# Testing E-OBS European high-resolution gridded dataset of daily precipitation and surface temperature

4 Nynke Hofstra[1,2,*], Malcolm Haylock[3], Mark New[1], Phil D. Jones[3]

6 [1] School of Geography and the Environment, Oxford University Centre for the Environment,

7 South Parks Road, Oxford, OX2 8JU, England

8 [2] Now at: Environmental Systems Analysis Group, Wageningen University, P.O. Box 47, 6700

9 AA Wageningen, The Netherlands

10 [3] Climatic Research Unit, School of Environmental Sciences, University of East Anglia,

11 Norwich, NR4 7TJ, England

12 [*] Corresponding author: nynke.hofstra@wur.nl

## Abstract

15 Gridded datasets derived through interpolation of station data have a number of potential

16 inaccuracies and errors. These errors can be introduced either by the propagation of errors in the

17 station data into derived gridded data or by limitations in the ability of the interpolation method

18 to estimate grid values from the underlying station network. Recently, *Haylock et al* [2008]

19 reported on the development of a new high-resolution gridded dataset of daily climate over

20 Europe (termed E-OBS). E-OBS is based on the largest available pan-European dataset and the

21 interpolation methods used were chosen after careful evaluation of a number of alternatives, yet

22 the dataset will inevitably have errors and uncertainties. In this paper we assess the E-OBS

23 dataset with respect to: 1) homogeneity of the gridded data; 2) evaluation of inaccuracies arising

24  from available network density, through comparison with existing datasets that have been

25  developed with much denser station networks; and 3) the accuracy of the estimates of

26  interpolation uncertainty that are provided as part of E-OBS.

27

28  We find many inhomogeneities in the gridded data that are primarily caused by inhomogeneities

29  in the underlying station data.  In the comparison of existing data with E-OBS we find that while

30  correlations overall are high, relative differences in precipitation are large, and usually biased

31  towards lower values in E-OBS.  From the analysis of the interpolation uncertainties provided as

32  part of E-OBS, we conclude that the interpolation standard deviation provided with the data

33  significantly underestimates the true interpolation error when cross-validated using station data,

34  and therefore will similarly underestimate the interpolation error in the gridded E-OBS data.

35  While E-OBS represents a valuable new resource for climate research in Europe, users of the

36  data need to be aware of the limitations in the dataset and use the data appropriately.

37

## 1. Introduction

38

39  Gridded climate data derived from meteorological station measurements underpin a wide range

40  of applications and research in climate science, including evaluation of global and regional

41  climate models, the construction of bias-corrected climate change scenarios and driving many

42  applications in climate impacts assessments [*Haylock et al.*, 2008].  Increasingly, there has been

43  a need for gridded data at higher spatial and temporal resolutions, as the focus of climate change

44  research has shifted from global to regional and local scales.  Recently, *Haylock et al.* [2008]

45  described the development of the first high-resolution gridded dataset of daily climate over

46  Europe (termed E-OBS), as part of the EU funded ENSEMBLES project.  The dataset,

47  comprising daily mean, minimum and maximum temperature and precipitation, was constructed

48   through interpolation of the most complete collection of station data over wider Europe [*Klok*

49   *and Klein Tank*, 2008].  The data are available on four different RCM grids (0.25 and 0.5 degree

50   regular lat-lon and 0.22 and 0.44 degree rotated-pole) and cover the period 1950-2006.

51   Additionally, estimates of interpolation uncertainties are included as part of the dataset [*Haylock*

52   *et al.*, 2008].

53

54   Gridded datasets derived through interpolation of station data have a number of potential

55   inaccuracies and errors.  Errors in the underlying station data can be propagated into the gridded

56   data; typical sources of error include incorrect station location information, individual erroneous

57   values or non-climatic breaks (inhomogeneities) in the station time series.  A second source of

58   uncertainty relates to the ability of the interpolation method to estimate grid values from the

59   underlying station network.  In general, interpolation accuracy decreases as the network density

60   decreases, is less accurate for variables with more variable spatial characteristics (e.g.

61   precipitation) and degrades in areas of complex terrain (e.g. mountain areas).  While E-OBS is

62   based on the largest available pan-European dataset and the interpolation methods used were

63   chosen after careful evaluation of a number of alternatives [*Hofstra et al.*, 2008], the dataset will

64   inevitably have errors and uncertainties.

65

66   The aim of this paper is to assess the E-OBS dataset with respect to some of the potential errors

67   that may be present.  Users can then familiarise themselves with the strengths and weaknesses of

68   the data and use them responsibly.  We have chosen three features of E-OBS to analyse in this

69   paper: 1) homogeneity of the gridded data; 2) inaccuracies due to the underlying station network

70   density, though comparison with existing datasets that have been developed with much denser

71   station networks; and 3) the accuracy of the estimates of interpolation uncertainty that are

72   provided as part of E-OBS.

73

74    Long-term station data are often influenced by non-climatic factors, such as changes in station

75    location or environment, instruments and observing practices. These so-called inhomogeneities

76    can often lead to misinterpretations of the climate data analysed [*Peterson et al.*, 1998]. The

77    station data used for E-OBS are not fully homogenised. Individual station series may have been

78    homogenised by the original custodians of each series, but the series provided by partner

79    organisations have been used directly, meaning potentially inhomogeneous stations may be

80    contributing to the interpolated grids. As station density strongly influences the interpolation

81    [*Hofstra et al.*, 2008], E-OBS was constructed using many *potentially* inhomogeneous stations,

82    as their exclusion would degrade the station network density and hence accuracy of the

83    interpolation. In addition, several studies explain that, for area averages of relatively large areas,

84    inhomogeneities balance out during interpolation [*Dai et al.*, 1997; *New*, 1999; *Peterson et al.*,

85    1998]. However, that may not be the case for the E-OBS high-resolution grids. Therefore, the

86    first out of three topics tested is the homogeneity of the dataset.

87

88    The second topic is a comparison with other gridded datasets that have been developed with

89    much denser station networks. These datasets are available, in the case of precipitation, for long

90    periods for the UK and the Alps and for the period October 1999 – December 2000 for Europe as

91    a whole. For temperature, unfortunately, we have only been able to secure data for the UK.

92    Datasets developed with denser station networks are assumed to be a better approximation of the

93    true area-averages. So if the E-OBS gridded dataset produces grid area-averages that are close to

94    those calculated from the higher quality grids, the E-OBS dataset can be deemed to be a

95    reasonable representation of the true area-average gridded values.

96

97   Because of the inevitable interpolation uncertainties, the E-OBS dataset is provided with

98   information on the interpolation uncertainty for each grid box and each day [*Haylock et al.*,

99   2008].  E-OBS interpolation uncertainty is derived by combining the Bayesian standard error

100  estimates of the monthly climatology [*Hutchinson*, 1995] and the interpolation standard

101  deviation for daily anomalies [*Yamamoto*, 2000] (see section 5 for more detail).  Here we

102  concentrate on the interpolation standard error estimates, and evaluate the accuracy of the

103  estimates through cross-validation against station data.  This represents the first evaluation of the

104  *Yamamoto* [2000] standard error method, which has to date only been applied to geological data.

105

106  The remainder of the paper is structured as follows.  Section 2 provides a more detailed

107  description of the E-OBS dataset, including the underlying station data and the interpolation and

108  gridding methodology.  We then cover each of the three evaluations in turn: inhomogeneities

109  (Section 3), comparison against regional gridded datasets based on denser station networks

110  (Section 4) and evaluation of the interpolation standard error estimates (Section 5).  We conclude

111  with a summary of results and a discussion of the implications of our assessment for use of the

112  E-OBS dataset.

113

## 114   2. The E-OBS dataset

115  The E-OBS gridded dataset is derived through interpolation of the ECA&D (European Climate

116  Assessment and Data) station data described in *Klok and Klein Tank* [2008].  The station dataset

117  comprises a network of 2316 stations, with the highest station density in Ireland, the Netherlands

118  and Switzerland, and lowest density in Spain, Northern Africa, the Balkans and Northern

119  Scandinavia.  The number of stations used for the interpolation differs through time and by

120  variable.  The full period of record used for interpolation is 1950 – 2006 , but the period 1961 –

121 1990 has the highest density. At any particular time, there are more precipitation than

122 temperature stations. Inhomogeneities in the station time-series have been flagged, but

123 potentially inhomogeneous stations are used for the interpolation, for reasons noted above.

124

125 The E-OBS dataset is derived through a three stage process [*Haylock et al.*, 2008]. Monthly

126 means (totals) of temperature (precipitation) are first interpolated to a 0.1 degree latitude by

127 longitude grid using three-dimensional (latitude, longitude, elevation) thin plate splines. Daily

128 anomalies, defined as the departure from the monthly mean (total) temperature (precipitation),

129 are interpolated to the same 0.1 degree grid, and combined with the monthly mean grid. For

130 temperature, daily anomalies are interpolated using kriging with elevation as an external drift

131 factor. For precipitation indicator kriging is first used, where the state (wet/dry) of precipitation

132 is first interpolated, after which the magnitude at 'wet' 0.1 degree grid points is interpolated

133 using universal kriging. Finally, the 0.1 degree points are used to compute area-average values

134 at the four E-OBS grid resolutions (0.25 and 0.5 degree regular latitude-longitude grid and 0.22

135 and 0.44 degree lat-long rotated-pole grids). In this paper, we use the 0.25 degree regular

136 latitude-longitude grid for further evaluation, as results for the other grids are essentially the

137 same.

138

139 Standard error estimates that accompany the gridded data are derived through combination of the

140 individual standard error estimates for monthly and daily interpolations. Standard error for the

141 monthly mean or total are the Bayesian standard error estimates, as available in the ANUSPLIN

142 package used for the spline interpolation [*Hutchinson*, 1995; *Wahba*, 1983]. Error estimates for

143 daily anomalies have been calculated using the method proposed by *Yamamoto* [2000] (see

144 Section 5). Both standard error estimates are calculated at the 0.1 degree master grid. For

145 temperature monthly and daily uncertainties are combined taking the square root of the sum of

146 the squares of the two uncertainties. For precipitation the relative uncertainty of the daily total is

147 the square root of the sum of the squares of the relative uncertainty of the monthly total and the

148 relative uncertainty of the daily proportion of monthly total precipitation. Uncertainties at the

149 0.1 degree grid have been averaged over the target grids allowing for spatial autocorrelation.

150 Details on the interpolation methods and how we implemented them as well as on the calculation

151 of the uncertainties are available in *Haylock et al.* [2008].

152

## 153 3. Homogeneity assessment

### 154 *3.1. Homogeneity testing*

155 To analyse the influence of inhomogeneities in station data on gridded time-series and to inform

156 the user about possible inhomogeneous areas within the dataset, we apply a homogeneity test to

157 the gridded dataset and compare results to the same test for station data. Numerous tests could

158 be used [e.g., *Peterson et al.*, 1998], but for this study we use the Wijngaard method [*Wijngaard*

159 *et al.*, 2003], which is the same test that was applied to the ECA&D station data used to construct

160 the E-OBS, where 39% of the precipitation and 25% of the temperature station series were found

161 to be potentially homogeneous over the period 1961 – 2006 [*Klok and Klein Tank*, 2008].

162

163 The Wijngaard method is an absolute test, as it does not use a supposedly homogeneous

164 reference series. This was appropriate for the version of the ECA&D dataset before the

165 ENSEMBLES project started, because of its sparse network [*Wijngaard et al.*, 2003]. It

166 comprises four homogeneity tests: the standard normal homogeneity test (SNHT) for a single

167 break [*Alexandersson*, 1986], the Buishand range test [*Buishand*, 1981], the Pettitt test [*Pettitt*,

168 1979] and the Von Neumann test [*Von Neumann*, 1941]. These location-specific tests have

169 different characteristics; for example, the SNHT test is more sensitive to inhomogeneities earlier

170 or later in the time-series, whereas the Buishand and Pettitt tests work better for breaks near the

171 middle of the series.  If zero or one of the tests detects a break at the 1% significance level the

172 time-series is classified 'useful'; if a break is detected by two tests the series is classified

173 'doubtful' and if three or four tests find a break, the series is classified 'suspect'.

174

175 For precipitation the annual wet day count is used for the analysis of breaks, as this statistic

176 generally has lower variance than total precipitation, enabling a better signal to noise ratio for

177 significance testing.  For temperature, the annual mean diurnal temperature range (mDTR) and

178 the annual mean of the absolute day-to-day differences of DTR (vDTR) are used for

179 homogeneity detection. DTR is used in preference to mean, maximum or minimum temperature,

180 as it has been shown that tests on DTR are more sensitive: breaks that are mainly radiation

181 related have different effects on minimum and maximum temperature and are, therefore, only

182 weakly apparent in these variables, but do appear clearly in DTR.  As the homogeneity tests are

183 applied to both mDTR and vDTR, a temperature station is classified according to the worst

184 outcome for the two variables.

185

186 We apply the Wijngaard tests to both station and E-OBS gridded data and compare the results.

187 We calculate the annual wet day count, mDTR and vDTR for each year if for each month no

188 more than 20% of the data are missing.  If less than 80% of the years in the period 1950-2006 are

189 present, the homogeneity test for that station or grid box is not performed, although these stations

190 may have been used for the interpolation. *Wijngaard et al.* [2003] concluded that a 1 mm

191 threshold should be applied to define a wet day because otherwise too many breaks were

192 detected, and we accordingly adopt this threshold.

## 3.2. Results and discussion

Figure 1 shows the stations and grid boxes that are potentially useful (green), doubtful (blue) or suspect (red), according the Wijngaard classification. For precipitation there are more many more useful stations and grid boxes than suspect ones. Suspect areas are mainly located in Northern Norway, Scotland, Italy, the Balkan, parts of Central Europe and in Northern Russia. For temperature most of Europe has a statistical significant inhomogeneity at some point in the gridded data, indicated by breaks in mDTR or vDTR (or both). However, if we only look at mDTR there are major differences (see Figure-S 1 in the supplementary material), with many more potential homogeneities in coastal areas, with remaining areas of central France, UK, Netherlands, parts of Spain and major parts of Ukraine, Northern Russia, Finland, southern Sweden, Czech Republic, Baltic States and Former Yugoslavia classified as useful in that case. That we find breaks in mDTR along the coast may be explained by a reduced variability in those areas due to the influence of the sea, making it easier to detect a break in mDTR. Inhomogeneities are much more widespread in vDTR with no clear difference between coastal and non-coastal areas.

Figure 1 also shows that the areas that have the most suspect stations often also have suspect grids, but sometimes even one suspect station may influence a whole area. An example of the latter is precipitation in northern Sweden where only one station is suspect, but has an influence over many grid boxes. Conversely, some stations have a smaller influence on the area, as, for example, in Russia where many stations are inhomogeneous, but only small areas are influenced. Many stations in this area have breaks in different years and these may be cancelled out in the gridded values. For temperature, inhomogeneous stations are present across the whole of Europe, which is reflected in the inhomogeneities of the gridded data.

218   In the case of precipitation many more areas of the grids are classified as potentially useful than

219   for temperature (78% for the wet day count versus 46% for mDTR and 28% for vDTR for the

220   grids, and 89% versus 49% and 56% for the stations, see Table 1), which is related to the fact

221   that the homogeneity test is less sensitive for the wet day count.  The percentage of stations that

222   are qualified useful is higher in this study than in the study of *Klok and Klein Tank* [2008] (89%

223   for the wet day count in this study vs. 39% in the *Klok and Klein Tank* study and 49% vs. 25%

224   for temperature).  The reason for this is most likely the time period used; we use the additional

225   first 11 years of the data, in which fewer stations have full data coverage.  When there are fewer

226   stations available, also fewer breaks are detected in the data.  mDTR has a much higher

227   percentage of useful grids than vDTR, whereas vDTR has a higher percentage of useful stations

228   than mDTR.  This indicates that in the station breaks are more strongly manifested in the mean

229   of the data, whereas in the grids breaks are more strongly manifested in the standard deviation.

230   That may be due to the fact that the variability of the grid values are dependent on the station

231   density of the network used for the interpolation and the distance to the grid centre [*Hofstra et*

232   *al.*, 2009].  A station network that does not have a constant density in time may introduce

233   inhomogeneities.

234

235   We also assessed the distribution of breaks in time and compare these between gridded and

236   station data (Figure 2).  As expected, the SNHT detects more inhomogeneities near the beginning

237   and end of the period than the Buishand and Pettitt tests.  SNHT also detects more breaks for any

238   one variable than the other tests (Table 1).  For wet day count the inhomogeneity in 1965

239   detected in the station data by the Pettitt test is also visible in the gridded data.  Breaks in the

240   1975-1985 period in the station data are mainly reflected in the gridded data close to 1980.  For

241   mDTR the breaks in station and gridded data do not show a specific pattern.  However, where for

242   vDTR the largest inhomogeneities in the station data are found around 1970, the largest breaks in

243 the gridded data are found in the early 1990s. The latter breaks may be due to a declining station

244 density around this time. We investigated whether inhomogeneities could be determined on a

245 decadal basis, by analysing each of the six decades separately, but the Wijngaard method is not

246 sensitive enough to find any inhomogeneities in these shorter periods at the 0.01 significance

247 level.

248

249 We also divided the calculated potential breaks for all three methods of the 57 year period into

250 six decadal groups and assess the inhomogeneities spatially (see Figures S2-S5 in supplementary

251 material). We can conclude, for example for precipitation, that most Italian and former

252 Yugoslavian stations around the Adriatic Sea with a break have this break in the period 1980-

253 1990 for all three tests; these breaks are also propagated through into the gridded data. For

254 precipitation, for all three tests in general, the timing of the breaks in the gridded and station data

255 compares quite well. For temperature, the agreement in timing of breaks between the station and

256 gridded data is smaller. For example, for vDTR a large part of Russia and the Ukraine has the

257 largest significant break between 1990 and 2000 for all three tests, whereas most stations in this

258 area suggest the largest break exists between 1960 and 1980. This indicates that there may be

259 multiple breaks in the station time-series of which one becomes more important in the gridded

260 data.

261

262 The inhomogeneities within the gridded data are important to keep in mind during any use of the

263 dataset. For example, when studying trends in the data, the results within the areas that are

264 suspect may not be meaningful. For those who require more detail on the inhomogeneities in the

265 gridded data, we have prepared a file that includes, for precipitation and temperature, the

266 potential classification of homogeneity of each 0.25 degree grid box (useful, doubtful, suspect)

267 and, for each of the four homogeneity tests, whether a statistical significant inhomogeneity has

11

268 been detected and if so the year of the largest break. The file can be downloaded from the E-

269 OBS download site (http://eca.knmi.nl/download/ensembles/ensembles.php).

270

# 4. Comparison with existing datasets

271

## 4.1. Existing datasets

272

273 In the second test of the dataset we compare E-OBS to existing datasets developed with much

274 denser station networks. Since station density is a very important factor in the interpolation and

275 the interpolation errors are smaller in areas with a dense station network [*Hofstra et al.*, 2008],

276 these existing datasets are deemed close to the 'true' areal average, and provide a useful

277 reference against which to judge the E-OBS dataset. The three existing datasets used are the UK,

278 Alps and ELDAS datasets. ELDAS and the Alps datasets only comprise precipitation data. The

279 UK dataset contains all four variables. We were unable to find or not allowed access to

280 additional datasets in other regions.

### 4.1.1. UK

281

282 The UK dataset, supplied by the UK Met Office, comprises a 5x5 km equal-area grid, covering

283 the period 1958 – 2002 for precipitation, 1995 – 2002 for minimum and maximum temperature

284 and 1995 – 2006 for mean temperature [*Perry and Hollis*, 2005]. This dataset is compiled from

285 a station network of 4400 stations for precipitation and 540 stations for temperature using

286 multiple regression with geographic factors as the independent variables, followed by inverse

287 distance weighting (IDW) of the residuals. In comparison, the ECA&D station network had 138

288 stations within this area, of which most had 70 - 85% of the data available for all variables. To

289 allow comparison with the E-OBS interpolations all grid-points within each 0.25 degree grid

290    used for the interpolation have been averaged.  We also compare this dataset to ELDAS (see

291    Section 4.1.3), for which a 1 degree grid is used.

## 4.1.2. Alps

293    The Alps dataset, comprising precipitation only, is an updated version of the climatology and

294    daily data described by *Frei and Schär [1998]* and *Schwarb* [2001], described in more detail by

295    *Hofstra et al.* [2008].  The data are available on a 0.25 by 0.1667 degree grid and cover the

296    period 1966 – 1999.  For the period 1966 – 1970 there are no data available over Austria and

297    after 1990 there are data quality issues with many of the Italian stations, so in our comparison,

298    we use the period 1966-1990, except for Austria, where the period 1971 – 1990 has been used.

299    The dataset is constructed through addition of daily anomalies to the long term climatological

300    mean.  Anomalies were interpolated from station data using a modified version of the Shepard

301    algorithm [an ADW technique, *Frei and Schär*, 1998; *Shepard*, 1984] and the long-term

302    climatology was derived with a local regression approach [PRISM, *Daly et al.*, 2002] specifically

303    calibrated for the Alps [*Schwarb et al.*, 2001].  The dataset is based on over 6500 station records.

304    In comparison, the E-OBS station network had 341 stations available within this area, with

305    majority having over 70% data presence.  To allow comparison with E-OBS on a common grid,

306    both datasets have been averaged to a 0.25 x 0.25 degree grid.

## 4.1.3. ELDAS

308    The ELDAS daily precipitation dataset was developed by *Rubel et al.* [2004] for the

309    Development of a European Land Data Assimilation System to predict Floods and Droughts

310    (ELDAS) project.  It covers Central and Northern Europe at 0.2 degree latitude by longitude and

311    covers the relatively short period of October 1999 to December 2000.  Some 21,600 stations

312    were used for the interpolation, compared to 2000 for E-OBS over the ELDAS domain.  Station

313    density is reasonably homogeneous, but areas such as Portugal, Belgium, Italy, the Balkan,

314      Czech Republic, the Baltic states and Scandinavia have a lower density than Spain, France, the

315      Netherlands, the UK, Denmark, Germany, Poland, Switzerland and Austria.  Interpolation was

316      done via the Precipitation Correction and Analysis method [*Rubel and Hantel*, 2001]; this

317      comprises a dynamical bias correction combined with an ordinary block kriging algorithm.  To

318      enable comparison, we averaged ELDAS and E-OBS to a common 1 degree latitude by

319      longitude grid.

## *4.2.  Comparison*

321      We compare E-OBS to the high-quality grids using five skill scores for temperature and six for

322      precipitation.  We calculate the skill scores for all data together to obtain overall scores, and also

323      on a grid-point basis to explore the spatial patterns in difference between grids.  We use the mean

324      absolute error (MAE), root mean squared error (RMSE), compound relative error (CRE) and

325      Pearson correlation (R) to assess temperature and the precipitation amount.  The Critical Success

326      Index (CSI) and Percent Correct (PC) are used to study precipitation state (wet or dry, where a

327      wet day is defined as having precipitation $\geq 0.5$ mm).  The skill scores are described in detail

328      elsewhere [*Hofstra et al.*, 2008], but we include an explanation of each score in the

329      supplementary material.  For precipitation we also divide the MAE and RMSE by the mean

330      precipitation for the grids in order to remove the influence of the amount of precipitation on

331      these two skill scores in each grid.

332

333      We note that the high-quality data are not true areal averages. However, given they are based on

334      order of magnitude denser networks than E-OBS, we expect them to be subject to smaller

335      interpolation errors. Thus we can only quantify differences between the datasets, which provide a

336      qualitative indication of potential errors in E-OBS, but should not be interpreted as errors of the

337      dataset.

## *4.3.  Results and discussion*

338

339  Table 2 provides an overview of the results of the skill scores, calculated 'globally' for each grid

340  pairing, as well as for each standard season.  At first sight, the datasets compare very well:

341  correlations, CSIs and PCs are high (for example, the global correlation coefficient for

342  temperature is approximately 0.99 and for precipitation 0.85-0.92), the CREs are small and

343  RMSEs are fairly small (for example, CRE is 0.02-0.04 and 0.18-0.36 for temperature and

344  precipitation).  However the mean differences between datasets are quite large.  RMSE is 0.7-0.9

345  for temperature and 2.2-2.4 for precipitation, apart from the Alps where it is larger, at 5.8.  MAE

346  shows similar, but smaller differences.  For precipitation, the relative RMSE varies between 0.73

347  (UK) to 1.3 over the Alps.  Relative difference between E-OBS precipitation and the other

348  datasets are smaller in winter (UK and ALPS) and autumn (ELDAS). The main reason for larger

349  differences between the datasets in summer is that in summer precipitation is mainly convective

350  rather than frontal.  During this season the correlation between stations is lower than in the other

351  seasons.  Interpolation with a larger station density will then produce better areal averages than

352  interpolation using a less dense network.  For mean and minimum temperature the datasets are

353  closer to each other in spring, whereas they compare better in winter for maximum temperature.

354

355

356  Figure 3 presents the results for precipitation spatially.  E-OBS compares best to the UK dataset,

357  as does the ELDAS dataset, suggesting that over the UK E-OBS is fairly reliable. The

358  differences are generally larger over the West of Scotland, where topography is an important

359  contributing factor to spatial variability in rainfall.  E-OBS does not agree as well with the Alps

360  dataset, where the topographic complexity means that the sparse E-OBS network does not result

361  in the same gridded data as the denser Alps network; although absolute errors are large because

362  precipitation is on average higher in the Alps, relative errors are also larger than in the UK.

363     Similarly, E-OBS compares poorly to ELDAS over Norway, due to the greater station density for

364     the ELDAS dataset in this topographically complex area. Finally, the E-OBS precipitation

365     dataset has virtually no stations available in northern Africa, which causes the poor agreement in

366     this area. Figure 4 shows the spatial pattern of skill for temperature over the UK. In general, the

367     agreement is good for all three temperature elements. Differences are greatest over Scotland

368     compared to the rest of the UK. That may be a result of the higher station density of the UK

369     network, which may have had more station data available at higher elevations in Scotland.

370     Differences in agreement between the grids are generally larger than differences between the

371     four seasons.

372

373     We also evaluate whether E-OBS shows a bias compared to the high density datasets, by

374     counting the frequency of days where E-OBS is more than $\pm 0.1$ standard deviations from the

375     high density dataset (Figure 5). For precipitation, E-OBS shows a negative bias at nearly all grid

376     boxes relative to the Alps and ELDAS datasets. Compared the ELDAS dataset, E-OBS is

377     positively biased over parts of Norway and at scattered locations elsewhere in Europe. Over the

378     UK, E-OBS rainfall tends to be negatively biased in areas of higher rainfall in the west, apart

379     from Northern Ireland where there is a positive bias (and also compared to ELDAS). For

380     temperature there are areas with a positive (too warm) and a negative (too cold) bias. One

381     striking feature is that areas such as Devon/Cornwall and Southern Wales, that are too warm for

382     minimum temperature, are often too cold for maximum temperature. The bias for temperature is

383     not consistent over the whole of the UK.

384

385     In Figure 6 we assess the difference between E-OBS and the high density datasets across the

386     distribution of precipitation amount and temperature. For this we calculate for each grid deciles

387     of temperature and precipitation (for all wet days). We then calculate for each day and each grid

388   the absolute difference between the E-OBS and the other datasets and plot the median, $5^{th}$, $25^{th}$,

389   $75^{th}$ and $95^{th}$ percentiles of these differences in each decile (Figure 6).  While precipitation is

390   biased towards smaller values in all deciles of the dataset, the bias is larger for more extreme

391   precipitation.  In the comparison of the $10^{th}$ decile for the Alps the error between the two datasets

392   can be as high as 16 mm, which is the median of the error when E-OBS is compared to the Alps

393   dataset (see median of 9-$10^{th}$ decile of E-OBS versus Alps comparison in Figure 6).  The reason

394   for this relates to the much higher station density in the other datasets.  For E-OBS, interpolation

395   typically occurs from more distant stations compared to the high density datasets; as extreme

396   precipitation events are usually more localised, they will be over-smoothed if a sparse network is

397   used.  For temperature, differences in error are similar for all deciles, with an average of around

398   0.5 °C.  The errors are slightly larger in the $1^{st}$ decile for minimum temperature and the $10^{th}$

399   decile for maximum temperature, which means that there are slightly larger errors in the

400   extremes, but overall extreme temperature events will be quite well represented [see also the

401   discussion of extremes in *Haylock et al.*, 2008].

402

403   We can conclude that the E-OBS shows quite large differences to the existing datasets based on

404   higher density station network.  While correlations overall, and on a grid-by-grid basis, are high,

405   relative differences in precipitation are large, and usually biased towards an underestimation. For

406   temperature (UK only), mean absolute differences are at least 0.5 °C.  The fact that the ELDAS

407   precipitation dataset shows a much better spatial match to the UK dataset than E-OBS underlines

408   the fact that E-OBS is fundamentally limited by its underlying station network.  As the E-OBS

409   network density over the UK is above average compared to density over the rest of Europe, we

410   can conclude that this issue is likely to be pervasive across much of the E-OBS domain.

411   Assessment of the agreement with existing datasets for all deciles of precipitation and

412   temperature shows that the errors are larger in the extremes than in the more average amounts of

413    precipitation or temperature. There seem to be significant problems with the underestimation of

414    precipitation extremes. Comparability is much higher for temperature than for precipitation, due

415    to the fact that temperature is a continuous variable as opposed to precipitation.

416

## 5. Uncertainty assessment

### 5.1. Calculation of uncertainties

419    *Brohan et al.* [2006] give an overview of all sources of all known and calculable uncertainty in

420    their HadCRUT3 gridded global monthly temperature dataset. Three groups of uncertainties

421    have been identified: 1) station error, 2) sampling error and 3) bias error. Station error includes

422    errors made during thermometer reading, possible adjustment of homogeneities, calculation of

423    the station normal, and processing of raw data. The sampling error is the difference between the

424    'true' spatial average and the interpolated estimate. It depends on, amongst others, the number

425    of stations in the grid box, the distribution of those stations and on the variability of the climate

426    in the grid box. The gridding method used by *Brohan et al.* [2006] is a simple area average of

427    the stations within a grid, which is different from the kriging method that we use, but the

428    sampling error of our gridding method will depend on the same factors. Two sources of bias

429    error are summarised by *Folland et al.* [2001]: urbanization effects [*Jones et al.*, 1990] and

430    thermometer exposure changes [*Parker*, 1994]. For precipitation a similar list of sources of

431    uncertainty can be made. Here we focus on sampling error as it is expected to be the largest

432    contributor to overall error. The objective here is to evaluate the accuracy of the estimates of

433    interpolation sampling error for daily anomalies used in E-OBS. As explained in the

434    introduction, these daily errors are estimated using the method proposed by *Yamamoto* [2000].

435

436    *Yamamoto* [2000] estimates the so-called 'interpolation standard deviation' at each grid point as

437    the weighted average of the squared differences between station and interpolated values as

438    follows:

439

440    $$s_0 = \sqrt{\sum_{i-1}^{n} \lambda_i \left[ z(x_i) - z*(x_0) \right]^2}$$    [1]

441

442    where $x_i$ (i =1,n) are the locations of the stations used for the interpolation and $\lambda_i$ are the weights

443    used in the kriging interpolation and z are the observed values at the i stations used for the

444    interpolation ($x_i$) and $z^*$ is the interpolated value at the location for the interpolation ($x_0$).

445

446    *Yamamoto* [2000] compared his interpolation standard deviation to the kriging standard

447    deviation and cross validation error. The kriging standard deviation is a standard by-product of

448    kriging and used widely as a measure of reliability of the kriging procedure. The interpolation

449    standard deviation has much larger correlation with cross-validation error than with the kriging

450    standard deviation. The reason for that is that the kriging standard deviation is not a true

451    estimate of uncertainty [*Journel and Rossi*, 1989; *Monteira da Rocha and Yamamoto*, 2000], as

452    it cannot properly measure local data dispersion [*Yamamoto,* 2000].

453

454    As we do not have the true grid values for evaluation, we adopt station cross-validation to test

455    the accuracy of the *Yamamoto* [2000] interpolation standard deviation. We estimate the daily

456    anomaly at each station in the ECA&D dataset used to construct E-OBS, using the same

457    interpolation approach used for E-OBS gridded data. Interpolation standard deviation is

458    calculated using equation [1] above and cross-validation error as the absolute difference between

459    the interpolated station value and the observed value:

19

460

461     $cve_0 = |z*(x_0) - z(x_0)|$                 [2]

462

463 We next transform the interpolation standard deviations into 95% confidence intervals by

464 multiplication with 1.96 (assuming a normal distribution) and addition to and subtraction from

465 the interpolated daily values for each station. We then count the number of times the observed

466 station value falls within the 95% confidence interval for the interpolated value, with the

467 expectation that if the confidence interval is an accurate estimate of interpolation uncertainty we

468 would expect the station value to fall outside the confidence interval approximately 5% of the

469 time.

470 ### *5.2. Results and discussion*

471 We first compare the cross-validation error (CVE) and interpolation standard deviation (ISD)

472 through scatter plots. Results are similar for all temperature variables, so we only show figures

473 for precipitation and minimum temperature.

474

475 Correlation between the CVE and ISD for both temperature and precipitation is positive (Figure

476 7). The relationship between CVE and ISD is stronger for precipitation (r=0.57) than minimum

477 temperature (r=0.33), which provides confidence that the spatial distribution of ISD will reflect

478 the spatial variability in interpolation error. The relationship is also closer to one-to-one for

479 precipitation, whereas for temperature, ISD tends to be too large at smaller CVE and vice versa.

480

481 However, a better test of the accuracy of the ISD is the count of the percentage of station values

482 falling outside the interpolation 95% confidence interval derived from the ISD (Figure 8). For

483 precipitation, the upper 95% limit is mostly exceeded between 5-10% of the time, while values

484 fall below the lower limit 10-25% of the time, indicating that while the upper limit is a

485 reasonable estimate, the lower limit is poorly defined, and that precipitation is frequently

486 significantly underestimated. For temperature, there are roughly equal numbers of values falling

487 above and below the 95% confidence interval, but as with precipitation, the number exceeds that

488 expected. Most stations have at least 10% of data falling outside the confidence interval, with

489 many stations having more than 25% of values outside the interval. There is also a clear north-

490 south gradient in the percentage of the precipitation values falling outside the confidence limits,

491 with the CI underestimation being much larger in the north. The main reason for this is the fact

492 that there are fewer rain days in the south of Europe, compared to the north. The error is smaller

493 when no or little precipitation is observed, compared to a situation when a lot of precipitation is

494 observed.

495

496 From this analysis, we can conclude that the interpolation standard deviation provided with the

497 data is a strong underestimation of the actual interpolation error and should be used with care.

498 Moreover, it has to be taken into account, that the confidence intervals available with the gridded

499 data only include interpolation sampling error and no station and bias errors.

500

## 6. Summary and Conclusions

502 We have analysed the new E-OBS European high-resolution gridded dataset of daily minimum,

503 maximum and mean temperature and precipitation in three ways. First, we assessed the

504 homogeneity of the gridded data and related this to the homogeneity of the station data.

505 Secondly, we compared the dataset to existing gridded datasets developed with denser station

506 networks. And finally, we evaluated the accuracy of the interpolation standard deviation, a

507 measure of interpolation error that is provided with the dataset. While the three issues we assess

508     do not give a complete overview of the reliability of the dataset, they do provide important

509     additional information for users of the dataset.

510

511     The results of the Wijngaard [2003] homogeneity tests show that there are many *potential*

512     inhomogeneities present in the gridded dataset. There are more statistically significant breaks

513     present in temperature than precipitation data, and within the temperature data, there are more

514     breaks for vDTR than mDTR variables. Inhomogeneities in the gridded data are often related to

515     inhomogeneities in the stations contributing to the value of the grid. However, this relation is not

516     the same for all areas. Sometimes an area is inhomogeneous even if there is only one

517     inhomogeneous station in the area (e.g. for precipitation in northern Sweden) and in other

518     occasions many stations are inhomogeneous, but the area is not effected (e.g. for temperature in

519     south-eastern France). The year of the break of inhomogeneous grids generally corresponds to

520     the year of the break of stations in the surrounding area, although the correspondence is better for

521     precipitation than for temperature. We provide a data file that contains, for temperature and

522     precipitation, information on the grid boxes where the data are potentially inhomogeneous. This

523     information will be critical when, for example, performing analyses of trends in extremes using

524     E-OBS. For a future update of the E-OBS dataset we recommend that the issue of

525     inhomogeneities is studied thoroughly. A balance will have to be found between the loss of

526     station data and the introduction of inhomogeneities and homogenisation of the station data

527     should be considered.

528

529     When compared to existing high-resolution regional gridded data for the UK, ALPS and Europe

530     (ELDAS) that are based on much denser station networks, E-OBS shows an excellent

531     correlation. However, mean absolute errors are significant, in the order 0.5 °C for temperature

532     and greater than 100% for precipitation. For both variables and all skill scores the datasets

533 compare worse in areas with more relief. For precipitation agreement is in general better in

534 winter, whereas for temperature agreement is mainly best in spring. In the case of precipitation,

535 E-OBS also shows a negative bias, indicating that E-OBS tends to be over-smoothed relative to

536 the high-density datasets. For temperature, E-OBS shows a small positive bias over quite large

537 areas, but some scattered areas have a stronger negative bias. Moreover, the E-OBS dataset

538 compares better to the mean of the variables of the existing datasets than to the extremes,

539 although differences are much larger for precipitation than for temperature. Consequently, the

540 dataset should be used with caution in comparison to RCM outputs, especially with respect to

541 evaluation of RCM precipitation extremes.

542

543 The uncertainty estimates available with the data only represent sampling, or interpolation,

544 errors. These are calculated by combining errors from both parts of the interpolation process,

545 namely interpolation of the monthly mean (temperature) or totals (precipitation) using thin plate

546 smoothing splines and the interpolation of daily anomalies using versions of kriging (see Section

547 2). We evaluated the daily interpolation error estimates, estimated using *Yamamoto*'s [2000]

548 interpolation standard deviation approach. A comparison of these errors with cross-validation

549 errors shows that for most of Europe cross-validation error is positively correlated with

550 interpolation standard deviation. However, the frequency with which the 95% interpolation

551 confidence interval is exceeded is much larger than expected, indicating that the interpolation

552 standard deviation significantly underestimates the actual interpolation error. The 95%

553 confidence limits are on average exceeded 25% and sometimes even over 50% of the time. In a

554 future update of the data we recommend that ensemble stochastic simulations, i.e. a set of

555 interpolated realisations should be considered for the estimation of uncertainties. These have

556 also been mentioned in *Haylock et al.* [2008] but have not been implemented due to time

557  constraints. *Bellerby and Sun* [2005] and *Teo and Grimes* [2007] suggest short-cuts that should

558  reduce the computing time required.

559

560  The E-OBS dataset is the first publically available dataset that covers the whole of Europe at a

561  very high spatial resolution for daily data. However, as this study reveals, there are some

562  potentially important limitations to the data. Inhomogeneities are present within the data, the

563  data show quite large absolute and relative differences and biases to existing datasets that have

564  been developed with very dense station networks, and the standard errors delivered with the data

565  appear to significantly underestimate the true interpolation error. This will have to be taken into

566  account when the data are used, e.g. for the evaluation of RCM outputs. Trends analysis may

567  also be affected by potential inhomogeneities in the data. In addition, the underestimation of

568  extremes within the data may, for instance, influence future predictions using RCM outputs

569  regarding flooding. Moreover, when using the standard errors that have been supplied with the

570  data it has to be taken into account that these errors only include interpolation sampling errors

571  and that they are an underestimation of the true error.

572

573  The E-OBS data will often be the only available dataset for studies of e.g. the comparison of

574  RCM outputs for the whole of Europe. With the collation of more data and hence better

575  availability, reconsideration of how to deal with inhomogeneities in station data and how to

576  improve the uncertainty estimates the data will improve in the future. However, users of the data

577  should take notice of the weaknesses mentioned in this paper and use the data appropriately.

578

579 **Acknowledgements**

584

585 **References**

586 Alexandersson, H. (1986), A homogeneity test applied to precipitation data, *Journal of*

587 *Climatology*, *6*, 661-675.

588 Bellerby, T. J., and J. Sun (2005), Probabilistic and ensemble representations of the uncertainty

589 in an IR/Microwave satellite precipitation product, *Journal of Applied Meteorology*, *6*, 1032-

590 1044.

591 Brohan, P., et al. (2006), Uncertainty estimates in regional and global observed temperature

592 changes: A new data set from 1850, *Journal of Geophysical Research*, *111*, D12106; 12101-

593 12121.

594 Buishand, T. A. (1981), The analysis of homogeneity of long-term rainfall records in the

595 Netherlands, KNMI Scientific Report WR 81-7, De Bilt, the Netherlands.

596 Dai, A., et al. (1997), Surface observed global land precipitation variations during 1900-88,

597 *Journal of Climate*, *10*, 2943 - 2962.

598 Daly, C., et al. (2002), A knowledge-based approach to the statistical mapping of climate,

599 *Climate Research*, *22*, 99-113.

600 Folland, C. K., et al. (2001), Global temperature change and its uncertainties since 1861,

601 *Geophysical Research Letters*, *28*(13), 2621-2624.

25

602     Frei, C., and C. Schär (1998), A precipitation climatology of the alps from high-resolution rain-

603         gauge observations, *International Journal of Climatology*, *18*, 873-900.

604     Haylock, M., et al. (2008), A European daily high-resolution gridded dataset of surface

605         temperature, precipitation and sea-level pressure, *Accepted by Journal of Geophysical*

606         *Research*.

607     Hofstra, N., et al. (2008), The comparison of six methods for the interpolation of daily, European

608         climate data, *Accepted by Journal of Geophysical Research*.

609     Hofstra, N., et al. (2009), The influence of interpolation and station network density on the

610         distribution and extreme trends of climate variables in gridded data, *Submitted to Journal of*

611         *Climate*.

612     Hutchinson, M. F. (1995), Interpolating mean rainfall using thin plate smoothing splines,

613         *International Journal of Geographical Information Systems*, *9*(4), 385-403.

614     Jones, P. D., et al. (1990), Assessment of urbanization effects in time series of surface air

615         temperature of land, *Nature*, *347*, 169-172.

616     Journel, A. G., and M. E. Rossi (1989), When do we need a trend model in kriging?,

617         *Mathematical Geology*, *21*, 715-739.

618     Klok, L., and A. M. G. Klein Tank (2008), Updated and extended European dataset of daily

619         observations, *Accepted by International Journal of Climatology*.

620     Monteira da Rocha, M., and J. K. Yamamoto (2000), Comparison between kriging variance and

621         interpolation variance as uncertainty measurements in the Capanema iron mine, state of

622         Minas Gerais - Brazil, *Natural Resources Research*, *9*, 223-235.

623     New, M. (1999), Uncertainty in representing observed climate, in *Representing uncertainty in*

624         *climate change scenarios and impact studies*, edited by T. Carter, et al., pp. 59-66, Climate

625         Research Unit, Norwich.

626 Parker, D. E. (1994), Effects of changing exposure of thermometers at land stations,

627      *International Journal of Climatology*, *14*, 1-31.

628 Perry, M., and D. Hollis (2005), The generation of monthly gridded datasets for a range of

629      climate variables over the UK, *International Journal of Climatology*, *25*, 1041-1054.

630 Peterson, T. C., et al. (1998), Homogeneity adjustments of in situ atmospheric climate data: a

631      review, *International Journal of Climatology*, *18*, 1493-1517.

632 Pettitt, A. N. (1979), A non-parametric approach to the change-point detection, *Applied*

633      *Statistics*, *28*, 126-135.

634 Rubel, F., and M. Hantel (2001), BALTEX 1/6-degree daily precipitation climatology 1996-

635      1998, *Meteorology and Atmospheric Physics*, *77*, 155-166.

636 Rubel, F., et al. (2004), Daily and 3-hourly Quantitative Precipitation Estimates for ELDAS,

637      edited, p. 32, Biometeorology Group, Univ. Vet. Med., Vienna.

638 Schwarb, M. C., et al. (2001), *Mean annual and seasonal precipitation in the European Alps*

639      *1971– 1990*, plates 2.6, 2.7 pp., Landeshydr. und Geol., Bern.

640 Shepard, D. S. (1984), *Computer mapping: The SYMAP interpolation algorithm*, 133-145 pp.,

641      Springer, New York.

642 Teo, C.-K., and D. I. F. Grimes (2007), Stochastic modelling of rainfall from satellite data,

643      *Journal of Hydrology*, *346*, 33-50.

644 Von Neumann, J. (1941), Distribution of the ratio of the mean square successive difference to the

645      variance, *Annals of Mathematical Statistics*, *12*, 367-395.

646 Wahba, G. (1983), Bayesian "Confidence Intervals" for the Cross-Validated Smoothing Spline,

647      *Journal of the Royal Statistical Society. Series B (Methodological)*, *45*(1), 133-150.

648 Wijngaard, J. B., et al. (2003), Homogeneity of 20th century European daily temperature and

649      precipitation series, *International Journal of Climatology*, *23*, 679-692.

650    Yamamoto, J. K. (2000), An alternative measure of the reliability of ordinary kriging estimates,

651        *Mathematical Geology*, *32*(4), 489-509.
652

653     **Table 1.** The fraction of stations or grids that are useful, doubtful or suspect and the

654     inhomogeneous fraction for each statistical test

655

| | | # stations or grids | Overall Fraction | | | Fraction with Breaks | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | | Useful | Doubtful | Suspect | SNHT | Buishand | Pettitt | Von Neumann |
| Wet day fraction | Stations | 836 | 0.892 | 0.044 | 0.064 | 0.123 | 0.072 | 0.114 | 0.087 |
| | Grids | 22176 | 0.781 | 0.078 | 0.141 | 0.219 | 0.164 | 0.216 | 0.166 |
| mDTR | Stations | 472 | 0.492 | 0.114 | 0.394 | 0.477 | 0.422 | 0.432 | 0.468 |
| | Grids | 21970 | 0.464 | 0.099 | 0.437 | 0.515 | 0.470 | 0.460 | 0.485 |
| vDTR | Stations | 472 | 0.555 | 0.097 | 0.348 | 0.434 | 0.388 | 0.400 | 0.381 |
| | Grids | 21970 | 0.275 | 0.113 | 0.612 | 0.738 | 0.630 | 0.580 | 0.697 |

656

657

658

**Table 2.** Skill scores for the comparison of the E-OBS gridded dataset with the UK, Alps, and

ELDAS gridded datasets for the four variables minimum, maximum and mean temperature and

precipitation. Skill scores have been calculated for each grid point and are then averaged.

**Annual**

|  |  | R | MAE | MAE/ mean | RMSE | RMSE/ mean | CRE | CSI | PC |
|---|---|---|---|---|---|---|---|---|---|
| UK | Minimum temperature | 0,984 | 0,687 | n/a | 0,895 | n/a | 0,041 | n/a | n/a |
|  | Maximum temperature | 0,991 | 0,597 | n/a | 0,780 | n/a | 0,024 | n/a | n/a |
|  | Mean temperature | 0,991 | 0,517 | n/a | 0,695 | n/a | 0,023 | n/a | n/a |
|  | Precipitation | 0,916 | 1,081 | 0,355 | 2,170 | 0,729 | 0,182 | 0,836 | 0,909 |
| Alps | Precipitation | 0,880 | 2,253 | 0,514 | 5,766 | 1,325 | 0,357 | 0,769 | 0,897 |
| Eldas | Precipitation | 0,846 | 1,159 | 0,457 | 2,419 | 1,009 | 0,316 | 0,744 | 0,874 |

**Winter**

|  |  | R | MAE | MAE/ mean | RMSE | RMSE/ mean | CRE | CSI | PC |
|---|---|---|---|---|---|---|---|---|---|
| UK | Minimum temperature | 0,971 | 0,700 | n/a | 0,918 | n/a | 0,082 | n/a | n/a |
|  | Maximum temperature | 0,977 | 0,507 | n/a | 0,680 | n/a | 0,056 | n/a | n/a |
|  | Mean temperature | 0,974 | 0,533 | n/a | 0,718 | n/a | 0,068 | n/a | n/a |
|  | Precipitation | 0,925 | 1,187 | 0,331 | 2,227 | 0,627 | 0,176 | 0,856 | 0,914 |
| Alps | Precipitation | 0,894 | 2,013 | 0,505 | 5,031 | 1,274 | 0,346 | 0,784 | 0,906 |
| Eldas | Precipitation | 0,848 | 1,256 | 0,458 | 2,360 | 0,926 | 0,373 | 0,759 | 0,869 |

**Spring**

|  |  | R | MAE | MAE/ mean | RMSE | RMSE/ mean | CRE | CSI | PC |
|---|---|---|---|---|---|---|---|---|---|
| UK | Minimum temperature | 0,973 | 0,663 | n/a | 0,860 | n/a | 0,069 | n/a | n/a |
|  | Maximum temperature | 0,981 | 0,640 | n/a | 0,822 | n/a | 0,051 | n/a | n/a |

| | | R | MAE | MAE/mean | RMSE | RMSE/mean | CRE | CSI | PC |
|---|---|---|---|---|---|---|---|---|---|
| | Mean temperature | 0,984 | 0,491 | n/a | 0,631 | n/a | 0,039 | n/a | n/a |
| | Precipitation | 0,916 | 0,893 | 0,359 | 1,803 | 0,730 | 0,181 | 0,828 | 0,908 |
| Alps | Precipitation | 0,881 | 2,237 | 0,514 | 5,345 | 1,231 | 0,365 | 0,775 | 0,888 |
| Eldas | Precipitation | 0,853 | 1,039 | 0,465 | 2,103 | 0,992 | 0,338 | 0,742 | 0,875 |

**Summer**

| | | R | MAE | MAE/mean | RMSE | RMSE/mean | CRE | CSI | PC |
|---|---|---|---|---|---|---|---|---|---|
| UK | Minimum temperature | 0,955 | 0,668 | n/a | 0,866 | n/a | 0,116 | n/a | n/a |
| | Maximum temperature | 0,970 | 0,709 | n/a | 0,896 | n/a | 0,087 | n/a | n/a |
| | Mean temperature | 0,965 | 0,520 | n/a | 0,700 | n/a | 0,082 | n/a | n/a |
| | Precipitation | 0,898 | 1,004 | 0,402 | 2,136 | 0,874 | 0,207 | 0,807 | 0,903 |
| Alps | Precipitation | 0,852 | 2,531 | 0,546 | 6,088 | 1,385 | 0,392 | 0,732 | 0,878 |
| Eldas | Precipitation | 0,826 | 1,026 | 0,514 | 2,003 | 1,334 | 0,577 | 0,690 | 0,870 |

**Autumn**

| | | R | MAE | MAE/mean | RMSE | RMSE/mean | CRE | CSI | PC |
|---|---|---|---|---|---|---|---|---|---|
| UK | Minimum temperature | 0,976 | 0,720 | n/a | 0,928 | n/a | 0,067 | n/a | n/a |
| | Maximum temperature | 0,987 | 0,518 | n/a | 0,667 | n/a | 0,035 | n/a | n/a |
| | Mean temperature | 0,983 | 0,526 | n/a | 0,709 | n/a | 0,042 | n/a | n/a |
| | Precipitation | 0,921 | 1,243 | 0,341 | 2,408 | 0,681 | 0,173 | 0,849 | 0,912 |
| Alps | Precipitation | 0,899 | 2,228 | 0,495 | 6,196 | 1,368 | 0,326 | 0,783 | 0,914 |
| Eldas | Precipitation | 0,863 | 1,226 | 0,431 | 2,511 | 0,911 | 0,306 | 0,765 | 0,879 |

662

**Figure 1.** Overall homogeneity, according to the Wijngard test, of the station network (top) and

the gridded data (bottom) for precipitation (left) and temperature (right). For temperature mDTR

and vDTR are combined, with the most negative outcome for the two variables used.

666

**Figure 2.** The fraction of stations and grid points with a statistically significant (0.01)

inhomogeneity in each year of the dataset. Inhomogeneities are calculated for the full 1950-2006

period.

670

**Figure 3.** A spatial overview of the skill scores R (-), MAE (mm), RMSE (mm), CRE (-) and

CSI for precipitation for the comparison of the E-OBS dataset with the datasets of the UK (top

row), Alps (2[nd] row) and ELDAS (3[rd] row) and the UK versus ELDAS (bottom row). MAE /

mean precipitation (-) and RMSE / mean precipitation (-) are added to remove the influence of

the average amount of precipitation in a grid cell on the skill score.

676

**Figure 4.** As

Figure **3**, but for the skill scores R (-), MAE (°C), RMSE (°C) and CRE (-) for minimum (top),

maximum (middle) and mean (bottom) temperature for the comparison with the UK dataset.

680

**Figure 5.** Spatial pattern of bias in the E-OBS dataset compared to higher quality data over the

Alps, ELDAS domain and UK, expressed: the percentage of days that E-OBS data are more than

0.1 standard deviations below the higher quality data, *subtracted* from the percentage of days the

E-OBS data are more than 0.1 standard deviation above the higher quality data. Thus, a positive

value indicates that E-OBS data tend to be biased greater than the higher quality data, and vice

versa. Precipitation is shown left, with UK top, Alps in the middle and ELDAS at the bottom.

687    Temperature (UK only) is shown right, with minimum temperature at the top, maximum

688    temperature in the middle and mean temperature at the bottom.

689

690    **Figure 6.** Absolute error in different deciles for each comparison with existing datasets for

691    precipitation (left) and temperature (right). In the left figure red is for the UK, green for the Alps

692    and blue for ELDAS, in the right figure red is for minimum temperature, green for maximum

693    temperature and blue for mean temperature. The box of absolute error shows the $0.25^{th}$, median

694    and $0.75^{th}$ percentile, the whiskers show the $0.05^{th}$ and $0.95^{th}$ percentile.  Deciles are calculated

695    for each grid separately.

696

697    **Figure 7.** Bivariate histograms showing the joint frequency distribution of cross validation error

698    and interpolation standard deviation for precipitation (left) and minimum temperature (right).

699    Both figures are on a log-log scale.

700

701    **Figure 8.** Spatial patterns of the percentage of interpolated data exceeding the lower (left) and

702    upper (right) limits of the 95% confidence interval for precipitation (top) and minimum

703    temperature (bottom) for all stations.  Insets display histograms of the frequency of the over- or

704    underestimation of the stations.

705