644 with a period of drought in parts of the U.S. leading  
 645 to the economically disastrous “Dust Bowl era”[54],  
 646 so it is quite plausible that this period of drought  
 647 also corresponded to a warm period. However, there  
 648 are many non-climatic biases present in long-term  
 649 temperature records like the U.S. Historical Clima-  
 650 tology Network. So, **if** the adjustments applied by  
 651 the National Climatic Data Center have successfully  
 652 removed these biases, then this would make the ad-  
 653 justed datasets more reliable. For this reason, we  
 654 will consider the implications of all three datasets.  
 655 We will assess the reliability of the National Climatic  
 656 Data Center’s adjustments in Sections 4.4 and 4.5.

657 The temperature trends of the estimates show

658 strongly non-linear behaviour. So, as we will discuss  
659 in Section 4.2, calculating the linear trends for the  
660 estimates is a very crude method for describing the  
661 temperature trends. Nonetheless, it offers a simple  
662 metric which allows us to make rough comparisons  
663 between the different estimates.

664 Table 1 lists the 1895-2011 linear trends for each  
665 of the estimates, as calculated by linear least squares  
666 fitting. The linear trends are also shown graphically  
667 in Figure 7. The corresponding  $r^2$  fitting parameters  
668 are also shown in Table 1, and it can be seen that  
669 most of the linear trends are very poor fits. This is  
670 as expected, since the trends are quite non-linear. So,  
671 the linear trend values should be treated cautiously.

672 We first note that all of the subsets have a long-  
673 term “warming” trend, i.e., they all have positive linear  
674 trends for the 1895-2011 period. However, as  
675 we mentioned above, this is not surprising, or particularly  
676 informative. The 1890s-1910s were a relatively cool period  
677 and the 1980s-2000s were a relatively warm period, so it  
678 would be expected that all estimates should have a positive  
679 trend.

680 Second, the net effect of each of the adjustments  
681 applied by the National Climatic Data Center is to  
682 substantially increase the linear trends of the estimates.  
683 For example, the *Unadjusted* good quality subset has a  
684 linear trend of  $+0.21^\circ\text{C}/\text{century}$ , but when the time-of-  
685 observation adjustments are applied, this doubles to  
686  $+0.42^\circ\text{C}/\text{century}$ , and it increases by a similar amount  
687 after the step-change adjustments to  $+0.64^\circ\text{C}/\text{century}$ .  
688 Similar increases occur for all subsets. The only exception  
689 is that the step-change adjustments reduce the linear trend  
690 of the bad quality subset.  
691

692 The next thing we note is that there are indeed substantial  
693 differences between the linear trends of the different  
694 subsets. For the *Unadjusted* and *Time-of-observation adjusted*  
695 datasets, the good quality and intermediate quality subsets  
696 have the lowest linear trends and the trend consistently  
697 increases from the intermediate quality to the poor quality  
698 to the bad quality subsets. This suggests that the biases  
699 due to inadequate station exposure introduce non-climatic  
700 warming trends into the temperature records.  
701

702 The step-change adjustments dramatically reduce the  
703 differences between subsets. There are at least two schools  
704 of thought on why this occurs. Some researchers have  
705 argued that this is because the step-change adjustments  
706 have managed to remove the siting biases at each station  
707 [2]. However, a second explanation is that the reduction  
708 arises because the sit-

ing biases have been blended or “homogenized” together,  
rather than removed. In Section 4.5, we argue in favour  
of this second argument. Watts has also favoured this  
second argument in on-line commentary on his website  
[55].

In Table 1, we also list the weighted mean linear  
trend of the rated stations, which is calculated by  
weighting the trend of each subset by the percentage  
of stations in that subset. Although these weighted  
mean trends are less than for the poor and bad quality  
subsets, they are still greater than the linear trends  
for the good quality subset. This suggests that the mean  
U.S. temperature trends are indeed significantly biased  
by inadequate station exposure.

Surprisingly, the linear trends for the *Unadjusted*  
and *Time-of-observation adjusted* datasets are actually  
slightly smaller for the intermediate quality subset  
than for the good quality subset. This appears to  
contradict the expectation that the bias should  
continuously decrease in going from the worst quality  
subsets to the best quality subsets.

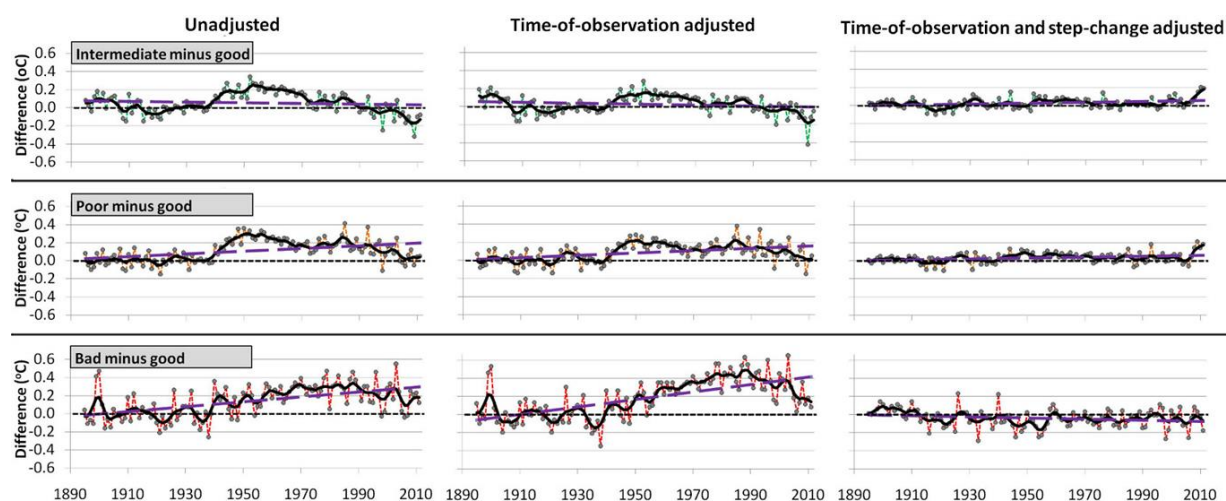
As Muller et al., 2013 suggest, it is plausible that  
the intermediate quality stations are of a high enough  
quality that they are unbiased [15]. After all, the only  
difference between Ratings 2 and 3 is the distance of  
the station from heating sources, and this distance is  
at least 10m for Rating 3 stations. So, it might be  
that 10m is a sufficient distance for a station to be  
unaffected by inadequate station exposure. In that  
case, since the sample size of the intermediate quality  
subset is considerably larger than the good quality  
subset (see Figure 1), it could be that the trends of  
the intermediate subset are the most reliable. However,  
this is not clear, and on physical grounds, we know  
that the good quality subset is the least likely of the  
subsets to have biases due to inadequate station  
exposure. Therefore we will assume that the good  
quality subset is the best representation of the  
“unbiased” subset, in terms of station exposure. This  
gives slightly lower estimates of the biases for the  
*Unadjusted* and *Time-of-observation adjusted* datasets.

On this basis, we estimate the siting biases for  
each subset by subtracting the good quality temperature  
trends of the appropriate dataset from the subset  
temperature trends. These difference trends are shown  
in Figure 8. The 1895-2011 linear trends and  $r^2$   
fitting parameters are listed in Table 2.

Again, the difference trends show substantial non-  
linearity. But, it can be seen from Figure 8, that  
the linear fits do at least offer a rough approximation  
of the 1895-2011 trend. On this basis, we can cal-

Subset	<i>Unadjusted</i>		<i>Time-of-observation adjusted</i>		<i>Time-of-observation and step-change adjusted</i>	
	$^{\circ}\text{C}/\text{century}$	$r^2$	$^{\circ}\text{C}/\text{century}$	$r^2$	$^{\circ}\text{C}/\text{century}$	$r^2$
Good (8%)	+0.21	0.03	+0.42	0.10	+0.64	0.23
Intermediate (22%)	+0.17	0.02	+0.37	0.09	+0.68	0.25
Poor (64%)	+0.36	0.10	+0.54	0.19	+0.69	0.28
Bad (6%)	+0.48	0.16	+0.83	0.37	+0.57	0.21
Weighted mean	+0.31		+0.51		+0.67	

**Table 1:** Linear trend estimates ( $^{\circ}\text{C}/\text{century}$ ), calculated by linear least squares fitting, of the four subsets and three different datasets shown in Figures 4, 5 and 6.  $r^2$  shows the fitting coefficients, which theoretically can vary from 0 (poor fit) to 1 (perfect fit)



**Figure 8:** Annual differences between the good quality subset and each of the other subsets for the three datasets. Solid lines correspond to the differences between the 11 point binomial smoothed versions. The coloured dashed lines correspond to the 1895–2011 linear trends of the differences.

760 calculate rough estimates for the siting bias of each of  
 761 the datasets from either Tables 1 or 2. For Table 1,  
 762 we can do so by subtracting the linear trend of each  
 763 subset from the linear trend for the good quality sub-  
 764 set. This gives the same values as in Table 2, since  
 765 the linear trend of the differences is equivalent to the  
 766 difference between the linear trends.

767 The biases in the *Unadjusted* dataset seem to be  
 768 about  $+0.15^{\circ}\text{C}/\text{century}$  for the poor quality sub-  
 769 set and  $+0.27^{\circ}\text{C}/\text{century}$  for the bad quality sub-  
 770 set (see Tables 1 or 2). For the set of all the  
 771 rated U.S. Historical Climatology Network stations,  
 772 i.e., the weighted mean, the siting bias seems to be  
 773 about  $+0.10^{\circ}\text{C}/\text{century}$ . As a percentage of the linear  
 774 trends, this means that siting biases account for  
 775 approximately 42% of the poor quality trends, 56%  
 776 of the bad quality trends and 32% for the entire rated

network.

777  
 778 The magnitude of the biases in the *Time-*  
 779 *of-observation adjusted* dataset are similar, i.e.,  
 780  $+0.13^{\circ}\text{C}/\text{century}$  for the poor quality subset,  
 781  $+0.41^{\circ}\text{C}/\text{century}$  for the bad quality subset and  
 782  $+0.09^{\circ}\text{C}/\text{century}$  for the entire rated network. How-  
 783 ever, since the time-of-observation adjustments con-  
 784 sistently increase the linear trends of all subsets by  
 785 about  $+0.2^{\circ}\text{C}/\text{century}$ , the fraction of the trends due  
 786 to siting bias is less. That is, the percentage of the  
 787 trends due to siting biases is reduced to 24% of the  
 788 poor quality subset, 49% of the bad quality subset  
 789 and 18% of the full rated network. Having said  
 790 that, the time-of-observation adjustments increase  
 791 the trends of the bad quality subset by more than  
 792 any of the other subsets ( $+0.35^{\circ}\text{C}/\text{century}$ ), which  
 793 increases the difference between the good and bad

Subset	Unadjusted			Time-of-observation adjusted			Time-of-observation and step-change adjusted		
	$^{\circ}C/100y$	$r^2$	Bias	$^{\circ}C/100y$	$r^2$	Bias	$^{\circ}C/100y$	$r^2$	Bias
Intermediate - good	-0.04	0.01	-24%	-0.05	0.03	-14%	+0.05	0.09	+7%
Poor - good	+0.15	0.16	+42%	+0.13	0.16	+24%	+0.05	0.09	+7%
Bad - good	+0.27	0.27	+56%	+0.41	0.40	+49%	-0.06	0.05	-11%

**Table 2:** Linear trends ( $^{\circ}C/century$ ) of the differences between the good quality subsets and the other subsets, calculated by least squares fitting from the data shown in Figure 8, for each of the datasets.  $r^2$  shows the fitting coefficients, which vary from 0 (poor fit) to 1 (perfect fit). “Bias” shows the ratio of the linear trend of the difference relative to the linear trend of the subset, as a percentage.

794 quality subsets, i.e., our estimate of the siting biases.  
795 So, this suggests that siting biases still account for  
796 about 49% of the trends of the bad quality subset.

797 The percentage of the trends in the *Time-of-*  
798 *observation adjusted* subsets that are due to the time-  
799 of-observation adjustments are as follows: Good qual-  
800 ity = 50% adjustments; intermediate quality = 54%  
801 adjustments; poor quality = 33% adjustments; bad  
802 quality = 42% adjustments; entire rated network  
803 = 39% adjustments. We will consider the time-of-  
804 observation adjustments in Section 4.4.

805 On average, the step-change adjustments also  
806 substantially increase the linear trends, by about  
807  $+0.16^{\circ}C/century$  (for the entire rated network).  
808 However, the step-change adjustments are quite dif-  
809 ferent for each of the subsets. In particular, the  
810 increase in trend is greatest for the good and in-  
811 termediate quality subsets ( $+0.22^{\circ}C/century$  and  
812  $+0.31^{\circ}C/century$ , respectively), and the adjustments  
813 actually decrease the trend for the bad quality sub-  
814 sets (by  $-0.26^{\circ}C/century$ ).

815 As a result, the differences between the subsets are  
816 substantially reduced, and the estimates of the siting  
817 biases for the *Time-of-observation and step-change*  
818 *adjusted* dataset are dramatically reduced. Indeed,  
819 since the trend of the bad quality subset is now less  
820 than the good quality subset, it could be argued, as  
821 Menne et al., 2010 did[2], that the bias becomes a  
822 slight “cooling” bias ( $-0.06^{\circ}C/century$ ) for the worst  
823 sited stations. The linear trend of the good quality  
824 subset is still less than for the entire rated network  
825 ( $+0.64^{\circ}C/century$  compared to  $+0.67^{\circ}C/century$  -  
826 see Table 1). So, nominally, it could be argued that  
827 there is still a bias of about  $+0.03^{\circ}C/century$ , but  
828 this is less than 5% of the linear trend of the entire  
829 rated network.

830 The problem in assessing the biases in the *Time-of-*  
831 *observation and step-change adjusted* dataset arises  
832 from uncertainties over the reliability of the step-

change adjustments. The step-change adjustments  
were developed in an attempt to identify and remove  
any non-climatic step-change biases in the tempera-  
ture records[27], e.g., station relocations or changes  
in instrumentation. So, a plausible explanation for  
the reduction in the differences between the subsets  
is that the adjustments have succeeded in removing  
the non-climatic siting biases, as Menne et al., 2010  
contend[2]. However, as Watts notes in on-line com-  
mentary[55], another plausible explanation is that the  
adjustments have blended together the biases from  
each of the subsets, leading to all subsets being bi-  
ased by about the same amount. In Section 4.5, we  
argue that the latter explanation is more likely.

In any case, it appears that siting biases have in-  
creased the temperature trends of the entire rated  
network by about 32% for the *Unadjusted* dataset  
and about 18% for the *Time-of-observation adjusted*  
dataset. This is a substantial non-climatic bias,  
which contradicts the claims of others, e.g., Muller  
et al., 2013[15]. The linear trends of the bad qual-  
ity subset are particularly heavily biased by siting  
bias (56% for the *Unadjusted* dataset and 49% for the  
*Time-of-observation adjusted* dataset), although it is  
true that the bad quality subset only makes up 6% of  
the U.S. Historical Climatology Network (Figure 1).

## 4 Discussion and comparison with the previous analyses

Initially, it might appear that there are serious con-  
tradictions between each of the attempts to use the  
Surface Stations results to estimate the biases intro-  
duced to U.S. temperature trends by inadequate sta-  
tion exposure. However, in this section, we will at-  
tempt to reconcile these apparent contradictions, and  
establish reasonable estimates for the actual siting bi-  
ases present in the U.S. Historical Climatology Net-

work. In Section 2.3, we discussed the findings of the Watts, 2009[1]; Menne et al., 2010[2]; Muller et al., 2013[15]; Martinez et al., 2012[16]; and Watts et al. (in preparation, 2012)[17], while in Section 3, we discussed the findings of our analysis.

All six of these studies have taken slightly different methodological approaches to analysing the Surface Stations results. So, in Section 4.1, we will first summarise these differences, and assess what influence they might have had on the conclusions of the different studies.

A major problem with analysing the Surface Stations results is that both the temperature trends of the individual subsets and the difference trends between the subsets are non-linear (as we discussed in Section 3). This means that using linear trend analysis is an overly simplistic approach. Nonetheless all six of the studies (including ours) use linear trend analysis. In Section 4.2, we will highlight some of the problems involved in using this approach.

Probably the most difficult challenge in estimating the siting biases lies in the fact that there are many other non-climatic biases present in the U.S. Historical Climatology Network. This means that it is difficult to establish how much of the differences between subsets are actually due to the siting biases, and how much are due to other factors. In Section 4.3, we will consider how one such factor (urbanization bias) could influence the siting bias estimates.

A popular approach to minimising this problem has been to rely on datasets that have been statistically “homogenized” in an attempt to reduce the magnitude of these non-climatic biases. However, this approach will only be successful if the homogenization techniques work. So, it is important to assess the reliability of these techniques. In Section 4.4, we will consider the time-of-observation adjustments, while in Section 4.5, we will assess the step-change adjustments.

#### 4.1 Use of different subsets and averaging approaches

In this section, we will consider some of the main methodological differences between each of the Surface Stations studies, and estimate what impact these differences have on the conclusions of the studies.

- Different approaches to regional averaging

Menne et al., 2010 calculated their regional trends for the contiguous U.S. by first assigning all of the

stations in each subset to  $0.25^\circ \times 0.25^\circ$  grid boxes, then calculating the mean monthly anomalies for each of those boxes. The grid anomalies were then area-weighted to yield regional anomalies for each of the subsets[2].

Instead of gridding the stations, Fall et al., 2011 divided the contiguous U.S. into nine climatic regions, then calculated the mean monthly anomalies of the different subsets for each of those regions. These regional means were then weighted by the areas of the regions, and averaged to give regional averages for the entire contiguous U.S.[14]. Watts et al. (in preparation, 2012) also adopted this approach[17].

Our approach is probably intermediate between those two approaches. Like Menne et al., 2010, we used an area-weighted uniform gridding. However, in keeping with our global analyses in other papers[11–13], we used a much larger grid size ( $5^\circ \times 5^\circ$ ). This larger grid size still yields about 40 grid boxes for the contiguous U.S. which is more than the 9 regions used by Fall et al. and Watts et al.

The Martinez et al., 2012 study just considered a single state (Florida) which only contained 22 stations. So, they did not attempt gridding and just calculated the monthly means of all the stations in each subset. The area of the state of Florida is less than one of our  $5^\circ \times 5^\circ$  grid boxes, so this is probably reasonable.

Fall et al., 2011 argue that the net differences between their approach and Menne et al.’s is very small when averaged over the contiguous U.S.[14]. So, it seems likely that the different gridding approaches used by us, Menne et al., Fall et al., Watts et al. and Martinez et al. are all reasonably equivalent, and therefore have probably not introduced much of a difference between the analyses.

Having said that, Muller et al., 2013 did not use a gridding approach, but instead used a particular statistical interpolation approach. In a separate study, we note that this approach to calculating regional (and global) temperature trends appears to dampen the differences between different station records[11]. We suspect that this is one reason why they were unable to detect much of a difference between their subsets, since other studies[14, 17] (including ours) were able to detect significant differences between the subsets in the *Unadjusted* datasets used by Muller et al.

- Different rating data sets

We downloaded the station ratings from the [Surface Stations](#) website on 8th March 2012. The rat-

968 ings available then were the ones used by Fall et al.,  
969 2011[14], and were determined from surveys carried  
970 out by the Surface Stations volunteers between 2nd  
971 June, 2007 and 23rd February, 2010. Ratings were  
972 available for 1007 out of the 1218 stations (82.7%).  
973 This was also the dataset used by Muller et al., 2013  
974 and Martinez et al., 2010[16]. However, the other  
975 studies used slightly different versions of the dataset.

976 Watts, 2009[1] was based on a preliminary analysis  
977 (2nd November, 2009) which only included ratings  
978 for 865 stations. However, although the Menne et  
979 al., 2010 study was carried out in response to the  
980 Watts, 2009 study, it was based on an even earlier  
981 (18th April, 2008) provisional dataset, in which only  
982 525 stations had been rated (43.1%), and these rat-  
983 ings had not undergone quality control. According  
984 to Watts, Menne et al. were in the process of ask-  
985 ing him for the more complete dataset used in the  
986 Watts, 2009 study, but stopped when Watts asked to  
987 be involved in their study[55].

988 For the Watts et al. (in preparation, 2012) stud-  
989 ies, further refinements and additional station surveys  
990 were available (15th June, 2011 - 1st July, 2012)[17].  
991 But, since the Watts et al. study is still in prepa-  
992 ration, the newer dataset had not been published on  
993 the Surface Stations website at the time of writing.

994 Although Menne et al., 2010 only used a provi-  
995 sional dataset, their results apparently broadly agree  
996 with those of the Fall et al., 2011 and Muller et al.,  
997 2013 studies[14, 15]. This suggests that, while their  
998 sample may have been somewhat biased by the “low-  
999 lying fruit problem”, as Watts warned[55], the prelim-  
1000 inary dataset was probably sufficiently complete for  
1001 an approximate estimate. For this reason, it seems  
1002 that the different sample sizes used by the studies did  
1003 not majorly influence their results.

1004 In addition, as a prelude to the Menne et al. study,  
1005 NOAA National Weather Service Forecast Office per-  
1006 sonnel investigated 276 of the 525 stations (52.6%)  
1007 and confirmed those ratings from the Surface Sta-  
1008 tions project[2]. Hence, the Menne et al., 2010 study  
1009 offers an independent confirmation of the poor qual-  
1010 ity of the station sitings.

1011 The Watts et al. (in preparation, 2012) study uses  
1012 the most up-to-date version of the dataset. However,  
1013 they used a new rating system described by Leroy,  
1014 2010[31], which is apparently more rigorous. As a re-  
1015 sult, they were only able to calculate ratings for 779  
1016 of the stations under the new system. This is still a  
1017 larger sample than the Menne et al., 2010 study, and  
1018 Watts et al. claim that the new system is able to more

1019 accurately distinguish between good and poor quality  
1020 stations[17]. The new system increases the percent-  
1021 age of stations identified as Rating 1 (6% compared to  
1022 1% under the old system); Rating 2 (14% compared to  
1023 7%); Rating 3 (32% compared to 22%) and Rating  
1024 5 (12% compared to 6%), but decreases the percent-  
1025 age of stations identified as Rating 4 (36% compared  
1026 to 64%).

1027 Since the ratings under the new system have not  
1028 been published yet, we have not been able to quan-  
1029 tify the effects such a change would have on our anal-  
1030 ysis. But, in general, if particular aspects of poor  
1031 siting lead to greater biases, then a rating system  
1032 which more rigorously distinguishes between these as-  
1033 pects should provide more accurate estimates of the  
1034 biases. We would expect the magnitude of our esti-  
1035 mated biases to increase under such a system, since,  
1036 under a less rigorous system, some stations which are  
1037 relatively unbiased would have been misidentified as  
1038 having a worse quality, while other stations which  
1039 are relatively biased would have been misidentified  
1040 as having a better quality.

1041 Watts et al., (in preparation, 2012) argue that the  
1042 new Leroy, 2010 rating system leads to an increase in  
1043 the differences between subsets[17], which suggests it  
1044 is an improved system. If so, it is encouraging that  
1045 the new system more than doubles the percentage of  
1046 stations identified as being of good quality, because it  
1047 suggests that it may be possible to isolate a relatively  
1048 large subset of high quality stations for estimating  
1049 U.S. temperature trends, after all.

- 1050 • Different sub-setting approaches

1051 As can be seen from Figure 1, the good and bad  
1052 quality subsets have relatively small sample sizes. In  
1053 particular, only 1% of the stations are of Rating 1,  
1054 and these stations are not evenly distributed across  
1055 the U.S. (see Figure 2). It is plausible that these small  
1056 sample sizes and uneven distributions could introduce  
1057 statistical artefacts. In an attempt to minimise these  
1058 potential artefacts, all of the studies have grouped at  
1059 least some of the ratings together.

1060 In the Menne et al., 2010 study, the small sample  
1061 size problems were accentuated, because they only  
1062 had ratings for 525 stations[2]. For that reason, they  
1063 decided to consider only two subsets - a “good” sub-  
1064 set comprising the 71 stations of their sample with  
1065 Ratings 1 or 2 (13.5%), and a “poor” subset compris-  
1066 ing the 454 stations with Ratings 3, 4 or 5 (86.5%).  
1067 We saw in Section 3 that while the Rating 5 sub-  
1068 set shows a large apparent warming bias (in the *Un-*  
1069 *adjusted* dataset at least), and the Rating 4 subset

1070 shows a considerable apparent warming bias, the linear  
1071 trends of the Rating 3 subset are actually the  
1072 lowest of all of the subsets, and that the apparent bi-  
1073 ases are “cooling”. So, by grouping Ratings 3, 4 and 5  
1074 together, the Menne et al. study would have reduced  
1075 the differences between their two subsets, making it  
1076 harder to detect the siting biases. In addition, their  
1077 “good” subset only included 71 stations, which were  
1078 not evenly distributed across the country. So, their  
1079 analysis might still have been affected by the problem  
1080 of inadequate sample sizes.

1081 Muller et al., 2013 argued from the fact that the  
1082 Rating 3 stations had the lowest linear trends, that  
1083 they were not affected by siting biases[15]. Hence,  
1084 for the main part of their analysis, they constructed  
1085 two subsets - an “OK” subset comprising stations  
1086 with Ratings 1, 2 and 3 and a “bad” subset com-  
1087 prising stations with Ratings 4 and 5. However, it  
1088 must be acknowledged that their preliminary analy-  
1089 sis involved a separate histogram of linear trends for  
1090 all five subsets, and they also included an additional  
1091 two subsets for their Table 1 (which listed the mean  
1092 linear trends of different subsets) - “good” and “poor”  
1093 using the same grouping as Menne et al.

1094 There was only two Rating 3 stations available for  
1095 the Martinez et al., 2012 study[16], and a total of  
1096 just 22 stations, so they created two subsets - one  
1097 from the 5 stations of Ratings 1 and 2 and the other  
1098 from the 13 stations of Ratings 4 and 5 (2 stations  
1099 were unrated). As we will discuss in Section 4.3, and  
1100 is apparent from Figure 1 and Table 1 of Martinez  
1101 et al., 2012[16], there are a number of climatic differ-  
1102 ences between the 22 stations in the Martinez et al.  
1103 study. So, it is probable that much of the apparent  
1104 differences between the subsets may be due to con-  
1105 founding factors other than siting bias. Martinez et  
1106 al., 2012 were conscious of this and cautioned against  
1107 relying on their comparisons too much[16].

1108 For our analysis, we wanted to avoid grouping too  
1109 many stations with different ratings together, but  
1110 considered the sample size of the Rating 1 stations  
1111 to be too small for a statistical analysis. Hence, we  
1112 grouped the Rating 1 and Rating 2 stations together  
1113 (as the “good quality” subset) and considered four  
1114 subsets, as described in Section 2.2. This means that  
1115 we were able to distinguish between the poor quality  
1116 and bad quality subsets, for instance. However, it  
1117 does mean it is probable that some of the apparent  
1118 differences between our subsets is related to the small  
1119 sample sizes and uneven distribution of the smaller  
1120 subsets, i.e., the good quality and bad quality sub-

sets, which only had 80 and 66 stations, respectively.

1121 Fall et al., 2011 also used our subsetting approach.  
1122 But, in addition, they attempted to assess the effects  
1123 of the uneven distribution of the good quality and  
1124 bad quality subsets, by carrying out a second analy-  
1125 sis testing the effects of using “proxies” for the good  
1126 and bad quality stations[14]. They constructed a re-  
1127 placement “proxy network” for both their good and  
1128 bad quality subsets by substituting each of the Rating  
1129 1 or 2 stations (for the good subset) and the Rating  
1130 5 stations (for the bad subset) with the nearest sta-  
1131 tions with a Rating of 3 or 4. They then compared  
1132 these proxy subsets with the full subset of all stations  
1133 with Ratings 3 or 4. They found that the 1979-2008  
1134 linear trends of the two proxy networks were differ-  
1135 ent from those of the full subset. This confirms that  
1136 at least some of the apparent differences between the  
1137 subsets is due to uneven distribution and small sam-  
1138 ple sizes[14].

1139 Watts et al. (in preparation, 2012) used a differ-  
1140 ent rating system, which divided the stations into  
1141 more evenly distributed subsets[17]. So, even though  
1142 the total number of rated stations was less than for  
1143 our analysis, the statistical problems of small sample  
1144 size and uneven distribution should have been con-  
1145 siderably reduced. At the time of writing, Watts et  
1146 al. had not yet archived these new ratings, but it is  
1147 likely that, when they are available, they would lead  
1148 to more reliable estimates of the siting biases.

- 1149 • Different homogenization steps of the tempera-  
1150 ture records

1151 As we mentioned in Section 2.3, the National Cli-  
1152 matic Data Center provide three different releases of  
1153 their U.S. Historical Climatology Network dataset.  
1154 The station records in two of their releases have been  
1155 adjusted in an attempt to remove the presence of non-  
1156 climatic biases[24]. We saw from Figures 4, 5 and 6,  
1157 as well as Table 1, that most of these adjustments  
1158 have the effect of increasing the long-term trends of  
1159 all of the subsets.

1160 If the adjustments are of similar magnitudes in all  
1161 subsets, then this should not overly affect the esti-  
1162 mates of the siting biases, since the differences be-  
1163 tween subsets would remain similar. However, we  
1164 saw in Section 3 that, while the time-of-observation  
1165 adjustments applied to each subset are quite simi-  
1166 lar, the step-change adjustments dramatically reduce  
1167 the differences between the subsets. As a result,  
1168 the estimates of the siting biases from the *Time-of-*  
1169 *observation and step-change adjusted* dataset will be  
1170



1171 much smaller than those determined from the other  
1172 datasets.

1173 The step-change adjustments were designed for re-  
1174 moving non-climatic biases[27]. So, it is plausible  
1175 that they may have removed (or reduced) the magni-  
1176 tude of the siting biases in the records. In this case,  
1177 the convergence of the station records would repre-  
1178 sent an improvement in their reliability, suggesting  
1179 the lower estimates are more realistic. However, if  
1180 the step-change adjustments are inadequate then this  
1181 convergence would actually represent a *reduction* in  
1182 their reliability. We will assess the reliability of the  
1183 step-change adjustments in Section 4.5.

1184 Menne et al., 2010 argue that the step-change ad-  
1185 justments could have removed many of the siting bi-  
1186 ases present in the unadjusted dataset[2]. They there-  
1187 fore claim that the *Time-of-observation and step-*  
1188 *change adjusted* dataset is the most reliable, and  
1189 hence their estimates of the siting biases are very  
1190 low. Also assuming that the adjustments are reli-  
1191 able, Martinez et al., 2012 only considered the *Time-*  
1192 *of-observation and step-change adjusted* dataset[16].

1193 In contrast, Muller et al., 2013 only considered the  
1194 *Unadjusted* dataset[15]. However, they carried out  
1195 their own independent data homogenization adjust-  
1196 ments. Like, Fall et al., 2011[14], we considered all  
1197 three releases in our analysis. We assess the validity  
1198 of the adjustments in Section 4.4 and 4.5.

1199 Watts et al. (in preparation, 2012) was scepti-  
1200 cal of the reliability of the adjustments, because  
1201 they noticed that the adjustments consistently in-  
1202 creased the 1979-2010 linear warming trends of *all*  
1203 subsets, even when they further sub-divided the sub-  
1204 sets by removing stations likely to be strongly af-  
1205 fected by non-climatic warming biases, i.e., urbanized  
1206 stations and airport stations[17]. As a result, they  
1207 considered the *Unadjusted* dataset to be more reli-  
1208 able than the *Time-of-observation and step-change*  
1209 *adjusted* dataset, and used the former for their main  
1210 conclusions. They did not analyse the *Time-of-*  
1211 *observation adjusted* dataset, as they considered such  
1212 an analysis beyond the scope of their study[17].

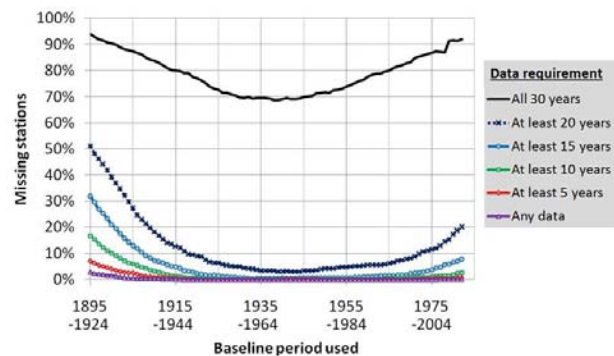
1213 Since the step-change adjustments reduce the dif-  
1214 ferences between subsets, in general, studies which  
1215 favour the *Time-of-observation and step-change ad-*  
1216 *justed* dataset should tend to obtain lower estimates  
1217 of the siting biases. However, we do note that Muller  
1218 et al., 2013 obtained a low estimate, even though they  
1219 used the *Unadjusted* dataset[15], while Martinez et  
1220 al., 2012 found evidence of strong siting biases, even  
1221 though they used the *Time-of-observation and step-*

*change adjusted* dataset[16].

- Use of different baseline anomaly periods.

1224 A common approach to constructing gridded tem-  
1225 perature trends is to convert each of the tempera-  
1226 ture records into “temperature anomalies” relative to  
1227 a fixed *baseline period*, e.g., 1961-1990. This means  
1228 that, for each station, the mean temperature of that  
1229 station over this period was subtracted from all of its  
1230 annual temperatures. This “common anomaly” ap-  
1231 proach[56] has the advantage that gridded averages  
1232 are not overly affected by changes in the numbers  
1233 of “cold” (e.g., high altitude) and “warm” (e.g., low  
1234 latitude) stations which have data for a given year.  
1235 Typically, a 30 year period is chosen, as this is gen-  
1236 erally considered long enough to reduce the noise of  
1237 inter-annual temperature changes, but short enough  
1238 that most records will have enough data in the pe-  
1239 riod.

1240 Unfortunately, the common anomaly approach has  
1241 the disadvantage that it artificially increases the ap-  
1242 parent agreement between stations near the anomaly  
1243 period (e.g., 1961-1990). This can be easily under-  
1244 stood by recognising that if the anomaly period is  
1245 shortened to one year, then the gridded average tem-  
1246 perature anomaly for that year will be exactly  $0.0^{\circ}\text{C}$ ,  
1247 by definition. A consequence of this is that the dif-  
1248 ferences between any two subsets may be artificially  
1249 reduced near the anomaly period.



1250 **Figure 9:** The percentage of U.S. Historical Clima-  
1251 tology Network stations which have to be dropped from  
1252 analysis when different 30 year baseline periods and data  
1253 requirements are used.

1250 To reduce the effect of this statistical artefact, we  
1251 chose the earliest possible 30 year period for our anal-  
1252 ysis, i.e., 1895-1924. This means that the difference  
1253 between the subsets should be low at the start of the

records, and any long-term divergences between the subsets should become more pronounced towards the end of the records.

A difficulty in using an early anomaly period is that many of the shorter records will not have enough data during that period to calculate an anomaly mean. From Figure 9 it can be seen that very few stations have complete records for the entire 30 year period, regardless of the period chosen<sup>2</sup>. However, if we are content to allow stations to be missing some data in the anomaly period, anomaly averages can be calculated for most of the stations during a lot of different baseline periods.

In our analyses in other studies, we typically use the 1961-1990 period and a requirement of at least 15 years data[11–13]. However, for this analysis, the requirement of at least 15 years, would have meant discarding more than 30% of the stations for the 1895-1924 period (Figure 9). As a result, we reduced this restriction to a minimum of 5 years during the 1895-1924 period. This reduced the number of stations which had to be discarded to 87, or about 7% of the available stations.

Each of the other studies used different anomaly periods. Menne et al., 2010[2] and Martinez et al., 2012[16] used the 1971-2000 period. Watts et al. (in preparation)[17] and Fall et al., 2011 used 1979-2008, but Fall et al. also considered a 36 year period, 1895-1930[14]. Muller et al., 2013 used a 31 year period, 1950-1980[15].

It is plausible that these different anomaly periods could alter the subset trends. To test this possibility, we repeated our analysis using the 1961-1990 period (which is the same one we use in our other papers[11–13]). The mean gridded temperature trends of each of the subsets using the two different baseline periods are shown in Figure 10. As there are a large number of plots shown, we have just shown the 11-point binomial smoothed trends, for clarity. The annual variability and standard errors of the mean for the non-smoothed plots are similar to those shown in Figures 4, 5 and 6.

Individually, the mean trends of each of the subsets are essentially the same under both baseline periods. The only major difference is that they have been rescaled to a different baseline. This is confirmed by Table 3, where we have calculated the 1895-2011 linear trends for each of the equivalent subsets. Al-

<sup>2</sup>The situation is better for the *Time-of-observation and step-change adjusted* dataset, since many of these data gaps have been filled in using inter-station and intra-station interpolation.

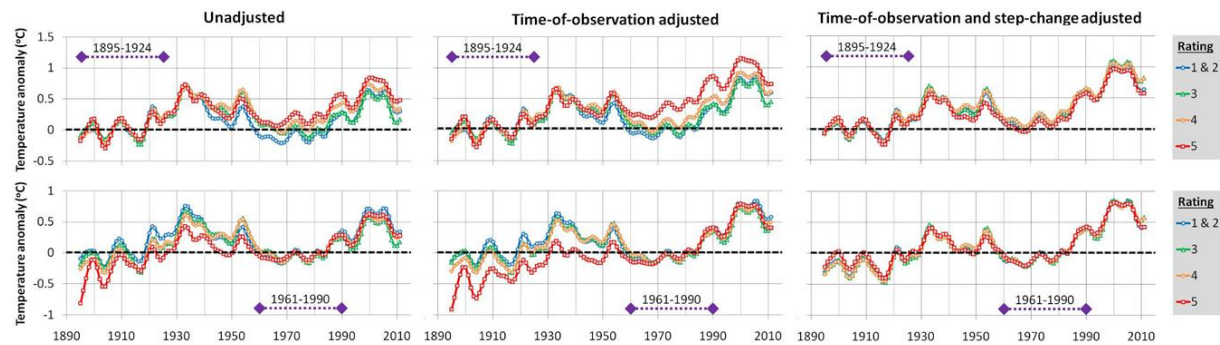
1895-2011 linear trends ( $^{\circ}C/century$ )		
Subset	Baseline period	
	1895-1924	1961-1990
<i>Unadjusted</i>		
Good quality	0.21	0.24
Intermediate quality	0.17	0.21
Poor quality	0.36	0.34
Bad quality	0.48	0.55
<i>Time-of-observation adjusted</i>		
Good quality	0.42	0.41
Intermediate quality	0.37	0.39
Poor quality	0.54	0.51
Bad quality	0.83	0.91
<i>Time-of-observation and step-change adjusted</i>		
Good quality	0.64	0.64
Intermediate quality	0.68	0.69
Poor quality	0.69	0.69
Bad quality	0.57	0.57

**Table 3:** 1895-2011 linear trends for each of the subsets using different baseline periods, calculated by linear least squares fitting.

though there are some differences between the trends calculated using the different baseline periods, they are quite small. Some of these differences are probably due to fact that some of the stations (87 or  $\sim 7\%$ ) did not have enough data during the 1895-1924 period, and so were not included in the gridded means.

Instead, the main differences arise when we are comparing different subsets. We can see that, as expected from our discussion above, all plots converge near the baseline period. This creates the false impression of extra agreement between subsets near the baseline period compared to other periods. We can see this extra agreement is merely a statistical artefact, because when we change the baseline period from 1895-1924 to 1961-1990, the apparent period of greatest “agreement” also shifts.

A less obvious statistical artefact occurs for the *Time-of-observation and step-change adjusted* subsets. For this dataset, the National Climatic Data Center has applied a series of pairwise step-change adjustments to each of the station records in an attempt to reduce the amount of non-climatic step biases. We will discuss these adjustments further in Section 4.5, but here it is worth noting that these adjustments are carried out by adjusting the earlier portions of the records relative to the most recent portions. In other words, the annual adjustments for each record converge towards zero for the most recent



**Figure 10:** Comparison of the 11 point binomial smoothed temperature trends of each of the subsets, using two different anomaly baseline periods - top panels: 1895-1924 and bottom panels: 1961-1990.

1330 year[24]. Since the adjustments can be of either sign, 1364  
 1331 this means that applying these adjustments should 1365  
 1332 lead to a partial convergence between all adjusted 1366  
 1333 station records for the most recent year (in this case, 1367  
 1334 2011). 1368

1335 Since the biases of the poor and bad quality subsets 1369  
 1336 are often estimated by considering the *differences be-* 1370  
 1337 *tween subsets*, it is important to recognise these artificial 1371  
 1338 statistical convergences will reduce the apparent 1372  
 1339 differences. In particular, the apparent biases (i.e., 1373  
 1340 the differences between subsets) will be substantially 1374  
 1341 reduced during whatever baseline period is chosen. 1375  
 1342 In addition, for the *Time-of-observation and step-* 1376  
 1343 *change adjusted* datasets, the biases will be artificially 1377  
 1344 reduced for recent years. For this reason, it is impor- 1378  
 1345 tant to estimate the extent of the poor station quality 1379  
 1346 biases over the entire 1895-now period, even if it is 1380  
 1347 believed that the biases were mostly introduced in 1381  
 1348 recent decades[1, 2]. Of the previous analyses, only 1382  
 1349 Fall et al., 2011 considered the effects of the first arti- 1383  
 1350 ficial convergence, i.e., the effect of changing baseline 1384  
 1351 periods[14]. The artificial convergence introduced by 1385  
 1352 the National Climatic Data Center’s step-change ad- 1386  
 1353 justments does not appear to have been considered 1387  
 1354 before now. 1388

## 1355 4.2 Over-reliance on linear trend 1389 1356 analysis 1390

1357 Wunsch, 1999[57] and Percival & Rothrock, 2005[58] 1391  
 1358 have illustrated how random, trend-less data series 1392  
 1359 frequently produce apparent “linear trends”, but that 1393  
 1360 these are spurious. This is well known in the field 1394  
 1361 of time series analysis, but Wunsch, 1999 correctly 1395  
 1362 notes that it is frequently overlooked by researchers 1396  
 1363 studying climate trends[57]. We would like to stress 1397

1364 another point, which while also well-known in the 1365  
 1366 field of regression analysis, is often overlooked by re- 1367  
 1368 searchers considering temperature trends - namely, 1368  
 1369 that, even if there are genuine trends in the data, 1369  
 1370 if the trends are non-linear, calculating their “linear 1370  
 1371 trends” can often be misleading. 1371

1372 All five of the previous studies attempting to quan- 1372  
 1373 tify the significance of the Surface Stations results use 1373  
 1374 linear trend analysis[2, 14–17], and even the prelim- 1374  
 1375 inary qualitative analysis of Watts, 2009 uses linear 1375  
 1376 trends to discuss one of their figures (Figure 23 of 1376  
 1377 Watts, 2009)[1]. We also presented much of our dis- 1377  
 1378 cussion in terms of linear trend analysis. But, as we 1378  
 1379 discussed in Section 3, the data we are discussing is 1379  
 1380 quite non-linear in nature. So, it is important to con- 1380  
 1381 sider the limitations of linear analysis. 1381

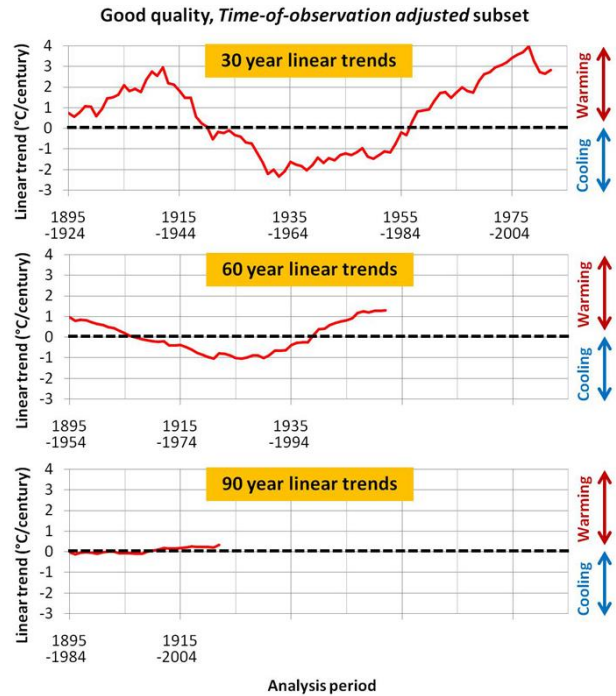
1382 There certainly is considerable convenience in cal- 1382  
 1383 culating the linear trend for a data series. It provides 1383  
 1384 a relatively simple, single value for describing an en- 1384  
 1385 tire data series, and is included as a standard tool 1385  
 1386 in most statistical data analysis packages, and even 1386  
 1387 spreadsheet software. However, the problem occurs 1387  
 1388 when it is used on data series that have substantial 1388  
 1389 non-linearities in their trends. A “linear trend” can 1389  
 1390 nominally be calculated for *any* two-dimensional se- 1390  
 1391 ries of data points, regardless of whether the series 1391  
 1392 is linear or not, but if the series is non-linear then 1392  
 1393 its “linear trend” might be completely meaningless. 1393  
 1394 Aside from the fact that, by framing their analysis 1394  
 1395 in terms of linear trends, researchers may become bi- 1395  
 1396 ased into prematurely expecting the trends to have an 1396  
 1397 underlying “linear behaviour”, this means that non- 1397  
 1398 linear trends will be overlooked. 1398

1399 Some researchers take a more sophisticated ap- 1399  
 1400 proach to linear trend analysis, by calculating the 1400  
 1401 significance of the linear fits, e.g., the Martinez et al., 1401

1400 2012 study used three separate statistical tests for  
 1401 assessing the significance of their linear trends[16].  
 1402 However, while this is preferable to using linear  
 1403 trends without testing, it is still an inadequate ap-  
 1404 proach, if the data series show non-linear trends. This  
 1405 is particularly the case, when the number of data  
 1406 points in the series is low (a point Martinez et al.,  
 1407 2012 concede[16]). For instance, if a data series con-  
 1408 tains only two points, then the “linear fit” of the  
 1409 points will nominally be a “perfect fit”. Of course, in  
 1410 reality, this tells us nothing about whether the data  
 1411 series is genuinely linear or not.

1412 We have intentionally plotted the trends of the  
 1413 twelve subsets in Figures 4, 5 and 6 without including  
 1414 their “linear trends”, as we have found that whenever  
 1415 a plot is shown with a linear trend, many readers will  
 1416 then mentally “see” a linear trend in the data, even  
 1417 if the actual trends is totally non-linear, e.g., cycli-  
 1418 cal. This is perhaps an artefact of the way the hu-  
 1419 man brain has evolved a high tendency towards *false*  
 1420 *pattern recognition*[57, 59]. In any case, the danger  
 1421 of this tendency can be illustrated by comparing the  
 1422 subsets in any of those figures. All of our subsets have  
 1423 positive linear trends when considered over the entire  
 1424 1895-2011 period (see Table 1). So, if a researcher re-  
 1425 lied on those linear trends for their comparison, they  
 1426 might find all of the subsets to be fairly similar, in  
 1427 that they all “showed warming”. However, it can be  
 1428 seen from Figure 8 that there are actually substan-  
 1429 tial differences between the subsets. For instance, for  
 1430 the time-of-observation adjusted datasets in Figure  
 1431 5, the good quality subset appears to have alternated  
 1432 between periods of “warming” and “cooling”, each  
 1433 lasting several decades. In contrast, the bad qual-  
 1434 ity subset suggests an almost continuous “warming”  
 1435 trend. Both subsets have a positive linear trend for  
 1436 the period 1895-2011, but for the good quality subset,  
 1437 this is merely a consequence of the 1890s-1900s being  
 1438 cooler than the 1990s-2000s, i.e., if a different start-  
 1439 ing point, and period length was chosen, both the  
 1440 sign and the magnitude of the “linear trend” could  
 1441 change.

1442 Figure 11 shows some of the different linear trends  
 1443 which can be calculated for the time-of-observation  
 1444 adjusted good exposure subset by merely varying the  
 1445 start year and the length of time over which the trend  
 1446 is calculated. Depending on the start year, periods of  
 1447 both “cooling” and “warming” can be found. In ad-  
 1448 dition, the magnitude of the trends depends strongly  
 1449 on the number of years included in the calculation.  
 1450 For instance, the maximum and minimum trends for

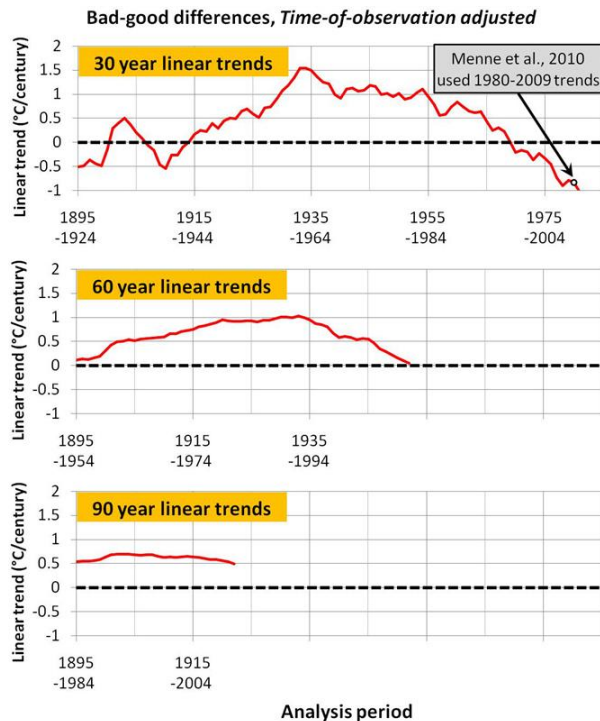


1451 **Figure 11:** Linear trends, using different starting years  
 1452 and period lengths, for the time-of-observation adjusted  
 1453 good quality subset.

1451 the 30 year calculations are an order of magnitude  
 1452 greater than the maximum and minimum trends for  
 1453 the 90 year calculations (+3.96 and -2.33 °C/century  
 1454 compared to +0.33 and -0.11 °C/century). If the  
 1455 temperature trends were truly linear then the values  
 1456 of the linear trends should not depend on either the  
 1457 start year or the number of years included.

1458 The *difference* between the bad quality and good  
 1459 quality subsets is slightly better characterised by a  
 1460 linear fit - see Figure 12. Although the shorter  
 1461 time-scale linear fits show considerable variability, the  
 1462 mean values for each of the time-scales are of the same  
 1463 order of magnitude (see Table 4). In addition, aside  
 1464 for a few of the 30 year estimates, the trends are all  
 1465 of the same sign, i.e., indicating a warming of the bad  
 1466 quality stations relative to the good quality stations.  
 1467 This slightly better fit is also confirmed by the better  
 1468 fitting coefficients (compare the  $r^2$  values in Tables 1  
 1469 and 2).

1470 However, even with these better fits, the 30 year  
 1471 estimates still show considerable variability over the  
 1472 entire period. The bias appears to be greatest for the  
 1473 1934-1963 period. Although most periods suggest a  
 1474 warming bias, for three short periods, the bias ap-



**Figure 12:** Linear trends, using different starting years and period lengths, for the difference between the good and bad quality time-of-observation adjusted subsets. For comparison with Menne et al., 2010[2], the 1980-2009 linear trend is highlighted.

pears to be a cooling one, i.e., the trends beginning during the period 1895-1900, 1907-1914 or 1970-1982 (ending in 1924-1929, 1936-1943 and 1999-2011 respectively) are negative. We note that much of the analysis by Menne et al., 2010 was based on 30 year linear trends during one of these negative trend periods - their analysis was on the 1980-2009 linear trends[2]. This might have contributed to their conclusion that poor station siting leads to a cooling bias, i.e., the opposite of our conclusions.

It can be seen that linear trends provide an overly simplistic description of non-linear time series, such as the subset temperature trends and difference trends we are considering in this study. Therefore, we do not recommend relying heavily on linear analysis in this case. Most of the previous analyses have relied very heavily on linear analysis, which limits the validity of their conclusions.

Having said that, linear trends do offer crude approximations of the overall trends, and can provide a helpful method for **crudely** comparing subsets and

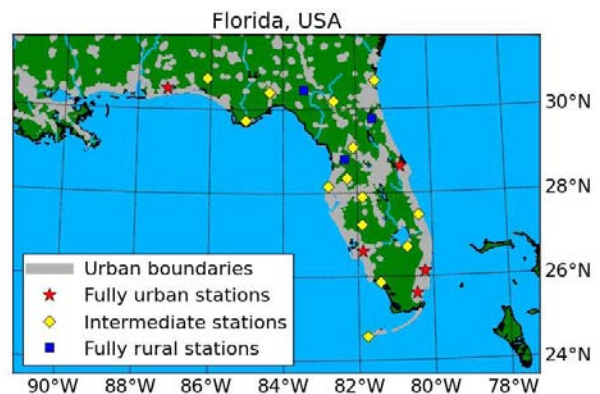
biases. Hence, we recognise they can be useful, *provided that the non-linear nature of the data is adequately discussed.*

Period of fit	Linear trends ( $^{\circ}\text{C}/\text{century}$ )		
	Mean	Minimum	Maximum
30 years	+0.39	-1.14	+1.55
60 years	+0.63	+0.05	+1.02
90 years	+0.62	+0.49	+0.70
Whole range	+0.41	N/A	N/A

**Table 4:** Mean, maximum and minimum linear trends of the difference between the bad and good subsets, for each of the time-scales used in Figure 12.

### 4.3 Neglecting urbanization bias

In a series of three papers, we discuss how urbanization bias has introduced artificial warming trends into many of the station records[11–13]. Although this problem seems to be more serious for the non-U.S. component of the Global Historical Climatology Networks, it is still a significant problem for the U.S. component.



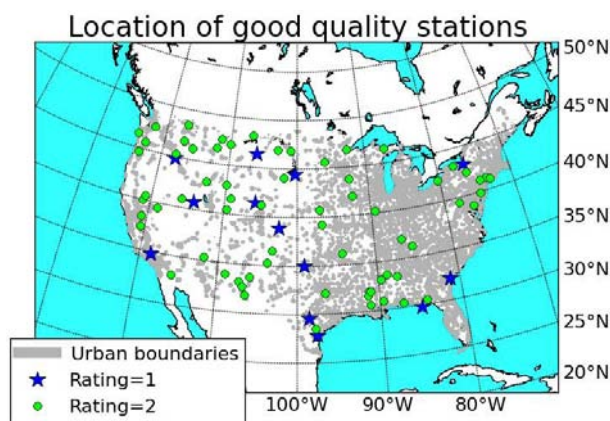
**Figure 13:** Location of the Florida stations used by the Martinez et al., 2012 study[16], and the degree of urbanization of those stations.

As the U.S. Historical Climatology Network is part of the Global Historical Climatology Network, we can use the estimates of station urbanization provided with the Global Historical Climatology Network dataset[60], to estimate the extent of urbanization in the U.S. network. The National Climatic Data Center provide two estimates of urbanization for their Global Historical Climatology Network stations - one based on the population in the general area of the

1516 station and one based on the 1994/1995 night light  
1517 brightness at the station. Under each estimate, a station  
1518 can be “rural”, “sub-urban” or “urban”. As we  
1519 did in Ref. [13], we will here define a station as “fully  
1520 rural” if it is identified as “rural” by both estimates,  
1521 “fully urban” if it is identified as “urban” by both  
1522 estimates, and otherwise “intermediate”.

1523 Using these definitions, of the 1218 U.S. Historical  
1524 Climatology Network stations, 277 (22.7%) are fully  
1525 rural, 842 (69.1%) are intermediate and 99 (8.1%)  
1526 are fully urban. As we discuss in Ref. [13], this  
1527 is relatively rural compared to the averages for the  
1528 Global Historical Climatology Network. However,  
1529 only 22.7% of the stations are fully rural, so it is  
1530 likely that a substantial fraction of the stations are  
1531 affected by at least some urbanization bias.

1532 Since urbanization bias generally leads to an artificial  
1533 warming trend, it can easily affect estimates of the  
1534 separate bias due to the station siting. For instance,  
1535 because the Martinez et al., 2012 study focused on a  
1536 single state (Florida), only two of the stations were  
1537 rated 5, and only five were rated 1 or 2. As can be  
1538 seen from Figure 13, only three of the 22 Florida  
1539 stations considered by Martinez et al., 2012 are  
1540 identified as fully rural (blue squares). So, it is  
1541 likely that much of the differences between the  
1542 different stations are actually due to different amounts  
1543 of urbanization bias. All three of the fully rural  
1544 stations (Federal Point, Inverness 3 SE and Madison)  
1545 had the same Surface Stations rating of 4, and so it  
1546 would not have been possible for them to just use the  
1547 fully rural stations for investigating the effects of the  
1548 different ratings.



**Figure 14:** Locations of good quality stations, showing urbanized regions. Urban boundaries were determined using the GRUMP dataset[61].

1549 The National Climatic Data Center have argued  
1550 that their data homogenization should have removed  
1551 much of the urbanization bias from the *Time of ob-*  
1552 *servatation and step-change adjusted* dataset[24, 62].  
1553 However, as we discuss in Ref. [13], their step-change  
1554 homogenization method is inadequate for removing  
1555 urbanization bias, especially if the neighbours they  
1556 use for homogenization are also affected by urban-  
1557 ization bias. Indeed, Pielke et al., 2007b note that  
1558 such step-change homogenization algorithms provide  
1559 inaccurate results when applied to stations with trend  
1560 biases, such as urbanization bias[10].

1561 For the COOP neighbours used by the National  
1562 Climatic Data Center for homogenizing the Historical  
1563 Climatology Network, we cannot use the Global  
1564 Historical Climatology Network estimates. However,  
1565 we can estimate their urbanization by assigning the  
1566 station co-ordinates to the gridded urbanization es-  
1567 timates from the GRUMP dataset (see Ref. [61] for  
1568 details on the GRUMP dataset). This approach iden-  
1569 tifies 45.0% of the COOP neighbours as being “ur-  
1570 ban” (for comparison, 55.6% of the U.S. Historical  
1571 Climatology Network stations are urban with this  
1572 approach). In other words, nearly half of the neigh-  
1573 bours used for homogenizing the stations are urban.  
1574 Since, this ratio is higher in urbanized areas, it is  
1575 likely that the step-change homogenization approach  
1576 used for generating the *Time-of-observation and step-*  
1577 *change adjusted* dataset has led to “urban blending”  
1578 into many of the rural stations which are near urban-  
1579 ized areas.

1580 We can see from Figure 14 that a large fraction of  
1581 the good quality stations are in urbanized areas. The  
1582 sample size of the good quality stations is already  
1583 quite low (8%), so there are probably not enough  
1584 stations to reliably estimate U.S. regional trends by  
1585 considering just the stations which are both of a good  
1586 quality and rural. If the siting quality of the Global  
1587 Historical Climatology Network is similar to that of  
1588 the U.S. network, then the problem of separating the  
1589 two confounding biases is likely to be even more chal-  
1590 lenging for the more urbanized global network.

1591 It may be difficult to obtain a large enough sample  
1592 of stations that are unaffected by either urbanization  
1593 bias or inadequate station exposure, for directly as-  
1594 sessing genuine climatic trends. Still, if the extent of  
1595 urbanization bias is comparable in each subset, then  
1596 it may still be possible to reliably estimate the bias in-  
1597 troduced by inadequate station exposure, separately.  
1598 This is because, we can then assume the urbanization  
1599 bias will be roughly the same magnitude in each of

1600 the subsets. In other words, while the temperature  
 1601 trends of the subsets may be affected by urbanization  
 1602 bias, the *differences between the trends* should not be  
 1603 overly affected.

1604 Table 5 shows the degree of urbanization of each of  
 1605 the subsets. There are some differences between the  
 1606 subsets. For example, both the good quality (Rat-  
 1607 ings 1 & 2) and the bad quality (Rating 5) stations  
 1608 have a relatively high percentage of fully urban sta-  
 1609 tions. Similarly, there is a relatively high percentage  
 1610 of fully rural stations of Rating 3, although a quite  
 1611 low percentage for Rating 5. However, in general, the  
 1612 ratios seem to be *roughly* similar to the average for  
 1613 the entire U.S. Historical Climatology Network (i.e.,  
 1614 the bottom row in Table 5). Therefore, it is probably  
 1615 still a reasonable approximation, *provided* we recog-  
 1616 nize that the estimates of the station exposure biases  
 1617 will be affected by the differences in urbanization  
 1618 bias between subsets. In some cases, these urbanization  
 1619 differences may be quite large, as in the Martinez et  
 1620 al., 2012 study.

Rating	Degree of urbanization		
	Fully rural	Intermediate	Fully urban
1 & 2	20.0%	61.3%	18.8%
3	27.2%	62.7%	10.1%
4	23.3%	69.8%	6.9%
5	8.2%	73.8%	18.0%
Unrated	21.8%	75.4%	2.8%
Average	22.7%	69.1%	8.1%

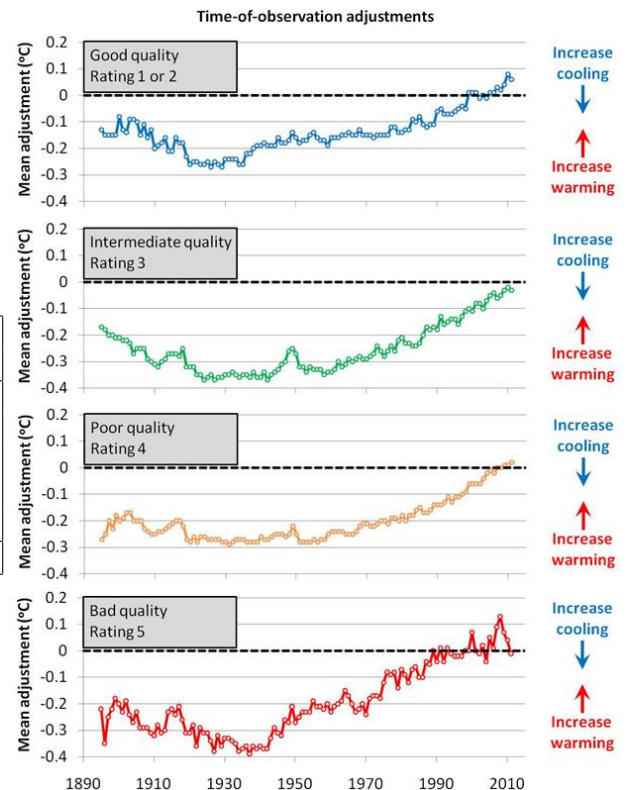
1621 **Table 5:** *The percentage urbanization of each of the*  
 1622 *Surface Stations subsets.*

#### 1621 4.4 Are the time-of-observations 1622 adjustments reliable?

1623 One of the main homogeneity adjustments applied to  
 1624 the U.S. Historical Climatology Network records is to  
 1625 correct for non-climatic biases introduced by docu-  
 1626 mented changes in the time at which observers made  
 1627 their measurements at individual stations. Karl et  
 1628 al., 1986 calculated that if observers were estimating  
 1629 daily mean temperatures from a single daily reading  
 1630 of a minimum-maximum thermometer, then the time  
 1631 of day in which they made their observation would  
 1632 have significantly affected the mean monthly temper-  
 1633 atures which they obtained[25]. They calculated sta-  
 1634 tistical estimates for this “time-of-observation bias”,  
 1635 which vary seasonally (from month to month), and  
 1636 regionally, e.g., the magnitude of the biases tends to

1637 increase with latitude because of the larger daily tem-  
 1638 perature range.

1639 With this in mind, when the National Climatic  
 1640 Data Center was compiling the U.S. Historical Clima-  
 1641 tology Network, they collected the observation times  
 1642 associated with each station, whenever they were re-  
 1643 ported. They then applied the corresponding Karl et  
 1644 al., 1986 adjustments to each of the station records for  
 1645 two of their releases (i.e., the *Time-of-observation ad-*  
 1646 *justed* and *Time-of-observation and step change ad-*  
 1647 *justed* datasets).



1648 **Figure 15:** *Mean time-of-observation adjustments ap-*  
 1649 *plied by the National Climatic Data Center to the grid-*  
 1650 *ded subsets.*

1648 Since the U.S. station histories indicate a general  
 1649 reduction in the number of evening observers and in-  
 1650 crease in the number of morning observers, the time-  
 1651 of-observation adjustments have the net effect of in-  
 1652 troducing a post-1930s warming trend into the U.S.  
 1653 temperature trends<sup>3</sup> - see Figure 15. In general, each  
 1654 of the homogeneity adjustments developed by the Na-  
 1655 tional Climatic Data Center for the U.S. Historical

<sup>3</sup>Although, the pre-1930s adjustments introduce a cooling trend.

1656 Climatology Network coincidentally seem to intro- 1707  
1657 duce “more warming” (e.g., see our analysis in Ref. 1708  
1658 [13]). This seems to have led to some cynicism about  
1659 their reliability[1, 17, 63]. In particular, Balling &  
1660 Idso, 2002 criticised the time-of-observation adjust-  
1661 ments and argued that the unadjusted data appeared  
1662 to better match satellite and radiosonde estimates of  
1663 surface temperature changes[63]. In response, Vose  
1664 et al., 2003[26] updated Karl et al., 1986’s study us-  
1665 ing a larger dataset. They found similar results to  
1666 Karl et al., 1986, and so concluded that their original  
1667 time of observation bias adjustments were reasonably  
1668 accurate.

1669 A detailed assessment of the time-of-observation  
1670 bias adjustments used by the National Climatic Data  
1671 Center is beyond the scope of this article. We  
1672 will note that our preliminary calculations suggest  
1673 that their adjustments are reasonable, provided the  
1674 archived station histories for the U.S. Historical Cli-  
1675 matology Network are accurate, and that the time-  
1676 of-observation changes do not involve other station  
1677 changes. For this reason, we will assume that the  
1678 *Time-of-observation adjusted* dataset is more reliable  
1679 than the *Unadjusted* dataset.

1680 The adjustments are of a relatively large magni- 1731  
1681 tude, however, e.g., doubling the linear trends of the 1732  
1682 good, intermediate and bad quality subsets - see Ta- 1733  
1683 ble 1. So, a more careful assessment of their reli- 1734  
1684 ability may be warranted. In particular, it might 1735  
1685 be worth checking whether the documented time-of- 1736  
1686 observation changes also coincided with other non- 1737  
1687 climatic changes, e.g., changes in instrumentation or 1738  
1688 station location. One way to assess the adjustments 1739  
1689 would be to compare the differences between the tem- 1740  
1690 perature record and the records of the station’s neigh- 1741  
1691 bours before and after adjustments. 1742

1692 With this in mind, the estimates of the biases in- 1743  
1693 troduced from inadequate station exposure which are 1744  
1694 determined from the *Time-of-observation adjusted* 1745  
1695 dataset should be more reliable than those deter- 1746  
1696 mined from the *Unadjusted* dataset. If the time-of- 1747  
1697 observation adjustments are similar for all stations, 1748  
1698 on average, then these estimates should be the same, 1749  
1699 in which case it would be irrelevant which dataset 1750  
1700 was used. However, it can be seen from Figure 15 1751  
1701 that there are slight differences in the mean adjust- 1752  
1702 ments applied to each of the subsets. Part of this 1753  
1703 may be a statistical artefact, due to some of the sub- 1754  
1704 sets having relatively small sample sizes. But, it is 1755  
1705 plausible that some of the factors influencing station 1756  
1706 quality may also be related to observer practice, in 1757

1707 which case there may be a genuine difference between 1708  
1709 subsets.

1709 The main difference between the subsets is that 1710  
1711 the adjustments during the 1930s are slightly greater 1712  
1713 for the bad quality subset. This increases the long- 1714  
1715 term trend of the difference between the bad quality 1716  
1717 and good quality subsets (see Figure 8), and therefore 1718  
1719 the estimate of the bias, as can be seen from Table 2. 1720  
1721 There is also a slight reduction in the estimate of the 1722  
1723 bias for the poor quality subset. 1724  
1725

1726 Watts et al. (in preparation, 2012) did not use 1727  
1728 the *Time-of-observation adjusted* dataset for their 1729  
1730 estimates of the biases, favouring the *Unadjusted* 1731  
1732 dataset[17]. As explained above, if the mean time-of- 1733  
1734 observation adjustments are similar for all subsets, 1735  
1736 then this should not matter. They used a differ- 1737  
1738 ent rating system which created a more evenly di- 1739  
1740 vided group of subsets, than the one we used. So, it 1741  
1742 is plausible that this could have reduced the differ- 1743  
1744 ences between the adjustments applied to each sub- 1745  
1746 set. If this is the case, then the application of *Time-* 1747  
1748 *of-observation* adjustments should not substantially 1749  
1750 alter the findings of Watts et al., but this would need 1751  
1752 to be tested. 1753  
1754

1755 As this article is focused on the reliability of the 1756  
1757 U.S. station records, the issue of time-of-observation 1758  
1759 biases in other networks is a subject for a separate 1759  
1760 study. However, we do think it is worth mentioning 1760  
1761 that systematic changes in observation time are *not* 1761  
1762 confined to the U.S. For instance, in August 1961, 1762  
1763 stations in the People’s Republic of China changed 1763  
1764 from taking measurements at 01, 07, 13 and 19h local 1764  
1765 mean solar time to 02, 08, 14 and 20h Beijing time, 1765  
1766 i.e., local mean solar time at 120°E [64]. On 1st July 1766  
1767 1961, the observation times at all Canadian airport 1767  
1768 weather stations were changed from 12h UTC (for 1768  
1769 maximum temperatures) and 00h UTC (for minimum 1769  
1770 temperatures) to 06h UTC[65]. In Japan, stations 1770  
1771 generally have a time-of-observation of midnight (00h 1771  
1772 JST). But, before 1940, it was 22h JST and during 1772  
1773 the period 1953-1963, minimum temperatures were 1773  
1774 observed separately at 09h JST[66]. 1774

1775 For this reason, it is likely that time-of-observation 1775  
1776 adjustments are also necessary for other countries. In 1776  
1777 the case of the U.S., the bias seems to have led to a 1777  
1778 significant artificial cooling trend over the 20th cen- 1778  
1779 tury[26]. Similar biases, either cooling *or* warming, 1779  
1780 may also exist for many of the stations in the Global 1780  
1781 Historical Climatology Network, and they may have 1781  
1782 led to systematic biases for individual countries. For 1782  
1783 instance, in Australia, a change in the official obser- 1783



1758 vation time in 1964 to 0900 local time appears to  
1759 have introduced an artificial warm bias in the records  
1760 relative to the pre-1964 periods[67]. Also, the U.S.  
1761 time-of-observation adjustments suggest a pre-1930s  
1762 warming bias (see Figure 15). Until the signs and  
1763 magnitudes of these biases can be reliably determined  
1764 for the Global Historical Climatology Network stations,  
1765 the records for those stations should be treated  
1766 cautiously.

1767 If future researchers attempt to track down obser-  
1768 vation times for stations outside of the U.S., it  
1769 may also be important to check if time of observa-  
1770 tion adjustments have already been applied to the  
1771 station records. Jones et al., 1986[56] suggest that  
1772 time-of-observation bias corrections may have already  
1773 been applied to some records, before incorporation in  
1774 datasets, e.g., “*Such adjustments were used in [the*  
1775 *World Weather Records dataset] for the United States*  
1776 *up to 1940 or 1950 (depending on the station)”*[56].  
1777 The U.S. Historical Climatology Network was com-  
1778 piled directly from the COOP archive[25]. So, it is  
1779 likely that their source records had not been corrected  
1780 for time of observation. However, the Global Histor-  
1781 ical Climatology Network dataset was compiled from  
1782 pre-existing datasets (including the World Weather  
1783 Records), so it is possible that some of its station  
1784 records have already had time-of-observation correc-  
1785 tions, while others have not.

#### 1786 4.5 Are the step-change adjustments 1787 reliable?

1788 From Figure 8, it is clear that the differences between  
1789 the subsets are considerably reduced in the *Time-of-*  
1790 *observation and step-change adjusted* dataset, com-  
1791 pared to the *Unadjusted* and *Time-of-observation ad-*  
1792 *justed* datasets. In this sense, the step-change adjust-  
1793 ments have led to a greater “homogeneity” between  
1794 individual station records. However, this does *not*  
1795 actually tell us whether the more homogeneous records  
1796 are more or less climatically representative.

1797 As we mentioned in Section 3, one explanation for  
1798 the greater homogeneity is that the step-change ad-  
1799 justments have succeeded in substantially reducing  
1800 the magnitude of the non-climatic biases in the sta-  
1801 tion records. If this is the case, then Menne et al.,  
1802 2010 are correct in concluding that the station quality  
1803 problems identified by the Surface Stations project  
1804 do not seriously affect the *Time-of-observation and*  
1805 *step-change adjusted* estimates of U.S. temperature  
1806 trends[2]. But, another explanation is that the

1807 greater homogeneity arises because the non-climatic  
1808 biases have been averaged between neighbouring sta-  
1809 tions through “blending” of the biases, rather than  
1810 removed.

1811 In other words, the records with the most non-  
1812 climatic biases will have some of their biases correctly  
1813 removed by the process, but the records with the  
1814 least non-climatic biases will have the non-climatic  
1815 biases of their neighbours introduced by the pro-  
1816 cess[55]. If this is the case, then Watts et al. (in  
1817 preparation, 2012) are correct in concluding that the  
1818 homogenization process has reduced the reliability of  
1819 the good quality station records, and that the non-  
1820 homogenized records are therefore more reliable, even  
1821 if they have not been corrected for non-climatic bi-  
1822 ases[17]. Essentially, Watts et al. are arguing that,  
1823 while the *Unadjusted* records probably contain non-  
1824 climatic biases, on average, they are still more re-  
1825 liable than the *Time-of-observation and step-change*  
1826 *adjusted* records.

1827 As we discussed in Section 4.4, we are assuming  
1828 the *Time-of-observation adjusted* dataset is more cli-  
1829 matically representative than the *Unadjusted* dataset  
1830 favoured by Watts et al. (in preparation, 2012). How-  
1831 ever, we also saw from Figure 15 that the net effect  
1832 of these adjustments is broadly similar for all of the  
1833 subsets. Indeed, the largest adjustments were for the  
1834 bad quality subset, and as a result, had the effect  
1835 of *increasing* the bias estimates for the bad quality  
1836 subset. Hence, the main question that needs to be  
1837 resolved is whether the step-change adjustments im-  
1838 prove or reduce the reliability of the records.

1839 With this in mind, it is worth carefully consider-  
1840 ing how the step-change adjustments are carried out.  
1841 When the Surface Stations project began, the Na-  
1842 tional Climatic Data Center were still using the Karl  
1843 & Williams, 1987 algorithm[68] for their step-change  
1844 homogenization of the U.S. Historical Climatology  
1845 Network dataset. However, in 2009, they updated the  
1846 dataset (see Menne et al., 2009[24]), and began using  
1847 the newer Menne & Williams, 2009 algorithm[27].

1848 There are several important differences between  
1849 the two algorithms, e.g., the Karl & Williams, 1987  
1850 algorithm only tested for non-climatic biases when  
1851 it was known from the station histories that a sig-  
1852 nificant change had occurred, while the Menne &  
1853 Williams, 2009 algorithm can test for both docu-  
1854 mented and undocumented changes. However, for  
1855 the purposes of the following discussion, both algo-  
1856 rithms are equivalent. Pielke et al., 2007b[10]’s criti-  
1857 cism of the Karl & Williams, 1987 algorithm is similar

1858 to ours, and so also applies to the current Menne & 1909  
1859 Williams, 2009 algorithm. 1910

1860 Both algorithms work by comparing the tempera- 1911  
1861 ture records of each “target” station to the records 1912  
1862 of a large number of neighbouring stations (20 for 1913  
1863 Karl & Williams, 1987; 40 for Menne & Williams, 1914  
1864 2009). In both algorithms, the years of potential 1915  
1865 station changes are identified from the station history 1916  
1866 files. But, the Menne & Williams, 2009 algo- 1917  
1867 rithm also looks for additional undocumented station 1918  
1868 changes, by comparing the target to each neighbour 1919  
1869 in turn, and noting the approximate years of any ap- 1920  
1870 parent discrepancies between the stations. If an ap- 1921  
1871 parent discrepancy occurs at around the same time 1922  
1872 in several neighbour-target comparison, then it is as- 1923  
1873 sumed that this corresponds to a station change, or 1924  
1874 “change-point” in the target station. 1925

1875 It should be noted that, if similar biases occur 1926  
1876 around the same time in enough neighbours, then this 1927  
1877 could cause the algorithm to incorrectly identify the 1928  
1878 target station as having the bias. This means that 1929  
1879 if biases are relatively frequent amongst the neigh- 1930  
1880 bours, any target stations which are genuinely unbi- 1931  
1881 ased will be mistakenly treated as non-climatic “out- 1932  
1882 liers”, while the biased records will be mistakenly 1933  
1883 treated as the climatic trends. Stations which are 1934  
1884 heavily biased should be correctly identified as “out- 1935  
1885 liers”, and their biases would be reduced, however 1936  
1886 only enough to match the average of the neighbours 1937  
1887 - which would themselves be biased. This would lead 1938  
1888 to a blending (or “homogenization”) of the biases 1939  
1889 amongst all the stations in the vicinity, whether bi- 1940  
1890 ased or not. 1941

1891 To estimate the sign and magnitude of any non- 1942  
1892 climatic biases introduced by the proposed change- 1943  
1893 points, the target station is again compared to the 1944  
1894 neighbours, one at a time. The temperature record of 1945  
1895 each neighbour is subtracted from the target’s record, 1946  
1896 to construct a “difference series”. The mean temper- 1947  
1897 ature difference of the segment after the proposed 1948  
1898 change-point is compared to the mean temperature 1949  
1899 difference of the segment before the proposed change- 1950  
1900 point. If the difference between the two means is 1951  
1901 greater than a certain threshold value, then it is con- 1952  
1902 sidered as a possible step-change bias. The differ- 1953  
1903 ent estimates of this proposed bias are compared for 1954  
1904 each neighbour, and if there is sufficient agreement 1955  
1905 between the estimates, then all of the years up to 1956  
1906 the year of the change-point are adjusted by that 1957  
1907 value[27, 68]. 1958

1908 As deGaetano, 2006[69] and Pielke et al., 2007b[10] 1959

note, and as we discuss in Ref. [13], these step- 1909  
change adjustments are inadequate for correcting 1910  
trend biases, such as urbanization bias. The Menne 1911  
& Williams, 2009 algorithm does actually con- 1912  
sider trend biases when identifying potential change- 1913  
points[27]. However, since they then treat the iden- 1914  
tified biases as step-change biases, the algorithm is 1915  
still inadequate. For instance, let us suppose that 1916  
the algorithm correctly identifies the start of a trend 1917  
bias, and the trend bias is roughly linear in nature. 1918  
In that case the *mean* value of the bias will only be 1919  
half of the bias at the end of the segment, and the 1920  
adjustment will be incomplete. 1921

Pielke et al., 2007b[10] tested the effectiveness of 1922  
the Karl & Williams, 1987[68] homogeneity adjust- 1923  
ments algorithm. They simulated 1000 temperature 1924  
records, and introduced artificial step change biases 1925  
and trend biases into these simulated records. When 1926  
they applied the Karl & Williams algorithm to these 1927  
records, they found the algorithm overestimated the 1928  
magnitude of positive step changes, and underesti- 1929  
mated the magnitude of negative step changes if the 1930  
station being homogenized also had a warming trend 1931  
bias[10]. This confirms that the statistical “alias- 1932  
ing” effect described above is a problem when station 1933  
records are affected by trend biases. 1934

An additional problem of trend biases is that, if a 1935  
neighbouring station suffers from a trend bias, then 1936  
this will increase (or decrease) the estimates of any 1937  
potential step-change adjustments. 1938

Next, we shall consider how these algorithms op- 1939  
erate on the U.S. Historical Climatology Network 1940  
datasets. However, before doing so, it is worth con- 1941  
sidering an additional change the National Climatic 1942  
Data Center adopted with their 2009 update[24]. 1943  
When they switched to using the Menne & Williams, 1944  
2009 algorithm, they also switched to using the entire 1945  
COOP Network for their station neighbours, rather 1946  
than just using the Historical Climatology Network 1947  
stations. 1948

There are at least three serious problems with this 1949  
decision: 1950

1. Because the National Climatic Data Center do 1951  
not currently provide easy access to the COOP 1952  
Network dataset, their decision makes it harder 1953  
for researchers to replicate and/or assess the re- 1954  
liability of their step-change adjustments. In 1955  
our case, we were only able to carry out our 1956  
analysis of the COOP neighbours because we 1957  
had downloaded a 2011 version of the COOP 1958  
datasets from a temporary folder on the National 1959

1960 Climatic Data Center's public ftp website (we  
1961 downloaded them from [ftp://ftp.ncdc.noaa.  
1962 gov/pub/data/williams/](ftp://ftp.ncdc.noaa.gov/pub/data/williams/)).

- 1963 2. The Historical Climatology Network stations  
1964 were partially chosen on the basis that they had  
1965 relatively long and complete station records[70].  
1966 As a result, the average period of overlap be-  
1967 tween neighbouring Historical Climatology Net-  
1968 work stations is much greater ( $78 \pm 12$  years) than  
1969 the overlap between Historical Climatology Net-  
1970 work stations and their nearest COOP stations  
1971 ( $29 \pm 6$  years) (see our discussion in Ref. [13]).  
1972 This shorter length means that the homogeniza-  
1973 tion algorithm is less likely to correctly identi-  
1974 fy the non-climatic biases in the records being  
1975 homogenized. The problem is also accentuated  
1976 by the Menne & Williams, 2009 algorithm, be-  
1977 cause they preferentially select neighbours which  
1978 have a higher correlation with the target sta-  
1979 tion[27]. This is a problem because poorly cor-  
1980 related station records are more likely to *appear*  
1981 well-correlated, if the period of overlap between  
1982 the records is short, as we illustrate in Ref. [13].

- 1983 3. The Surface Stations project only investigated  
1984 the U.S. Historical Climatology Network sta-  
1985 tions, and as a result, Surface Stations ratings  
1986 are not available for the 90% of COOP stations  
1987 which are not in the Historical Climatology Net-  
1988 work.

1989 On this last point, it seems reasonable to assume  
1990 that the COOP stations have a similar statistical dis-  
1991 tribution of ratings to the Historical Climatology Net-  
1992 work, i.e., something similar to Figure 1. The His-  
1993 torical Climatology Network was constructed from  
1994 COOP stations, so it is likely that it has a similar  
1995 distribution of ratings. Indeed, Fall et al., 2011 note  
1996 that several station observers were unaware whether  
1997 their station was part of the Historical Climatology  
1998 Network or not[14]. This suggests that the differ-  
1999 ences between COOP stations and Historical Clima-  
2000 tology Network stations were not even clear to the  
2001 observers. However, this does not tell us specifically  
2002 which neighbours are of good, intermediate, poor or  
2003 bad quality.

2004 Let us now consider the effects of the Menne &  
2005 Williams, 2009 step change adjustments on the U.S.  
2006 Historical Climatology Network stations. From our  
2007 discussion throughout this article, we can make sev-  
2008 eral predictions as to what these effects would be.

2009 We saw in Section 4.3 that a substantial percentage  
2010 of the U.S. stations are at least partially urbanized,  
2011 and hence many of the COOP neighbours used for ho-  
2012 mogenizing the Historical Climatology Network sta-  
2013 tions would be affected by urbanization bias. For this  
2014 reason, we would expect that stations with urbanized  
2015 neighbours would be affected by urban blending.

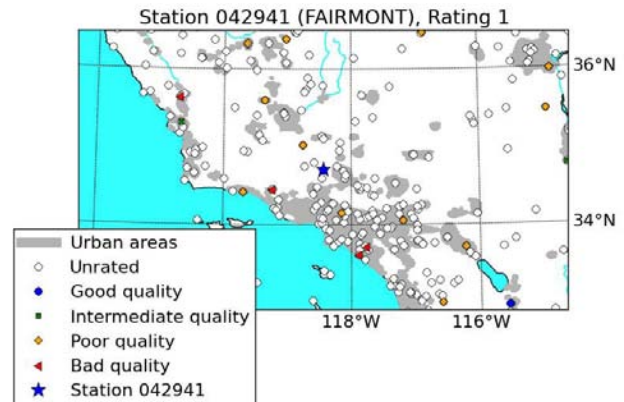


Figure 16: Locations of the neighbouring COOP stations in the vicinity of a good quality station, Fairmont.

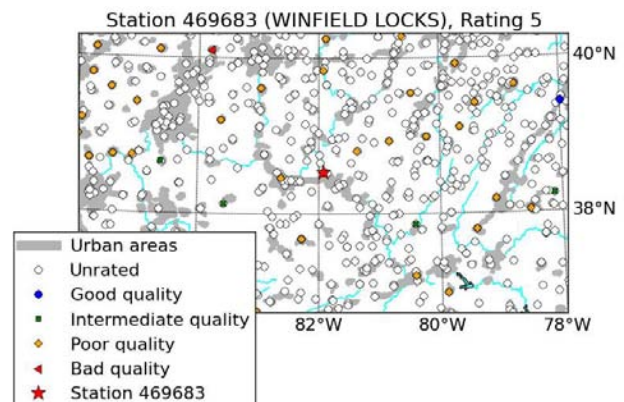


Figure 17: Locations of the neighbouring COOP stations in the vicinity of a bad quality station, Winfield Locks.

2016 The urbanization (warming) bias of the most heav-  
2017 ily affected stations would be slightly reduced, but in  
2018 a similar manner, many of the rural stations with no  
2019 urban bias would have an artificial warming trend in-  
2020 troduced by urban blending. Our analysis in Ref. [13]  
2021 suggests that this is the case. As mentioned earlier,  
2022 the application of step-change adjustments, such as  
2023 the Menne & Williams, 2009 algorithm, to trend bi-  
2024 ases substantially underestimates the biases because

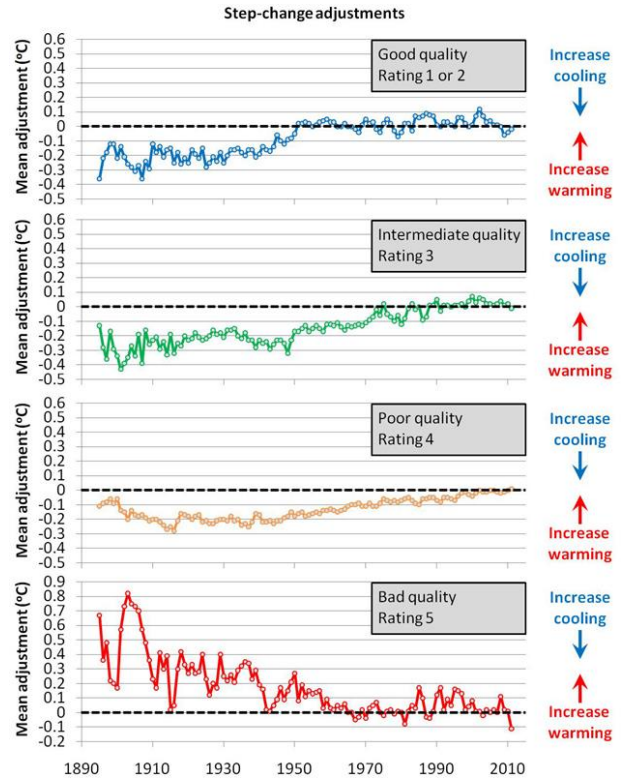
of statistical “aliasing” effects[10, 13, 69]. Hence, we would expect that, on average, the net effect of the adjustments would be the introduction of an artificial warming trend into all subsets that is roughly half the magnitude of the urbanization bias present in the *Time-of-observation adjusted* dataset.

Figures 16 and 17 show the COOP neighbouring stations in the vicinity of a typical rating 1 and rating 5 station, respectively. As expected, most of the neighbouring stations are unrated. But, in neither of the cases, are many of the rated neighbours of a good quality. If we assume that the COOP stations have a similar distribution of ratings to the Historical Climatology Network, then statistically, we would expect that only about 8% of the neighbours will be of a good quality. In other words, the good quality stations are likely to be considered as statistical outliers. The 6% of bad quality stations would also be considered statistical outliers, in a similar manner.

We would therefore expect that the trends of the “outlier” stations will be adjusted to better match those of the most common stations. Hence, we would expect that, while the warming trends of the bad quality stations would be partially reduced to better match those of the poor quality stations, the trends of the good quality would *incorrectly* be increased to better match those of the poor quality stations. The net effect of these adjustments would be to reduce the differences between the subsets.

As we mentioned in Section 4.1, the step-change adjustments are calculated retrospectively, i.e., when a non-climatic bias is identified, the most recent temperatures are assumed to be “correct”, and the earlier temperatures are adjusted to match. A consequence of this is that the net adjustments of all subsets will converge towards zero for the most recent years, introducing an artificial convergence between the subsets.

If these predictions are accurate, then the step-change adjustments would actually reduce the reliability of the records. At any rate, they would lead to unreliable estimates of the station exposure biases. The net adjustments applied to each subset are shown in Figure 18. They are certainly consistent with all of the predictions described above. For this reason, we suggest that the bias estimates from the *Unadjusted* and *Time-of-observation adjusted* datasets are currently the most accurate. This would invalidate the conclusions of Menne et al., 2010, who assumed that the estimates from the *Time-of-observation and step-change adjusted* dataset were the most reliable[2].



**Figure 18:** Mean step-change adjustments applied by the National Climatic Data Center to the gridded subsets using the Menne & Williams, 2009[27] algorithm.

Having said that, we must also acknowledge that there are also many other non-climatic biases in station records, e.g., station moves and changes in instrumentation. Many of these biases are probably the step-change biases which the Menne & Williams, 2009 algorithm was designed to remove[27]. Naïvely, it might be expected that, statistically, these biases should “cancel” each other out over time. However, non-climatic biases are often of the same sign and can lead to long-term apparent trends in the station records[71, 72].

Runnalls & Oke, 2006 suggest it is possible that this could have a tendency to introduce an apparent warming trend[72], in which case, a successful step-change adjustment algorithm would lead to net “cooling” adjustments. But, they point out that it could also lead to the opposite, in which case, successful step-change adjustments would lead to net “warming” adjustments. In other words, it is plausible that at least *some* of the “warming” introduced to the Historical Climatology Network by the step-change adjustments *may* be genuine. However, since

2098 it is expected for the reasons described above that  
2099 the step-change adjustments will mistakenly intro-  
2100 duce substantial artificial warming trends, it seems  
2101 unwise to make any assumptions about the correct  
2102 sign of the remaining non-climatic biases, yet.

## 2103 5 Recommendations to 2104 improve the reliability of 2105 climate records

2106 We saw that the majority of thermometer stations  
2107 in the U.S. Historical Climatology Network currently  
2108 suffer from inadequate siting, and that this has in-  
2109 troduced artificial warming trends into estimates of  
2110 U.S. temperature trends. It is quite likely that a sig-  
2111 nificant fraction of stations in the global network are  
2112 also affected by poor-siting problems. Indeed, from  
2113 preliminary assessments<sup>4</sup> of the Global Historical Cli-  
2114 matology Network stations from the rest of the world,  
2115 there seem to be a mixture of poor quality and good  
2116 quality stations.

2117 With this in mind, we believe it would be useful  
2118 if a global extension of the Surface Stations project  
2119 were carried out so that the magnitude of the siting  
2120 biases in the current estimates of global temperature  
2121 trends could be determined. Since such a project  
2122 would require considerable organization and interna-  
2123 tional collaboration, it might be useful to carefully  
2124 discuss how to assess the stations. In their initial  
2125 assessments, the Surface Stations team used Leroy,  
2126 1999[22]’s rating system, but Watts et al. (in prepara-  
2127 tion, 2012) suggest that the rating system proposed  
2128 by Leroy, 2010[31] is more appropriate. Decisions on  
2129 whether to use one of these two systems, or another  
2130 should be agreed at the start of the project. It would  
2131 be desirable to choose a rating system that would al-  
2132 low comparable estimates of past station exposures to  
2133 be at least approximated from the site documentation  
2134 which is sometimes included by station histories[10].

2135 According to Watts et al. (in preparation,  
2136 2012)[17], the Surface Stations team did not collect  
2137 enough information to adequately rate the shading  
2138 and ground cover associated with stations and their  
2139 ratings were instead predominantly based on heat

source/sink proximity. Changes in shading and/or  
2140 ground cover, e.g., from growth/removal of neigh-  
2141 bouring trees, can substantially alter micro-climate[3,  
2142 4]. So, future investigations should also investigate  
2143 these factors. 2144

2145 It would probably be sensible if the team involved  
2146 in the original Surface Stations project were con-  
2147 sulted before this type of project is begun on a global  
2148 scale, as they already have several years of valu-  
2149 able experience in carrying out station quality assess-  
2150 ments.

2151 Such a global project would probably have been  
2152 too overwhelming a task during the 1990s, when  
2153 the Global Historical Climatology Network was be-  
2154 ing first compiled. But, with modern improvements  
2155 in the globalization of communication, data archiving  
2156 and data sharing such a project is probably now  
2157 feasible.

2158 It would also be desirable if repeat station surveys  
2159 could be carried out every few years, as Leroy, 2010  
2160 recommends[31]:

*The rating of each site should be reviewed  
2161 periodically as environmental circumstances  
2162 can change over a period of time. A sys-  
2163 tematic yearly visual check is recommended:  
2164 if some aspects of the environment have  
2165 changed, a new classification process is nec-  
2166 essary. 2167*

*A complete update of the site classes should  
2168 be done at least every 5 years. 2169*

2170 A major difficulty in resolving the issue of how  
2171 station exposure have biased long term temperature  
2172 trends is that the Surface Stations project only pro-  
2173 vides us with information about the *current* station  
2174 quality. Repeat surveys every few years, as suggested  
2175 above, would reduce this issue, going forward. But,  
2176 if we want to continue to use station data for earlier  
2177 periods, we should probably attempt to collect more  
2178 information on the quality of these stations for earlier  
2179 times.

2180 Pielke et al., 2007b[10] note that for the early and  
2181 middle part of the 20th century, COOP observers  
2182 were encouraged to provide hand-drawn schematics  
2183 illustrating the site exposure associated with their  
2184 stations. Although these schematics were generally  
2185 not as comprehensive as the descriptions compiled  
2186 by the Surface Stations project, and station moves  
2187 were not always documented[10], it may be possi-  
2188 ble that partial estimates of the historical station  
2189 exposures can be reconstructed from these schemat-

---

<sup>4</sup>See the “surfacestation” posts on the Tallbloke’s  
Talkshop blog, the following posts on Roger Pielke Sr.’s  
blog: 2011/08/11, 2011/08/16, 2011/08/17, 2011/08/23,  
2011/08/29, 2011/09/08, 2011/09/16,2011/09/27 and  
2011/09/28, or these Watts Up With That posts on the  
Sydney, Australia station: 2008/07/02 and 2013/01/14.

2190 ics. In the 1980s, these schematics were replaced  
2191 with more cryptic, and less informative, shorthand  
2192 notation so that they could be entered into simple  
2193 computer forms. Nonetheless, when combined with  
2194 the earlier schematics and photographic evidence of  
2195 the current site exposure, it may be possible to re-  
2196 construct reasonable approximations of the site ex-  
2197 posures during that period.

2198 We saw in Section 4.5 that the step-change homogen-  
2199 ization algorithm currently used by the National  
2200 Climatic Data Center is unable to correctly deal with  
2201 the siting biases in the U.S. Historical Climatology  
2202 Network. In Ref. [13], we also showed that this algo-  
2203 rithm breaks down when a large number of stations  
2204 are affected by urbanization bias (or any non-climatic  
2205 *trend* biases). All automated statistical homogeniza-  
2206 tion algorithms are by their nature imperfect, since  
2207 they result in a mixture of true positives, true neg-  
2208 atives, false positives and false negatives. However,  
2209 while the Menne & Williams, 2009 algorithm appears  
2210 to perform relatively well on synthetic temperature  
2211 records[24, 73], it appears to be seriously problematic  
2212 when dealing with the extensive non-climatic biases  
2213 present in the U.S. Historical Climatology Network,  
2214 and presumably also the Global Historical Climatol-  
2215 ogy Network.

2216 We appreciate that it is important to remove these  
2217 non-climatic biases, as well as others, e.g., station  
2218 moves, time-of-observation biases or changes in in-  
2219 strumentation. A considerable amount of work has  
2220 already been carried out in assessing these step-  
2221 change homogenization techniques, e.g., see Refs. [69,  
2222 73, 74]. However, the current step-change homoge-  
2223 nization approaches still seem to introduce at least  
2224 as many problems as they remove. For this reason, it  
2225 seems unwise to apply these approaches to the tem-  
2226 perature records in an automated manner, as is cur-  
2227 rently done.

2228 We suggest that, if researchers want to continue  
2229 to use these automated homogenization approaches,  
2230 they should only be used for *identifying* potential  
2231 step changes (and possibly estimating their magni-  
2232 tudes), and flagging them for the researcher to man-  
2233 ually check, on a case-by-case basis. Records of the  
2234 researcher's justifications for accepting/rejecting in-  
2235 dividual adjustments should be kept, and provided  
2236 to users of the data, who want to make their own  
2237 assessments. This would require more work on the  
2238 part of the researcher, e.g., a typical station record  
2239 may contain about 6 such flags per 100 years[27], and  
2240 regional station networks may contain hundreds or

2241 thousands of station records. But, from our discus-  
2242 sion in Section 4.5, it seems manual checks are still  
2243 necessary.

2244 Rather than focusing on developing automated sta-  
2245 tistical methods for identifying undocumented non-  
2246 climatic biases, maybe we should first attempt to col-  
2247 lect (and digitize, where necessary) as much informa-  
2248 tion on individual station histories as possible. When  
2249 we have a more complete knowledge of the potential  
2250 non-climatic biases for which documentation is avail-  
2251 able, and have developed a more accurate knowledge  
2252 of the impact of the documented biases, we may be  
2253 in a better position to assess how to best treat the  
2254 undocumented biases.

2255 While station histories were collected during the  
2256 compilation of the U.S. Historical Climatology Net-  
2257 work, the same was *not* done for the Global His-  
2258 torical Climatology Network. When Peterson et al.,  
2259 1997 were compiling the Global Historical Climatol-  
2260 ogy Network[60], they were mainly combining pre-  
2261 compiled temperature (and precipitation) datasets.  
2262 These datasets often just contained mean monthly  
2263 temperature (and/or precipitation) values, and so Pe-  
2264 terson et al. decided not to consider detailed station  
2265 histories and metadata.

2266 Although it is true that standards differ between  
2267 meteorological organisations, and the level of doc-  
2268 umentation may vary from station to station, it is  
2269 likely that there is a considerable amount of use-  
2270 ful station documentation out there, which could be  
2271 useful. We suggest that, if the station observers  
2272 and/or meteorological organisations associated with  
2273 given stations are contacted directly, relatively de-  
2274 tailed station histories may be available. If a global  
2275 equivalent of the Surface Stations project were to be  
2276 organised the collection of these details could be car-  
2277 ried out simultaneously.

2278 Better station histories and documentation could  
2279 allow researchers to dramatically increase the reli-  
2280 ability of the global climate network. For instance, if  
2281 time-of-observation biases are as serious a problem  
2282 for U.S. stations as is claimed[25, 26], then it is likely  
2283 that similar problems exist for other networks in the  
2284 Global Historical Climatology Network. Station his-  
2285 tories could provide information on the observation  
2286 times and averaging methods used at different sta-  
2287 tions, which could help researchers assess these prob-  
2288 lems.

2289 As we discuss in Refs. [11–13], many of the long  
2290 station records in the Global Historical Climatology  
2291 Network contain a large number of unexplained data

gaps. These also are frequently associated with large temperature changes. The station observers, or station histories may be able to explain those gaps, and also whether any temperature changes between the gaps are likely to be climatic or not. In some cases, it may be possible to update the temperature records, if extra data is available. We note that it may be useful to also collect daily measurements too, if possible, as daily measurements often provide more insight than monthly measurements[75].

In terms of how such a global project could be organized, we note that, in the U.S. alone, approximately \$2.5 billion/year is currently allocated for “federal research on global change and climate change” through the United States Global Change Research Program. In such terms, a relatively small outlay of funding into systematic methodical on-site assessments and updates of station records and station histories would drastically improve the reliability of our climate records. Even in the absence of direct funding, the success of the Surface Stations project as well as other climate-related volunteer projects, such as the Old Weather and Climate Prediction projects, suggests that there could be sufficient interest in “citizen science” research[76] for a global extension of the original Surface Stations study of U.S. stations. If this approach was taken, it would be useful if the project received explicit approval and cooperation from the various international meteorological agencies whose stations would be considered. The backing of the World Meteorological Organization would probably help such a project.

Finally, we note that government-funded, high quality, carefully located and sited, modern climate networks like the National Climatic Data Center’s U.S. Climate Reference Network (USCRN), if properly managed, have the potential to offer future researchers climate records which would be unaffected by many of the main non-climatic biases plaguing current studies, e.g., urbanization bias, siting bias, time-of-observation bias and possibly instrumental bias[77, 78]. The U.S. Climate Reference Network already has a decade of data[77], and if similar networks were set-up and maintained in enough other countries, then it may be possible in the future to compile enough stations for global climate studies. We recognise that this would not be of immediate benefit to current researchers, but are reminded of the proverb: “The best time to plant a tree is twenty years ago. The second best time is now.”

## 6 Conclusions

Recent surveys of the current thermometer station exposure of the weather stations in the U.S. Historical Climatology Network have identified that about 70% of stations are poorly or badly sited[1, 2]. As a result, it is likely that many of the thermometer records currently used for estimating climate trends (both regional U.S. temperature trends and global temperature trends) have been biased by changes in the localised micro-climate in the immediate vicinity of the thermometer stations, rather than representing purely climatic trends. However, attempts to quantify these siting biases have until now been somewhat contradictory and controversial[1, 2, 14–17].

In this study, we estimate that siting biases have artificially increased the mean temperature trends of the *Unadjusted* station records by about 32%, with the subset of the worst-sited stations being increased by about 56%. When time-of-observation adjustments were applied to the station records, this led to a mean increase in temperature trends of about 39%, as has been previously noted[26, 63]. For this reason, although the siting biases remained similar in magnitude, the *relative* fraction of the temperature trends in the *Time-of-observation adjusted* station records that was due to siting biases was reduced - an artificial warming of about 18% for all rated stations, and about 49% for the worst-sited stations. This still represents a substantial bias.

When the *Time-of-observation adjusted* station records were also homogenized with step-change adjustments using the Menne & Williams, 2009 algorithm[27], the differences between all subsets of stations were substantially reduced. There was no longer much difference in temperature trends between the good quality stations and the bad quality stations. Previously it had been argued that this was because the step-change adjustments had removed most of the siting biases in the station records[2]. However, our analysis shows that the Menne & Williams, 2009 adjustments actually lead to considerable blending/mixing of siting biases between stations, rather than their removal. As a result, the current *Time-of-observation and step-change adjusted* version of the U.S. Historical Climatology Network is unreliable, and its use should be discontinued.

It seems likely that similar siting biases also exist in the Global Historical Climatology Network, which has been frequently used for estimating global temperature trends. This suggests that current estimates

of the amount of “global warming” since the 19th century (e.g., Ref. [23]) have been significantly over-estimated.

Siting biases are not the only non-climatic biases affecting the thermometer stations being used for studying climatic trends. For instance, in Refs. [11–13] we discuss the widespread problem of urbanization bias in current thermometer records. We agree that thermometer station records could be an invaluable resource for evaluating global and regional temperature trends. However, many of the problems identified by earlier researchers, such as Mitchell, 1953[79] have still not been adequately resolved. In Section 5, we offer a number of recommendations which we believe would lead to more reliable climate records.

## Acknowledgements

No funding was received for this research.

We are grateful to Anthony Watts and the other volunteers who worked on the Surface Stations project for generating and archiving the station exposure ratings for the U.S. Historical Climatology Network.

We would like to thank Don Zieman for some useful comments.

## References

[1] A. Watts. “Is the U.S. surface temperature record reliable?” The Heartland Institute. Chicago, IL. 2009. URL: <http://surfacestations.org/>.

[2] M. J. Menne, C. N. Jr. Williams, and M. A. Palecki. “On the reliability of the U.S. surface temperature record”. *J. Geophys. Res.* 115 (2010), p. D11108. doi: 10.1029/2009JD013094.

[3] D. A. Robinson. “Gathering climatic data of the highest quality” (1997), pp. 305–308.

[4] K. I. Scott, J. R. Simpson, and E. G. McPherson. “Effects of tree cover on parking lot microclimate and vehicle emissions”. *J. Arboriculture* 25 (1999), pp. 129–142. URL: <http://joa.isa-arbor.com/>.

[5] H. Yilmaz et al. “Determination of temperature differences between asphalt concrete, soil and grass surfaces of the city of Erzurum, Turkey”. *Atmósfera* 21 (2008), pp. 135–146. URL: [http://www.scielo.org.mx/scielo.php?pid=S0187-62362008000200002&script=sci\\_arttext](http://www.scielo.org.mx/scielo.php?pid=S0187-62362008000200002&script=sci_arttext).

[6] R. Mahmood, S. A. Foster, and D. Logan. “The Geoprofile metadata, exposure of instruments, and measurement bias in climatic record revisited”. *Int. J. Clim.* 26 (2006), pp. 1091–1124. doi: 10.1002/joc.1298.

[7] G. Manley. “Central England temperatures: monthly means 1659–1973”. *Quart. J. R. Met. Soc.* 100 (1974), pp. 389–405. doi: 10.1256/smsqj.42510.

[8] C. A. Davey and R. A. Sr. Pielke. “Microclimate exposures of surface-based weather stations: Implication for the assessment of long-term temperature trends”. *Bull. Amer. Meteor. Soc.* 86 (2005), pp. 497–504. doi: 10.1175/BAMS-86-4-497.

[9] R. A. Sr. Pielke et al. “Unresolved issues with the assessment of multidecadal global land-surface temperature trends”. *J. Geophys. Res.* 112 (2007), D24S08. doi: 10.1029/2006JD008229.

[10] R. A. Sr. Pielke et al. “Documentation of uncertainties and biases associated with surface temperature measurement sites for climate change assessment”. *Bull. Amer. Meteor. Soc.* 88 (2007), pp. 913–928. doi: 10.1175/BAMS-88-6-913.

[11] R. Connolly and M. Connolly. “Urbanization bias I. Is it a negligible problem for global temperature estimates?” 28 (*Clim. Sci.*). Ver. 0.1 (non peer reviewed draft). 2014. URL: <http://oprj.net/articles/climate-science/28>.

[12] R. Connolly and M. Connolly. “Urbanization bias II. An assessment of the NASA GISS urbanization adjustment method”. 31 (*Clim. Sci.*). Ver. 0.1 (non peer reviewed draft). 2014. URL: <http://oprj.net/articles/climate-science/31>.

[13] R. Connolly and M. Connolly. “Urbanization bias III. Estimating the extent of bias in the Historical Climatology Network datasets”. 34 (*Clim. Sci.*). Ver. 0.1 (non peer reviewed draft). 2014. URL: <http://oprj.net/articles/climate-science/34>.

[14] S. Fall et al. “Analysis of the impacts of station exposure on the U.S. Historical Climatology Network temperatures and temperature trends”. *J. Geophys. Res.* 116 (2011), p. D14120. doi: 10.1029/2010JD015146.

[15] R. A. Muller et al. “Earth atmospheric land surface temperature and station quality in the contiguous United States”. *Geoinfor. Geostat.* 1 (2013). doi: 10.4172/2327-4581.1000107.

[16] C. J. Martinez, J. J. Maleski, and M. F. Miller. “Trends in precipitation and temperature in Florida, USA”. *J. Hydrol.* 452–453 (2012), pp. 259–281. doi: 10.1016/j.jhydrol.2012.05.066.

[17] A. Watts et al. “An area and distance weighted analysis of the impacts of station exposure on the U.S. Historical Climatology Network temperatures and temperature trends”. *In preparation* (2012). URL: <http://wattsupwiththat.com/2012/07/29/press-release-2/>.

[18] D. A. Robinson. “The United States Cooperative climate-observing systems: Reflections and recommendations”. *Bull. Amer. Meteor. Soc.* 71 (1990), pp. 826–831. doi: 10.1175/1520-0477(1990)071<0826:TUSCCO>2.0.CO;2.

[19] R. S. Vose et al. “Comments on ‘Microclimate exposures of surface-based weather stations: Implication for the assessment of long-term temperature trends’”. *Bull. Amer. Meteor. Soc.* 86 (2005), pp. 504–506. doi: 10.1175/BAMS-86-4-504.

[20] T. C. Peterson. “Examination of potential biases in air temperature caused by poor station locations”. *Bull. Amer. Meteor. Soc.* 87 (2006), pp. 1073–1089. doi: 10.1175/BAMS-87-8-1073.

[21] NOAA/NESDIS. “Climate Reference Network (CRN) Site Information Handbook”. NOAA/NESDIS CRN Ser. X030, CRN Rep. NOAA-CRN/OSD-2002-0002R0UD0, National Climatic Data Center, Asheville, N.C., U.S.A. 2002. URL: <http://www1.ncdc.noaa.gov/pub/data/uscrn/documentation/program/X030FullDocumentD0.pdf>.

[22] M. Leroy. “Classification d’un site”. Note Technique #35. Direction des Systèmes d’Observation, Météo-France, Trappes, France. 1999. URL: [http://www.ccrom.org/ccrom/IMG/pdf/note\\_technique35-2.pdf](http://www.ccrom.org/ccrom/IMG/pdf/note_technique35-2.pdf).

[23] J. Hansen et al. “A closer look at United States and global surface temperature change”. *J. Geophys. Res.* D 106 (2001), pp. 23947–23963. doi: 10.1029/2001JD000354.

[24] M. J. Menne, C. N. Jr. Williams, and R. S. Vose. “The U.S. Historical Climatology Network monthly temperature data, version 2”. *Bull. Amer. Meteor. Soc.* 90 (2009), pp. 993–1007. doi: 10.1175/2008BAMS2613.1.

[25] T. R. Karl et al. “A model to estimate time of observation bias associated with monthly mean maximum, minimum and mean temperatures for the United States”. *J. Clim. Appl. Meteor.* 25 (1986), pp. 145–160. doi: 10.1175/1520-0450(1986)025<0145:AMTETT>2.0.CO;2.

[26] R. S. Vose et al. “An evaluation of the time of observation bias adjustment in the U.S. Historical Climatology Network”. *Geophys. Res. Lett.* 30 (2003), p. 2046. doi: 10.1029/2003GL018111.



2516 [27] M. J. Menne and C. N. Jr. Williams. "Homogenization of  
2517 temperature series via pairwise comparisons". *J. Clim.* 22 (2009),  
2518 pp. 1700–1717. DOI: [10.1175/2008JCLI2263.1](https://doi.org/10.1175/2008JCLI2263.1).

2519 [28] T. R. Karl et al. "A new perspective on recent global warm-  
2520 ing: Asymmetric trends of daily maximum and minimum temper-  
2521 ature". *Bull. Amer. Meteor. Soc.* 74 (1993), pp. 1007–1023. DOI:  
2522 [10.1175/1520-0477\(1993\)074<1007:ANPORG>2.0.CO;2](https://doi.org/10.1175/1520-0477(1993)074<1007:ANPORG>2.0.CO;2).

2523 [29] K. Gallo, D. R. Easterling, and T. C. Peterson. "The influ-  
2524 ence of land use/land cover on climatological values of the diur-  
2525 nal temperature range". *J. Clim.* 9 (1996), pp. 2941–2944. DOI:  
2526 [10.1175/1520-0442\(1996\)009<2941:TIOJUC>2.0.CO;2](https://doi.org/10.1175/1520-0442(1996)009<2941:TIOJUC>2.0.CO;2).

2527 [30] R. S. Vose, D. R. Easterling, and B. Gleason. "Maximum and  
2528 minimum temperature trends for the globe: An update through  
2529 2004". *Geophys. Res. Lett.* 32 (2005), p. L23822. DOI: [10.1029/  
2530 2005GL024379](https://doi.org/10.1029/2005GL024379).

2531 [31] M. Leroy. "Siting classification for surface observing stations  
2532 on land". JMA/WMO Workshop on Quality Management in Sur-  
2533 face, Climate and Upper-air Observations in RA II (Asia). Tokyo,  
2534 Japan 27-30 July 2010. 2010. URL: [http://www.jma.go.jp/jma/en/  
2535 Activities/qmws\\_2010/CountryReport/CS202\\_Leroy.pdf](http://www.jma.go.jp/jma/en/Activities/qmws_2010/CountryReport/CS202_Leroy.pdf).

2536 [32] M. Chenoweth. "A possible discontinuity in the U.S. histor-  
2537 ical temperature record". *J. Clim.* 5 (1992), pp. 1172–1179. DOI:  
2538 [10.1175/1520-0442\(1992\)005<1172:APDITU>2.0.CO;2](https://doi.org/10.1175/1520-0442(1992)005<1172:APDITU>2.0.CO;2).

2539 [33] M. Chenoweth. "Nonstandard thermometer exposures at U.S.  
2540 Cooperative weather stations during the late nineteenth century".  
2541 *J. Clim.* 6 (1993), pp. 1787–1797. DOI: [10.1175/1520-0442\(1993\)  
2542 006<1787:NTEAUC>2.0.CO;2](https://doi.org/10.1175/1520-0442(1993)006<1787:NTEAUC>2.0.CO;2).

2543 [34] D. E. Parker. "Effects of changing exposure of thermometers  
2544 at land stations". *Int. J. Climatol.* 14 (1994), pp. 1–31. DOI: [10.  
2545 1002/joc.3370140102](https://doi.org/10.1002/joc.3370140102).

2546 [35] W. S. Hughes. "Comment on D. E. Parker, 'Effects of chang-  
2547 ing exposure of thermometers at land stations'". *Int. J. Climatol.*  
2548 15 (1995), pp. 231–234. DOI: [10.1002/joc.3370150208](https://doi.org/10.1002/joc.3370150208).

2549 [36] D. E. Parker. "Effects of changing exposure of thermometers  
2550 at land stations: A reply". *Int. J. Climatol.* 15 (1995), p. 235. DOI:  
2551 [10.1002/joc.3370150209](https://doi.org/10.1002/joc.3370150209).

2552 [37] N. Nicholls et al. "Historical thermometer exposures in Aus-  
2553 tralia". *Int. J. Climatol.* 16 (1996), pp. 705–710. DOI: [10.1002/  
2554 \(SICI\)1097-0088\(199606\)16:6<705::AID-JOC30>3.0.CO;2-S](https://doi.org/10.1002/(SICI)1097-0088(199606)16:6<705::AID-JOC30>3.0.CO;2-S).

2555 [38] W. S. Hughes. "Comment on 'Historical thermometer expo-  
2556 sures in Australia' by N. Nicholls et al." *Int. J. Climatol.* 17  
2557 (1997), pp. 197–199. DOI: [10.1002/\(SICI\)1097-0088\(199702\)17:  
2558 2<197::AID-JOC113>3.0.CO;2-R](https://doi.org/10.1002/(SICI)1097-0088(199702)17:2<197::AID-JOC113>3.0.CO;2-R).

2559 [39] B. Trewin. "Another look at Australia's record high temper-  
2560 ature". *Aust. Met. Mag.* 46 (1997), pp. 251–256. URL: [http://www.  
2561 bom.gov.au/amm/docs/1997/trewin.pdf](http://www.bom.gov.au/amm/docs/1997/trewin.pdf).

2562 [40] P. Ø. Nordli et al. "The effect of radiation screens on Nordic  
2563 time series of mean temperature". *Int. J. Climatol.* 17 (1997),  
2564 pp. 1667–1681. DOI: [10.1002/\(SICI\)1097-0088\(199712\)17:15<1667::  
2565 AID-JOC221>3.0.CO;2-D](https://doi.org/10.1002/(SICI)1097-0088(199712)17:15<1667::AID-JOC221>3.0.CO;2-D).

2566 [41] A. Moberg et al. "Were southern Swedish summer temper-  
2567 atures before 1860 as warm as measured?" *Int. J. Climatol.* 23  
2568 (2003), pp. 1495–1521. DOI: [10.1002/joc.945](https://doi.org/10.1002/joc.945).

2569 [42] D. Frank et al. "Warmer early instrumental measurements  
2570 versus colder reconstructed temperatures: shooting at a moving  
2571 target". *Quat. Sci. Rev.* 26 (2007), pp. 3298–3310. DOI: [10.1016/  
2572 j.quascirev.2007.08.002](https://doi.org/10.1016/j.quascirev.2007.08.002).

2573 [43] R. Böhm et al. "The early instrumental warm-bias: a solution  
2574 for long central European temperature series 1760–2007". *Clim.  
2575 Change* 101 (2010), pp. 41–67. DOI: [10.1007/s10584-009-9649-4](https://doi.org/10.1007/s10584-009-9649-4).

2576 [44] M. Brunet, J. Asin, J. Sigró, et al. "The minimization of the  
2577 screen bias from ancient Western Mediterranean air temperature  
2578 records: an exploratory statistical analysis". *Int. J. Climatol.* 31  
2579 (2011), pp. 1879–1895. DOI: [10.1002/joc.2192](https://doi.org/10.1002/joc.2192).

2580 [45] K. G. Hubbard et al. "Air temperature comparison between  
2581 the MMTS and the USCRN temperature systems". *J. Atmos.  
2582 Oceanic Tech.* 21 (2004), pp. 1590–1597. DOI: [10.1175/1520-  
2583 0426\(2004\)021<1590:ATCBTM>2.0.CO;2](https://doi.org/10.1175/1520-0426(2004)021<1590:ATCBTM>2.0.CO;2).

2584 [46] K. G. Hubbard and X. Lin. "Reexamination of instrument  
2585 change effects in the U.S. Historical Climatology Network". *Geo-  
2586 phys. Res. Lett.* 33 (2006), p. L15710. DOI: [10.1029/2006GL027069](https://doi.org/10.1029/2006GL027069).

2587 [47] R. G. Quayle et al. "Effects of recent thermometer changes  
2588 in the Cooperative Station Network". *Bull. Amer. Meteor. Soc.*  
2589 72 (1991), pp. 1718–1723. DOI: [10.1175/1520-0477\(1991\)072<1718:  
2590 EORTC1>2.0.CO;2](https://doi.org/10.1175/1520-0477(1991)072<1718:EORTC1>2.0.CO;2).

2591 [48] R. Connolly and M. Connolly. "Global temperature changes  
2592 of the last millennium". 16 (Clim. Sci.). Ver. 0.1 (non peer reviewed  
2593 draft). 2014. URL: <http://oprj.net/articles/climate-science/16>.

2594 [49] R. Wilson et al. "A matter of divergence: Tracking recent  
2595 warming at hemispheric scales using tree ring data". *J. Geophys.  
2596 Res.* 112 (2007), p. D17103. DOI: [10.1029/2006JD008318](https://doi.org/10.1029/2006JD008318).

2597 [50] R. D'Arrigo et al. "On the 'divergence problem' in northern  
2598 forests: A review of the tree-ring evidence and possible causes".  
2599 *Glob. Planet. Change* 60 (2008), pp. 289–305. DOI: [10.1016/j.  
2600 gloplacha.2007.03.004](https://doi.org/10.1016/j.gloplacha.2007.03.004).

2601 [51] B. Walsh. "2012 was the hottest year in U.S. history. And yes  
2602 - it's climate change". Time magazine, Ecocentric blog post on 8  
2603 January 2013. Accessed: 2013-03-12. (Archived by WebCite® at  
2604 <http://www.webcitation.org/6F4ayyon3>). URL: [http://science.  
2605 time.com/2013/01/08/2012-was-the-hottest-year-in-u-s-  
2606 history-and-yes-its-climate-change/](http://science.time.com/2013/01/08/2012-was-the-hottest-year-in-u-s-history-and-yes-its-climate-change/).

2607 [52] J. Gillis. "Not even close: 2012 was hottest ever in U.S." New  
2608 York Times on 8 January 2013. Accessed: 2013-03-12. (Archived  
2609 by WebCite® at <http://www.webcitation.org/6F4adRVHa>). URL:  
2610 [http://www.nytimes.com/2013/01/09/science/earth/2012-was-  
2611 hottest-year-ever-in-us.html?\\_r=0](http://www.nytimes.com/2013/01/09/science/earth/2012-was-hottest-year-ever-in-us.html?_r=0).

2612 [53] K. Than. "2012: hottest year on record for continental U.S." National Geographic News. 9 January 2013. Accessed: 2013-03-12. (Archived by WebCite® at [http://www.webcitation.org/  
2613 6F4blW8i2](http://www.webcitation.org/6F4blW8i2)). URL: [http://news.nationalgeographic.com/news/2013/  
2614 01/130109-warmest-year-record-2012-global-warming-science-  
2615 environment-united-states/](http://news.nationalgeographic.com/news/2013/01/130109-warmest-year-record-2012-global-warming-science-environment-united-states/).

2616 [54] R. L. Baumhardt. "Dust Bowl Era". In: *Encyclopedia of Water Science*. Marcel Dekker, Inc. New York., 2003, pp. 187–191. 2618 2619

2620 [55] A. Watts. "Rumours of my death have been greatly exagger-  
2621 ated" (2010). URL: [http://wattsupwiththat.com/2010/01/27/  
2622 rumours-of-my-death-have-been-greatly-exaggerated/](http://wattsupwiththat.com/2010/01/27/rumours-of-my-death-have-been-greatly-exaggerated/).

2623 [56] P. D. Jones et al. "Northern Hemisphere surface air temper-  
2624 ature variations: 1851–1984". *J. Clim. Appl. Meteor.* 25 (1986),  
2625 pp. 161–179. DOI: [10.1175/1520-0450\(1986\)025<0161:NHSATV>2.0.  
2626 CO;2](https://doi.org/10.1175/1520-0450(1986)025<0161:NHSATV>2.0.CO;2).

2627 [57] C. Wunsch. "The interpretation of short climate records, with  
2628 comments on the North Atlantic and Southern Oscillations". *Bull.  
2629 Amer. Meteor. Soc.* 80 (1999), pp. 245–255. DOI: [10.1175/1520-  
2630 0477\(1999\)080<0245:TIOSCR>2.0.CO;2](https://doi.org/10.1175/1520-0477(1999)080<0245:TIOSCR>2.0.CO;2).

2631 [58] D. B. Percival and D. A. Rothrock. "Eyeballing" trends  
2632 in climate time series: A cautionary note". *J. Clim.* 18 (2005),  
2633 pp. 886–891. DOI: [10.1175/JCLI-3300.1](https://doi.org/10.1175/JCLI-3300.1).

2634 [59] M. Shermer. "Turn me on, dead man". *Sci. Am.* 292 (May)  
2635 (2005), p. 37. URL: [http://www.scientificamerican.com/article.  
2636 cfm?id=turn-me-on-dead-man](http://www.scientificamerican.com/article.cfm?id=turn-me-on-dead-man).

2637 [60] T. C. Peterson and R. S. Vose. "An overview of the Global  
2638 Historical Climatology Network temperature database". *Bull.  
2639 Amer. Meteor. Soc.* 78 (1997), pp. 2837–2849. DOI: [10.1175/1520-  
2640 0477\(1997\)078<2837:A00TGH>2.0.CO;2](https://doi.org/10.1175/1520-0477(1997)078<2837:A00TGH>2.0.CO;2).

2641 [61] Columbia University Center for International Earth Science  
2642 Information Network (CIESIN) et al. "Global Rural-Urban Map-  
2643 ping Project, Version 1 (GRUMPv1): Urban Extents Grid". Date  
2644 of download: 6th January 2013. 2011. URL: [http://sedac.ciesin.  
2645 columbia.edu/data/dataset/grump-v1-urban-extents](http://sedac.ciesin.columbia.edu/data/dataset/grump-v1-urban-extents).

2646 [62] Z. Hausfather et al. "Quantifying the effect of urbanization  
2647 on U.S. Historical Climatology Network temperature records". *J.  
2648 Geophys. Res.* in press (2013). DOI: [10.1029/2012JD018509](https://doi.org/10.1029/2012JD018509).

2649 [63] R. C. Jr. Balling and C. D. Idso. "Analysis of adjustments  
2650 to the United States Historical Climatology Network (USHCN)  
2651 temperature database". *Geophys. Res. Lett.* 29 (2002), p. 1387.  
2652 DOI: [10.1029/2002GL014825](https://doi.org/10.1029/2002GL014825).

- 2653 [64] D. A. Portman. “Identifying and correcting urban bias in  
2654 regional time series: Surface temperature in China’s northern  
2655 plains”. *J. Clim.* 6 (1993), pp. 2298–2308. DOI: [10.1175/1520-0442\(1993\)006<2298:IACUBI>2.0.CO;2](https://doi.org/10.1175/1520-0442(1993)006<2298:IACUBI>2.0.CO;2).
- 2657 [65] L. A. Vincent et al. “Bias in minimum temperature intro-  
2658 duced by a redefinition of the climatological day at the Canadian  
2659 synoptic stations”. *J. Appl. Meteor. Clim.* 48 (2009), pp. 2160–  
2660 2168. DOI: [10.1175/2009JAMC2191.1](https://doi.org/10.1175/2009JAMC2191.1).
- 2661 [66] F. Fujibe et al. “Long-term changes of temperature extremes  
2662 and day-to-day variability in Japan”. *Pap. Met. Geophys.* 58  
2663 (2007), pp. 63–72. DOI: [10.2467/mripapers.58.63](https://doi.org/10.2467/mripapers.58.63).
- 2664 [67] B. Trewin. “A daily homogenized temperature data set for  
2665 Australia”. *Int. J. Climatol.* In press (2012). DOI: [10.1002/joc.3530](https://doi.org/10.1002/joc.3530).
- 2667 [68] T. R. Karl and C. N. Jr. Williams. “An approach to adjusting  
2668 climatological time series for discontinuous inhomogeneities”. *J.*  
2669 *Clim. Appl. Meteor.* 26 (1987), pp. 1744–1763. DOI: [10.1175/1520-0450\(1987\)026<1744:AATACT>2.0.CO;2](https://doi.org/10.1175/1520-0450(1987)026<1744:AATACT>2.0.CO;2).
- 2671 [69] A. T. DeGaetano. “Attributes of several methods for detect-  
2672 ing discontinuities in mean temperature series”. *J. Clim.* 19 (2006),  
2673 pp. 838–853. DOI: [10.1175/JCLI3662.1](https://doi.org/10.1175/JCLI3662.1).
- 2674 [70] T. R. Karl, H. F. Diaz, and G. Kukla. “Urbanization: Its  
2675 detection and effect in the United States climate record”. *J. Clim.*  
2676 1 (1988), pp. 1099–1123. DOI: [10.1175/1520-0442\(1988\)001<1099:UIDAEI>2.0.CO;2](https://doi.org/10.1175/1520-0442(1988)001<1099:UIDAEI>2.0.CO;2).
- 2678 [71] J. R. Lanzante. “Resistant, robust and non-parametric tech-  
2679 niques for the analysis of climate data: Theory and examples, in-  
2680 cluding applications to historical radiosonde station data”. *Int.*  
2681 *J. Clim.* 16 (1996), pp. 1197–1226. DOI: [10.1002/\(SICI\)1097-0088\(199611\)16:11<1197::AID-JOC89>3.0.CO;2-L](https://doi.org/10.1002/(SICI)1097-0088(199611)16:11<1197::AID-JOC89>3.0.CO;2-L).
- 2683 [72] K. E. Runnalls and T. R. Oke. “A technique to detect micro-  
2684 climatic inhomogeneities in historical records of screen-level air  
2685 temperature”. *J. Clim.* 19 (2006), pp. 959–978. DOI: [10.1175/JCLI3663.1](https://doi.org/10.1175/JCLI3663.1).
- 2687 [73] V. K. C. Venema et al. “Benchmarking homogenization algo-  
2688 rithms for monthly data”. *Clim. Past* 8 (2012), pp. 89–115. DOI:  
2689 [10.5194/cp-8-89-2012](https://doi.org/10.5194/cp-8-89-2012).
- 2690 [74] J. Reeves et al. “A review and comparison of changepoint  
2691 detection techniques for climate data”. *J. Appl. Meteor. Clim.* 46  
2692 (2007), pp. 900–915. DOI: [10.1175/JAM2493.1](https://doi.org/10.1175/JAM2493.1).
- 2693 [75] A. T. DeGaetano. “A method to infer observation time  
2694 based on day-to-day temperature variations”. *J. Clim.* 12 (1999),  
2695 pp. 3443–3456. DOI: [10.1175/1520-0442\(1999\)012<3443:AMTIOT>2.0.CO;2](https://doi.org/10.1175/1520-0442(1999)012<3443:AMTIOT>2.0.CO;2).
- 2697 [76] C. Franzoni and H. Sauermann. “Crowd science: The organi-  
2698 zation of scientific research in open collaborative projects”. *SSRN*  
2699 (November 14, 2012). DOI: [10.2139/ssrn.2167538](https://doi.org/10.2139/ssrn.2167538). URL: <http://ssrn.com/abstract=2167538>.
- 2701 [77] H. J. Diamond et al. “U.S. Climate Reference Network after  
2702 one decade on operations: Status and assessment”. *Bull. Amer.*  
2703 *Meteor. Soc.* in press (2012). DOI: [10.1175/BAMS-D-12-00170](https://doi.org/10.1175/BAMS-D-12-00170).
- 2704 [78] K. P. Gallo. “Evaluation of temperature differences for paired  
2705 stations of the U. S. Climate Reference Network”. *J. Clim.* 18  
2706 (2005), pp. 1629–1636. DOI: [10.1175/JCLI3358.1](https://doi.org/10.1175/JCLI3358.1).
- 2707 [79] J. M. Mitchell. “On the causes of instrumentally observed  
2708 secular temperature trends”. *J. Meteor.* 10 (1953), pp. 244–261.  
2709 DOI: [10.1175/1520-0469\(1953\)010<0244:OTC0IO>2.0.CO;2](https://doi.org/10.1175/1520-0469(1953)010<0244:OTC0IO>2.0.CO;2).