

Description of the COST-HOME monthly benchmark dataset and the submitted homogenized contributions

Victor Venema¹ (1), Olivier Mestre (2), Enric Aguilar (3), Ingeborg Auer (4), José A. Guijarro (5), Peter Domonkos (3), Gregor Vertacnik (6), Tamás Szentimrey (7), Petr Stepanek (8), Pavel Zahradnicek (8), Julien Viarre (3), Gerhard Müller-Westermeier (9), Monika Lakatos (7), Claude N. Williams (10), Matthew Menne (10), Ralf Lindau (1), Dubravka Rasol (11), Elke Rustemeier (1), Kostas Kolokythas (12), Tania Marinova (13), Lars Andresen (14), Fiorella Acquaotta (15), Simona Fratianni (15), Sorin Cheval (16), Matija Klančar (6), Michele Brunetti (17), Christine Gruber (4), Marc Prohom Duran (3), Tanja Likso (11), Pere Esteban (18, 19), Theo Brandsma (20)

- (1) Meteorological institute of the University of Bonn, Germany;
- (2) Météo France, Ecole Nationale de la Météorologie, Toulouse, France;
- (3) Center on Climate Change (C3), Tarragona, Spain;
- (4) Zentralanstalt für Meteorologie und Geodynamik, Wien, Austria;
- (5) Agencia Estatal de Meteorología, Palma de Mallorca, Spain;
- (6) Environmental Agency of the Republic of Slovenia, Meteorology, Ljubljana, Slovenia;
- (7) Hungarian Meteorological Service, Budapest, Hungary;
- (8) Czech Hydrometeorological Institute, Brno, Czech Republic;
- (9) Deutscher Wetterdienst, Offenbach, Germany;
- (10) NOAA/National Climatic Data Center, USA;
- (11) Meteorological and Hydrological Service, General Meteorology Division, Zagreb, Croatia;
- (12) Laboratory of Atmospheric Physics, University of Patras, Greece;
- (13) National Institute of Meteorology and Hydrology, Sofia, Bulgaria;
- (14) Norwegian Meteorological Institute, Oslo, Norway;
- (15) Department of Earth Science, University of Turin, Italy;
- (16) National Meteorological Administration, Bucharest, Romania;
- (17) Institute of Atmospheric Sciences and Climate (ISAC-CNR), Bologna, Italy;
- (18) Centre d'Estudis de la Neu i de la Muntanya d'Andorra (CENMA-IEA)
- (19) Grup de Climatologia, Universitat de Barcelona, Spain
- (20) Royal Netherlands Meteorological Institute, De Bilt, The Netherlands.

Version: 18 Juli 2011

¹ Corresponding author: Meteorological institute of the University of Bonn, Auf dem Huelgel 20, 53121 Bonn, Germany, Victor.Venema@uni-bonn.de.

Introduction

As part of the COST¹ Action HOME² a dataset has been generated that will serve as a benchmark for homogenisation algorithms. Members of the Action and third parties have been invited and are still welcome to homogenise this dataset³. The results of this exercise was analysed to obtain recommendations for a standard homogenisation procedure and are described in an upcoming article (Venema et al., 2011). Chapter two discusses the generation of this benchmark dataset, the climate variables considered, which types of data are in the benchmark dataset, how they have been produced, the ways to introduce artificial inhomogeneities, and the additional specifications such as length, missing data and trends. This chapter is an updated version of a report (Venema et al., 2009), which was available to the participants. A draft for the properties of the benchmark was developed at a WG1 meeting and was approved by the Management Committee of the Cost Action.

The homogenized data that was returned by the participant is described in Chapter 3. In total 25 contributions have been returned before the deadline at which the truth was revealed. Multiple late contributions, submitted after the deadline, are also described. The descriptions of the contributions were written by the respective participants.

References

- Venema, V.K.C., O. Mestre, and E. Aguilar. *Description of the COST-HOME monthly benchmark dataset with temperature and precipitation data for testing homogenisation algorithms*. Report, ftp://ftp.meteo.uni-bonn.de/pub/victor/costhome/monthly_benchmark/description_monthly_benchmark_dataset.pdf, 2009.
- Venema, V.K.C., O. Mestre, E. Aguilar, I. Auer, J.A.Guijarro, P. Domonkos, G. Vertacnik, T. Szentimrey, P. Stepanek, P. Zahradnicek, J. Viarre, G. Müller-Westermeier, M. Lakatos, C.N. Williams, M. Menne, R. Lindau, D. Rasol, E. Rustemeier, K. Kolokythas, T. Marinova, L. Andresen, F. Acquaotta, S. Fratianni, S. Cheval, M. Klancar, M. Brunetti, Ch. Gruber, M. Prohom Duran, T. Likso, P. Esteban, Th. Brandsma. *Benchmarking monthly homogenization algorithms*. Submitted, 2011.

¹ European Cooperation in the field of Scientific and Technical Research

² COST ACTION-ES0601: Advances in homogenisation methods of climate series: an integrated approach

³ For more information on the Action or on how to homogenise the benchmark dataset please read Section 4 of Chapter 2 and have a look on our homepage for the most current information: <http://www.homogenisation.org>.

Description of the monthly benchmark dataset

1. Introduction

This section explains the methodology behind the benchmark dataset. The origin of the input data is described in Section 2, the generation of the dataset in Section 3. Finally the organisation of the intercomparison experiment is detailed in Section 4.

The benchmark dataset contains both real inhomogeneous data as well as artificial data. Tests on inhomogeneous artificial networks allow analyse performances of homogenisation as a whole. In this case the truth (the unperturbed surrogate or synthetic series) is known. Tests on real inhomogeneous climate networks allow the analysis of the spread of solutions and thus give an idea of uncertainties involved in homogenisation. Furthermore, this data section is important to study the properties of the detected inhomogeneities, which will allow future benchmarking studies to insert more realistic inhomogeneities.

1.1. Climate elements

Based upon a survey among homogenisation experts conducted in the frame of HOME by its WG1, we chose to start our work with monthly values for temperature and precipitation. This is done as most detection/correction methods are prepared to work at this resolution¹. The focus on monthly data is because currently monthly data is utilized in a many climate research studies, e.g. related to change in the mean climate and trend estimation. Furthermore, for monthly data there is a wide range of homogenisation algorithms and thus more need for an intercomparison experiment. Another dataset for daily data is planned.

The selection of only two variables, allows us to keep the benchmark dataset reasonably small. As most homogenisation methods are at least partially manual and time consuming, limiting its size is important. Temperature and precipitation are selected because most participants of the survey consider these elements the most relevant in their studies and are considered a priority to have them homogenised in their respective datasets. Furthermore, they represent two important types of statistics. In homogenisation algorithms temperature is typically modelled as an additive process and precipitation as a multiplicative process. It is thus expected that algorithms that perform well for these two climate variables will also be suited for many others.

1.2. Types of datasets

The benchmark will have three different types of datasets discussed in detail below:

1. inhomogeneous climate networks;
2. surrogate networks and
3. idealised synthetic time series.

The data is distributed via an ftp server: `ftp://ftp.meteo.uni-bonn.de`. The benchmark dataset can be found in the directory: `/pub/victor/costhome/monthly_benchmark/inho`. You can download all files in this directory tree or the file `inho.zip`, with all precipitation and temperature data, respectively. The benchmark can also be accessed through HOME's web site (`http://www.homogenisation.org`).

1.2.1. Real datasets

Real datasets will allow comparing the different homogenisation methods with the most realistic type of data and inhomogeneities. Thus this part of the benchmark is important for a

¹ Daily data homogenisation, which is also part of HOME, relies very often - again according to our survey - on breakpoints detected over monthly to annual data. So - although this benchmark does not cover daily data it is also partially useful for the homogenisation of daily data.

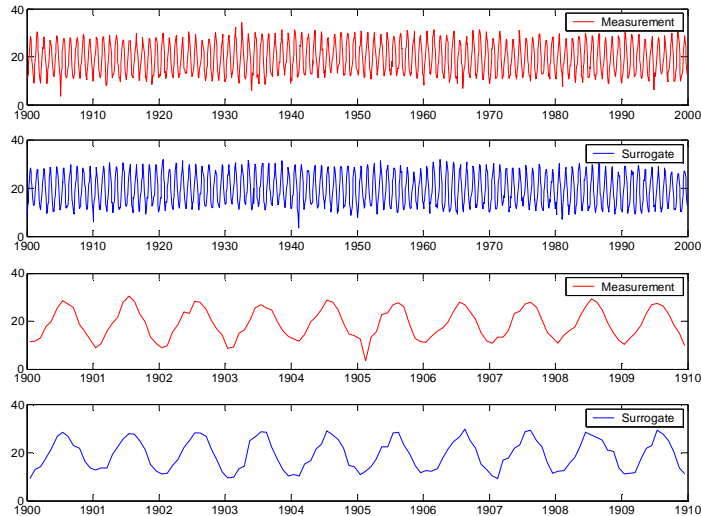


Figure 1. Example of a temperature measurement (red; first and third panel) and its surrogate (blue; second and fourth panel). The top two rows display the full 100-year time series, the lower two rows a zoom of a part of these time series.

faithful comparison of algorithms with each other. However, as in this case the truth is not known, it is not possible to quantify the improvements due to homogenisation. Therefore, the benchmark also has two datasets with artificial data.

The real data section of the benchmark dataset is located in the directories:

`/pub/victor/costhome/monthly_benchmark/inho/precip/real` and
`/pub/victor/costhome/monthly_benchmark/inho/temp/real`

1.2.2. Surrogate data

The aim of surrogate data is to reproduce the structure of measured data accurately enough that it can be used as substitute for measurements. Surrogate data can be useful when, e.g., not sufficient real measurements are available. The IAAFT¹ surrogates utilised in this work have the measured empirical distribution function and the (cross-)periodograms of a climate network. The (cross-)periodogram is equivalent to the (cross-correlation) autocorrelation function. In other words, the surrogate climate networks have the spatial and temporal correlations of real *homogenised* networks as well as the (possibly non-Gaussian) exact distribution of each station.

The periodogram is an estimate of the power spectrum, which describes how much variance is available at a certain time scale (wavenumber); it does not describe how much the variance itself varies (intermittence), the variance (at a certain scale) could be due to one large jump and many smaller jumps, or due to many medium sized jumps. The IAAFT algorithm tends to generate time series of the latter case that are not very intermittent (Venema et al., 2006a). This is also illustrated in Figure 4, which shows a so-called bounded cascade time series and its surrogate. The bounded cascade is easy to implement and generates fractal (self similar) time series; see e.g. Davis et al. (1996). The small-scale variability that belongs to the strong jumps in the bounded cascade time series is converted to small-scale variability all over the surrogate, effectively removing the jumps. This has as advantage that remaining inhomogeneities in the input dataset may not be too much of a problem. On the other hand the reduced number of large jumps may make the surrogates a bit too easy for the homogenisation

¹ IAAFT: Iterative Amplitude Adjusted Fourier Transform algorithm, developed by Schreiber and Schmitz (1996, 2000).

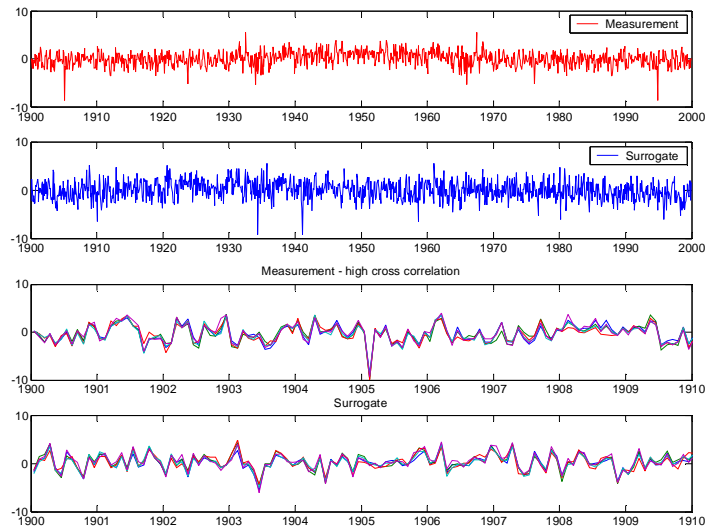


Figure 2. Example of the temperature anomaly, i.e. without the annual cycle, (red; first and third panel) and its surrogate (blue; second and fourth panel). The top two rows display the full 100-year time series, the lower two rows a zoom of a part of these time series.

algorithms compared to real data. See Venema et al. (2006a) for a more extensive discussion of structure of surrogate climate records.

Examples of surrogate time series are displayed in Figure 1 and 2. The ability of the surrogate networks to model the cross correlations is illustrated in Figure 3. To these surrogate networks known inhomogeneities are added that are described in Section 3.

Instead of surrogates one could use real data. However, less homogenised networks are available as needed. Furthermore, as discussed above surrogate data may have fewer problems with remaining inhomogeneities. Surrogate data will be much more realistic than Gaussian white noise. Both Markovian and Non-Markovian structures (e.g. long range correlations) are automatically accounted for. The implementation details of the algorithm are described in Section 3.

The surrogate data section of the benchmark dataset can be found on the ftp server: `ftp://ftp.meteo.uni-bonn.de` (and through `http://www.homogenisation.org/links.php`) in the directories:

```
/pub/victor/costhome/monthly_benchmark/inho/precip/sur1 and
/pub/victor/costhome/monthly_benchmark/inho/temp/sur1
```

1.2.3. Idealised synthetic time series

The idealised synthetic data is based on the surrogate networks. However, the differences between the stations have been modelled as uncorrelated Gaussian white noise. The idealised dataset is valuable because its statistical characteristics are assumed in most homogenisation algorithms and Gaussian white noise is the signal most used for testing the algorithms. This study will thus allow us to see if the surrogate data leads to different results. The implementation details are described in Section 2.

The surrogate data section of the benchmark dataset can be found on the ftp server: `ftp://ftp.meteo.uni-bonn.de` (and through `http://www.homogenisation.org/links.php`) in the directories:

```
/pub/victor/costhome/monthly_benchmark/inho/precip/syn1 and
/pub/victor/costhome/monthly_benchmark/inho/temp/syn1
```

For the homogenized data, see Chapter 3.

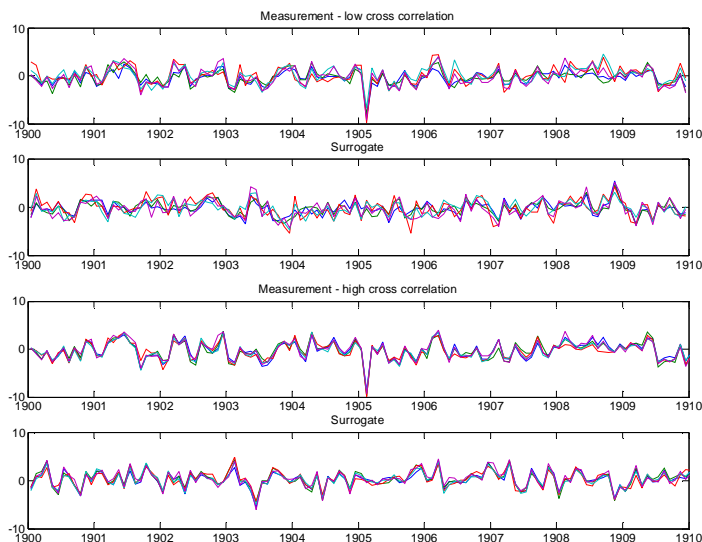


Figure 3. Example of first ten years of a network of measured temperature anomalies (first and third panel) and its surrogate (second and fourth panel). The top two rows are an example for a network with low cross correlations, the lower two rows for strong correlations.

2. Climate records

This benchmark dataset is limited to monthly temperature and precipitation records. Monthly data has pleasant statistical properties and is most used to study variations and trends in mean climate variables. Temperature and precipitation are both focus of many studies and represent two very different, additive and multiplicative, process models. In future other benchmarks may be established, e.g. for statistically more challenging daily data and other meteorological variables.

The input data is provided by Andorra, Austria, France, The Netherlands, Norway, Romania, and Spain; see Table 1. Data from Spain, Andorra and southern France is combined to one dataset called Catalonia. The inhomogeneous data is utilised for the real data section of the benchmark dataset. The homogenised data is used for the surrogate and synthetic data sections.

3. Benchmark generation

The benchmark dataset contains three types of data: inhomogeneous climate measurements (Section 3.1), surrogate data (Section 3.2) and synthetic data (Section 3.3). The basic surrogate and synthetic data represent homogenous climate networks. The inhomogeneities that are added to this data are described in Section 3.4.

3.1. Real data section

The inhomogeneous data sets from Table 1 have been used for the real data section; see Sec. 1.2.1. For most networks all stations are used, except for Norway. Because the Norwegian network is very large, this dataset is split into smaller subnetworks: 2 precipitation networks with 9 and 10 stations and two temperature networks with 7 stations.

For some of the networks metadata is available. Some of the metadata files are in a free text format, whereas other files have been formatted in the format for detected breaks (this data format was designed to return the results of the benchmark and not for metadata).

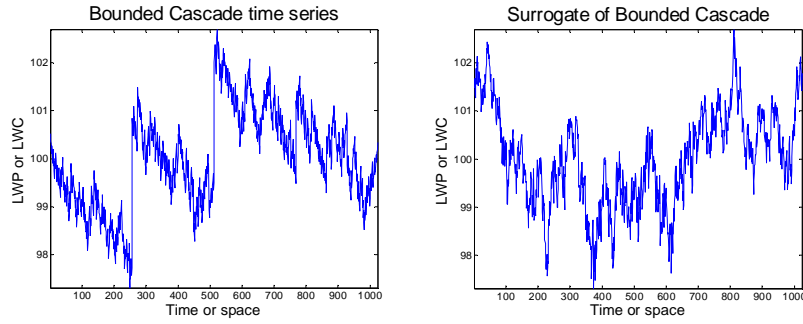


Figure 4. A realisation of the fractal bounded cascade method (left) and its surrogate (right).

The Austrian data originates from the Zentralanstalt für Meteorologie und Geodynamik (ZAMG) and was processed and homogenised as part of the HISTALP (Auer et al., 2007) and StartClim projects (Auer et al., 2008).

The dataset marked Catalonia contains data from Spanish Catalonia, Andorra and Southern France. The French data and the French part of the Catalonian data comes from the database BDCLIM of Météo-France. It was homogenised with PRODIGE (Caussinus and Mestre, 2004). The Andorran part of the Catalonian dataset comes from Pere Esteban (Centre d’Estudis de la Neu i la Muntanya d’Andorra (CENMA) de l’Institut d’Estudis Andorrans). The Spanish part of the Catalonian dataset comes from Servei Meteorològic de Catalunya (Meteocat). The real datasets are about one century long, except for Romania and Brittany, which contains about half a century of data.

The Austrian data was prepared by Christine Gruber (ZAMG). Marc Prohom (Meteocat) processed the data from Spanish Catalonia, Andorra and Southern France, called the region Catalonia. Olivier Mestre (Meteo France) delivered the French data sets. Theo Brandsma (Koninklijk Nederlands Meteorologisch Instituut, KNMI) provided the Dutch data. Lars Andresen (Meteorologisk institutt) supplied the Norwegian data and Sorin Cheval (Administratiei Nationale de Meteorologie) send in the Romanian data.

3.2. Surrogate data section

The surrogate data is used to multiply the number of available homogenised cases. As input it therefore needs homogenized data, which limits us to four datasets; see Table 1. To these surrogate climate networks we have added several breaks of known sizes and at known positions. The pre-processing of the climate data is described in Section 3.2.1 and the generation of the surrogates in 3.2.2. The addition of missing data, local and global trends and breaks points is described after the section on synthetic data in Section 3.4, because the same inhomogeneities are added to both surrogate and synthetic data.

Country	Inhomogeneous variables	Homogeneous variables
Austria		rr(43), tm(35)
France Bourgogne	rr(9)	rr(9)
France Brittany	tn(17), tx(17)	tn(17), tx(17)
The Netherlands	rr(11), tm(9)	
Norway	rr(189), tm(100)	
Romania	rr(22), tm(23)	
Catalonian region	rr(40), tn(30), tx(30)	tn(30), tx(31)

Table 1. Available input data for the benchmark dataset. Abbreviations used are rr: mean of daily rain sums, tn: monthly mean daily minimum temperature, tx: monthly mean daily maximum temperature and tm: monthly mean daily mean temperature. Between brackets is indicated the number of stations.

In total 20 surrogate temperature and 20 precipitation networks were generated. During the analysis it was found that some of the input stations were not well homogenised. As a consequence also the surrogate networks had too strong long term variability in the difference time series. Therefore 5 temperature and 5 precipitation networks needed to be removed. A stronger selection would not have led to difference results and was thus judged to be not necessary.

3.2.1. Pre-processing

In the pre-processing the annual cycle of the temperature records are removed, as well as linear trends for both elements. A number of stations is selected and the statistical input (periodogram and distribution) for the IAAFT algorithm that generates the surrogates is computed.

3.2.1.1. ANNUAL CYCLE

It would be possible to generate temperature surrogates including the annual cycle, but the annual cycle is not exactly a harmonic. Because the annual cycle is dominating for temperature, a more accurate reproduction of the temperature signal is expected by treating it explicitly. For every individual station, the annual cycle is computed as the mean temperature for every month. This annual cycle is subsequently subtracted. For precipitation no seasonal cycle is subtracted, only the mean values of the stations. After the generation of the surrogates the annual cycle (temperature) or mean (precipitation) is added again.

3.2.1.2. TRENDS

A trend in the time series would result in an artificially strong variability in the surrogate at large time scales. Therefore we have chosen to subtract the trend. This will also remove trend-like behaviour due to natural climate variability.

The trend is estimated by a linear regression. In case of temperature the predicted trend is subtracted from the time series. In case of precipitation the values are divided by the predicted trend and multiplied by the mean of the original measurement. After subtracting the trend, the mean is subtracted, which is added again after the generation of the surrogate.

3.2.1.3. STATION SELECTION

To reduce the total number of stations and to vary the strength of the correlations between the stations in a network we select a given number of stations from the measured records. The number of stations per network is listed in Table 2.

The selection algorithm starts with a random station. In the first stage, stations are added as long as the network is smaller as the wanted network and the mean correlation between the stations is larger than the wanted mean correlation; the stations with the highest mean

Variable	Setting	No. networks
<i>Input network</i>	Austria	7 12
	France	6 8
	Catalonia	7 0
<i>No. stations</i>	5	10
	9	5
	15	5
<i>No values</i>	100×12	

Table 2. Settings used for the generation of the surrogate and synthetic networks. The first column indicates the variable and the second column its value or its possible settings. In case more than one setting is available, the third column lists in how many of the surrogate or synthetic networks this setting is utilised. In case more than one value is indicated in a line with a vertical line in between, the first value is for temperature and the second one for precipitation.

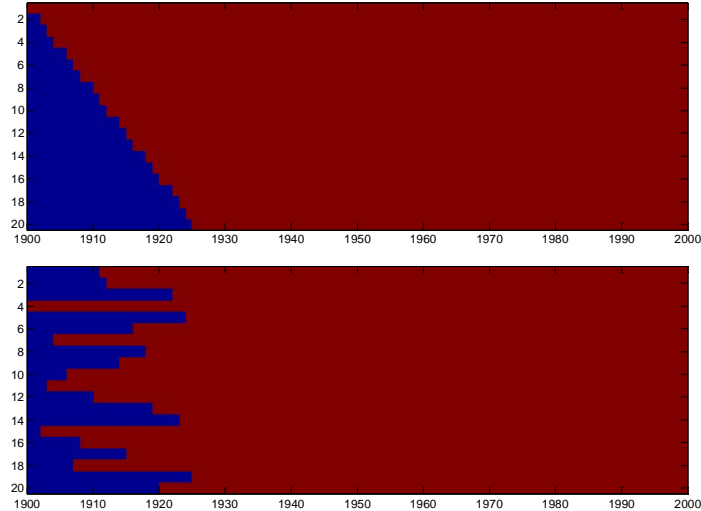


Figure 5. Missing data (blue) at the beginning of the network. The top panel shows the stations sorted for their length. The lower panel shows the same in a random arrangement.

correlation with the network are added. Adding stations will tend to increase the distance between the stations and decrease the mean correlation. After the first stage the network may have the wanted size and mean correlation that is smaller as wanted; in this case the selection is finished. Alternatively the network may have a mean correlation that is a little larger as wanted and a larger number of stations. Therefore, in the next stage, the possibly too large network is reduced while trying to reduce the mean correlation as little as possible. This is performed by removing each time one of the stations that are strongest correlated with one other station.

3.2.1.4. COMPUTATION OF STATISTICS

The surrogate (and synthetic) data contain data for one century from 1900 to 1999. Some homogenised temperature networks are shorter. A 100-year network is derived from the measurement by mirroring the measured networks multiple times until it is longer than 100 a. For example if the measured network contains data from 1940 to 2000, its mirrored version contains data in reversed order from 2000 to 1940 plus the data from 1940 to 2000. This mirrored network is then truncated to 100 a.

From the mirrored network the complex Fourier coefficients are computed by means of a 1-dimensional Fast Fourier Transform (FFT). A Hamming Fourier window is applied to reduce crosstalk (leakage); this causes a small reduction of the variance on large temporal scales (just below 100 a), but improves the estimate of the Fourier coefficients at short time scales (high frequencies). The magnitude of a Fourier coefficient describes the variability of the time series for a certain frequency. The phase difference (between two stations) of the Fourier coefficients determines the strength of their cross-correlation.

3.2.2. Generation of surrogate networks

The surrogates are computed with the multivariate Iterative Amplitude Adjusted Fourier Transform (IAAFT) algorithm developed by Schreiber and Schmitz (1996, 2000). Below the algorithm is described in words; the exact formulation for this version with all the needed equations can be found in Schreiber and Schmitz (2000), with a small modification of the second iterative step described in Venema et al. (2006b). To work with the benchmark you do not have to understand the algorithm. All you need to know is that the surrogates have exactly the distribution of the input and almost the same Fourier spectrum. As the Fourier coefficients determine the cross- and auto-correlation functions, this means that the surrogates have some

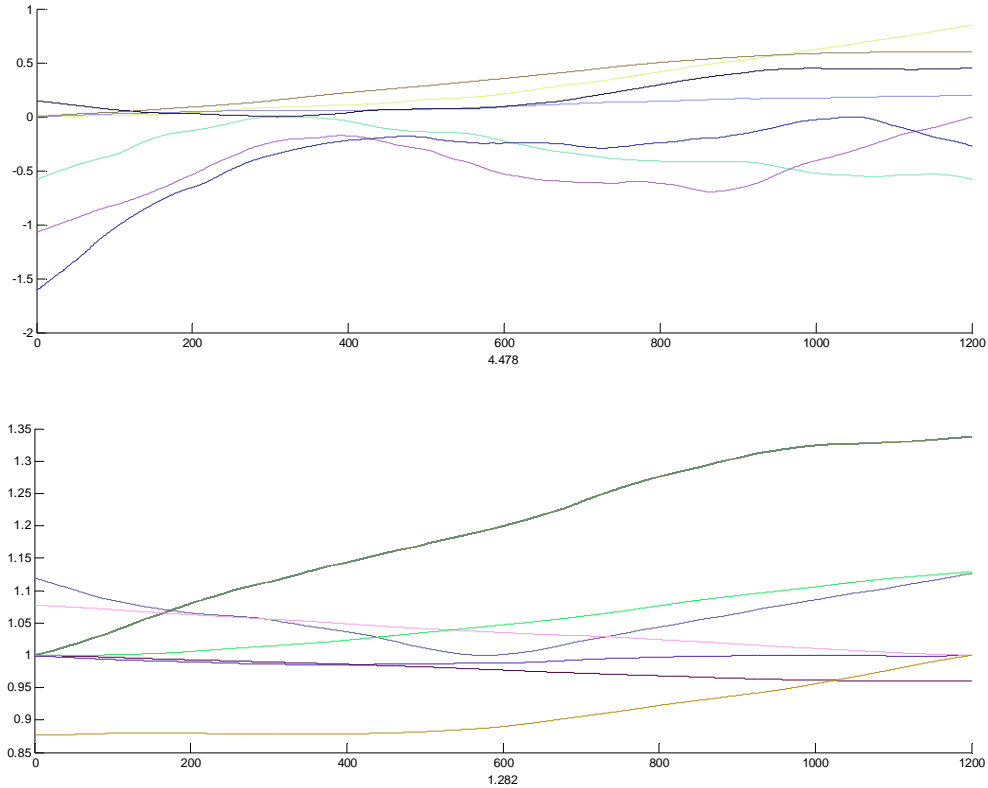


Figure 6. Examples of the signal that was added as a global trend to every station in the networks. The top panel are the functions added to the temperature networks, the lower panel depicts the factor for precipitation.

differences in their spatial and temporal structure. These deviations are expected not to be important for this application.

This iterative algorithm starts with white noise and has two iterative steps. In the first step, the spectral adjustment, the Fourier coefficients are adjusted. First, the Fourier transform of the surrogate from the previous iteration (or from the initial white noise) is computed. The magnitudes of the Fourier coefficients are changed to the magnitudes of the coefficients of the input network, i.e. the mirrored measured network. The phases are adjusted such that the phase differences between the Fourier coefficients of two stations are the same as the ones from the input network. This leaves one set of absolute phases free for the algorithm. These are computed such, that the time series are changed as little as possible.

In the second iterative step, called the amplitude adjustment, the values of the surrogate are set to be identical to those of the input network. In other words, after this step, the temperature or precipitation distribution of the surrogate is identical to the distribution of the input. This adjustment is performed by setting the highest value after the spectral adjustment to the highest value of the input. The same is performed for the second highest, etc. This is computationally solved by sorting the time series (while remembering their positions), substituting the values and reversing the sorting. As this adjustment is performed for every station separately, after this step every stations has the distribution of its corresponding input station.

Due to the amplitude adjustment, the Fourier coefficients are also changed some. For this reason the algorithm is iterative. Typically every iteration needs less and less adjustments. However the algorithm normally does not converge fully, but gets stuck in a local minimum. After the deadline we discovered that not all input networks were well homogenized and that thus also some surrogate networks contained large scale variability (in its difference or ratio

time series). As a consequence, we have removed 5 temperature and 5 precipitation networks. Selecting stronger did not clearly change any of the studied validation scores anymore.

3.3. Synthetic data section

The synthetic data aims to have difference (or ratio) time series that are Gaussian white noise. This type of fluctuation is assumed in most homogenisation algorithms. By analysing the difference in homogenisation performance between surrogate and synthetic data we want to investigate how important deviations from this idealisation are. Therefore, we tried to make the synthetic data as similar to the surrogate data in all other aspects. Thus we have created pairs of networks where the statistical properties of the synthetic network are based on its surrogate counterpart.

The generation of the synthetic data start by computing a time series with the network mean precipitation or temperature. The difference (temperature) or ratio (precipitation) of this mean to each station is computed. This relative time series is converted to Gaussian white noise (with the same mean, standard deviation and a similar spatial cross correlation matrix) and added (or multiplied) to the network mean time series. The cross correlation matrix (R) is reproduced in the following way. The spectral decomposition of R is $R=Q\Lambda Q^T$, where $\Lambda=\text{diag}[\lambda_1,\dots,\lambda_N]$ is the diagonal matrix of Eigen values of R and Q holds the Eigen functions (N is the number of stations). In this way we can compute $\sqrt{R} = Q\sqrt{\Lambda}Q^T$, where $\sqrt{\Lambda}=\text{diag}[\sqrt{\lambda_1},\dots,\sqrt{\lambda_N}]$. To compute noise with the spatial cross correlations R , one needs to multiply \sqrt{R} with white noise.

After the transformation to a Gaussian distribution, negative rain rates may occur; these values are explicitly set to zero. The cross correlations between the precipitation stations are unfortunately biased towards too low correlations. The cross correlation matrix of the ratio time series of the synthetic data is near the one of the surrogate data, but after multiplying the ratio time series to network mean time series the cross correlations are perturbed.

3.4. Inhomogeneities

The surrogate and synthetic data represent homogeneous climate data. To this data known inhomogeneities are added: outliers (Section 3.4.3) as well as break inhomogeneities (Section 3.4.4) and local trends (Section 3.4.5). Furthermore missing data is simulated (Section 3.4.1) and a global trend is added (Section 3.4.2). For a summary see Table 3.

Unknown to the participants before the deadline was that one of the surrogate and one of the synthetic networks did not contain any breaks or local trends (only missing data and outliers).

3.4.1. Missing data

Two types of missing data are added. Firstly, missing data at the beginning (older part) of the dataset, reproducing a gradual build up of the network, a very common situation in real datasets. Secondly, missing data in a large part of the network due to World War II, this is typical for European datasets.

This missing data at the beginning is modelled by a linear increase in the number of stations from three stations available in 1900 to all stations present in 1925. New stations always start in January. See Figure 5 for an example for a network of 20 stations.

There is a 50% chance that the data is missing in 1945. For the years preceding, 1945 down to 1940, the stations with missing data have a probability of 50 % that the data for the previous year is also missing.

3.4.2. Global trend

A global trend is added to every station in a network to simulate climate change. A different trend is used for every network. It is generated as very smooth fractal Fourier “noise” with a

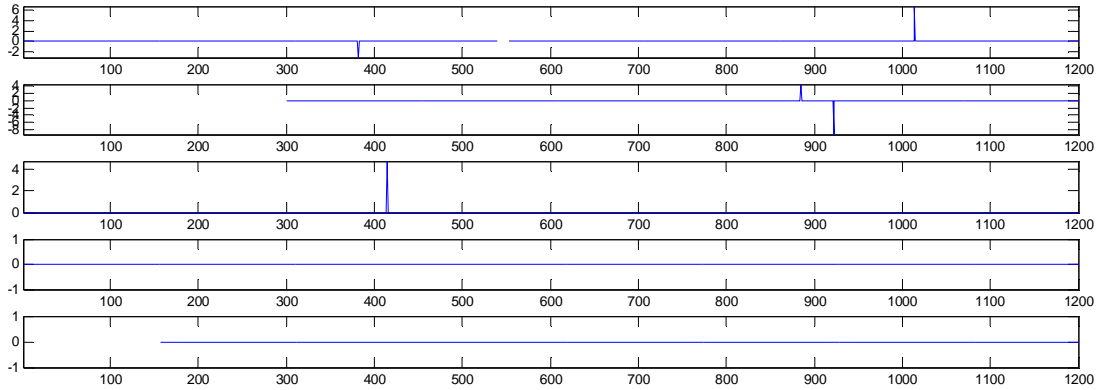


Figure 7. Example of perturbations due to outliers in a network of five temperature stations.

power spectrum proportional to k^{-4} ; only part of the signal is used to avoid the Fourier periodicity. The noise is normalised to a minimum of zero and a maximum of unity and then multiplied by a random Gaussian number. The width of the Gaussian distribution is one degree or ten percent. Both the temperature and the precipitation have a trend; see Figure 6. The precipitation trend is added by multiplying the time series with $1+\text{noise}$. Relative homogenisation algorithms should not be affected by this trend.

3.4.3. Outliers

Outliers are generated with a frequency of 1 per 100 a per station. The outliers are added to the anomaly time series, i.e. without the annual cycle for temperature. The value of the outliers is determined by a random value from the tails of the measurement. For temperature values below the 1st and above the 99th percentile are used. For precipitation values below the 0.1th and above the 99.9th percentile are used. See Figure 7 for an example of outliers added to a temperature network.

3.4.4. Break points

Breaks points are inserted by multiplying the precipitation with a factor or adding a constant to temperature; these constants and factors are different for every month. Two types of break points are inserted: random and clustered ones.

Random breakpoints are inserted at a rate of one per twenty years, i.e. on average five per station. This ratio is very close to what is found in most European datasets. To vary the quality station by station, first a random number between 2 and 8 % is chosen from a uniform distribution. The probability of a break event in any month is Poisson process with a break probability given by the random number divided by 12. It is required that there is data on both sides of the break, i.e. no breaks are inserted in the missing data at the beginning of the networks, but (multiple) breaks may be inserted in other missing data periods, close to each other or close to the edge.

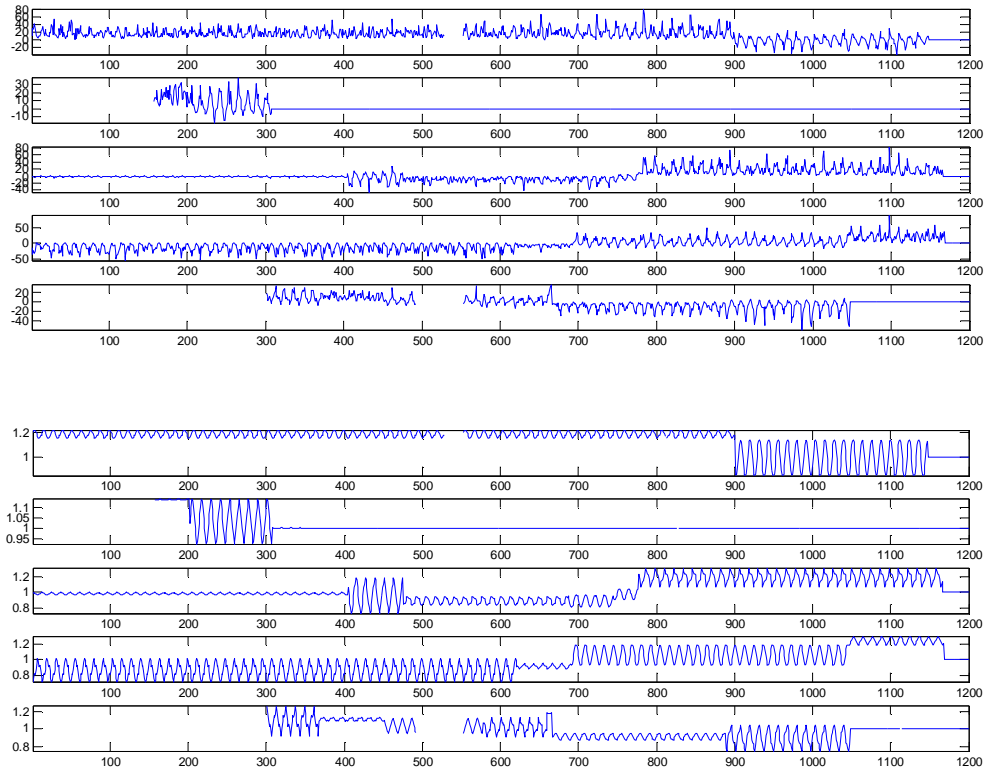


Figure 8. Example of perturbations due to random breaks in a network of five temperature (top) or precipitation (lower panel) stations.

The size of the break points is selected at random from a Gaussian distribution with a standard deviation of 0.8° for temperature and 15 % for rain. These mean break sizes have a seasonal cycle with standard deviation 0.4° and 7.5 %. The seasonal cycle is generated starting with white Gaussian noise. This noise is smoothed twice by a rectangular window of six month (using periodic boundary conditions), its mean is subtracted and its standard deviation set to a random number with the above mentioned standard deviations. Finally the minimum or maximum of the seasonal cycle is shifted towards the summer. These breaks are added to the temperature series. The break size plus one is multiplied with the precipitation time series. See Figure 8 for an example of the perturbations due to random breaks added to a network of 5 stations.

In 30 percent of the surrogate and synthetic networks we additionally add clustered breaks. In the networks affected by these spatially clustered breaks 30 % of the stations have a break point at the same time. The mean size and seasonal cycle of these breaks have the same distribution and properties as the random breaks. They are generated by computing one "network-wide" break and one temporary station-specific break for every affected station. These two breaks are then combined for the final break of a certain station using the equation: $0.8 \cdot \text{network-wide break} + 0.2 \cdot \text{temporary station break}$. At both sides of a clustered break there should be at least 10 % of all data.

3.4.5 Local trends

In 10 % of the temperature stations a local linear trend is introduced. The beginning of the trend is selected randomly. The length of the trend has a uniform distribution between 30 and 60 a. The beginning and the length are drawn anew until the end point is before 2000. The size of the trend at the end is randomly selected from a Gaussian distribution with a standard deviation of 0.8° . In half of the cases the perturbation due to the local trend continues at the

end, in the other half the station returns to its original value. An example of how the temperature can be affected by local trends is given in Figure 9.

3.5. Quicklook plots

In the directories with the data, there are also a number of quicklook plots to help get a quick overview of the data in this network. The file `overview_data.png` shows the data itself (including the annual cycle). The file `overview_data_anomalies.png` shows the same data, but without the annual cycle (temperature) or mean precipitation rate.

The file `correlations.png` shows the cross correlation matrix of the anomalies of the inhomogeneous data. The mean cross correlation and the mean maximum correlation shown in the title does not include the (unity) auto-correlations. The file `simple_map.png` shows the position of the stations in geographical coordinates. Every station is indicated by a red cross and its name. After the name in brackets the height in meter is given. Important cities are indicated by a black square and large rivers by a blue line; land is green and the sea is grey.

4. Organisation, homogenisation and analysis

The inhomogeneous benchmark dataset and the above mentioned quicklooks can be found on the ftp-server of the Meteorological Institute of the University of Bonn: `ftp://ftp.meteo.uni-bonn.de` in the directory: `/pub/victor/costhome/monthly_benchmark`. The homogenized benchmark dataset with all the returned contributions is found in the directory: `ftp://ftp.meteo.uni-bonn.de/pub/victor/costhome/homogenized_monthly_benchmark/`. The description of the data format can be found via the HOME homepage: `http://www.homogenisation.org` or at: `http://www.meteo.uni-bonn.de/venema/themes/homogenisation/`.

Participants are asked to homogenise the data and send it back in the same data format and directory structure. Only the top directory has to be renamed from `inho` (inhomogeneous) to `h000` (h for homogeneous, we will convert 000 to a running number). Furthermore, to every directory with data the participants will have to add a file with detected breaks points and a file with quality flags for every station. Preferably the returned data should be made available on an ftp-site or homepage and the link send to us. If this is really not possible please zip the data (only data no pictures) with maximum compression and send it as attachment to: `victor.venema@uni-bonn.de`.

Everyone is still welcome to participate for some time. Explicitly, scientists outside of the COST HOME Action are invited. We realise that the benchmark dataset is large, we thus would like to emphasis that also partial contributions (even if it is just one or a few networks)

Variable	Setting
<i>Outlier frequency</i>	1 %
<i>Freq. random breaks</i>	1/20 a
<i>Freq. clustered breaks</i>	30 % networks
<i>Fraction of stations with clustered breaks</i>	30 %
<i>Standard deviation of the size of the breaks</i>	0.8 0.1
<i>Fraction local trends</i>	10 % 0 %
<i>Standard deviation of the size local trends</i>	0.8 0
<i>Length local trends</i>	30 to 60 a

Table 3. Settings used for the generation of inhomogeneities in the surrogate and synthetic networks. The first column indicates the variable and the second column its value. In case more than one value is indicated in a line with a vertical line in between, the first value is for temperature and the second one for precipitation.

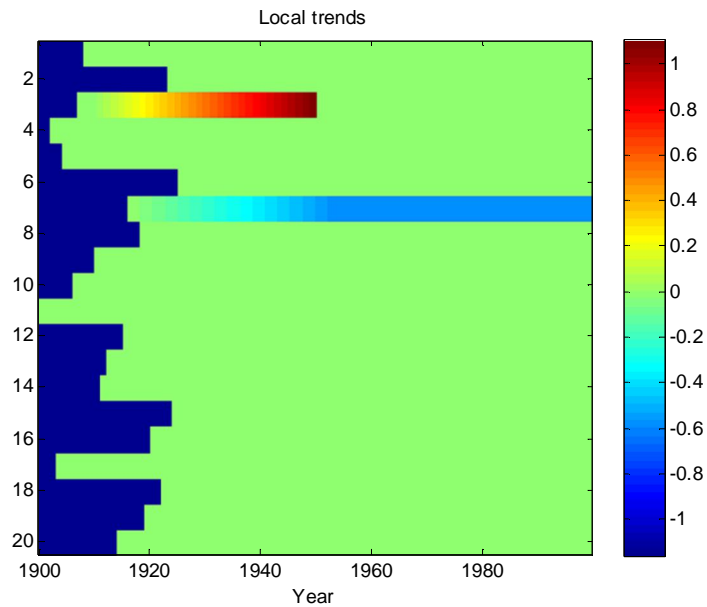


Figure 9. Example of perturbations due to local trends in a network of twenty temperature stations.

are valuable to us. In this case, please start with the surrogate networks in ascending order, so that we can study these networks in more detail with the largest number of algorithms. For the comparison between surrogate and synthetic data, it is sufficient if only the people with fully automatic algorithms homogenise the synthetic data.

Acknowledgements

The specifications of this benchmark dataset were derived based on discussions with the core group of the COST Action HOME, especially Olivier Mestre (Meteo France) and Enric Aguilar (University Rovira i Virgili of Tarragona, Spain). This dataset would have been impossible without contributed climate records. Many thanks thus go to Christine Gruber (for the Austrian data), Marc Prohom (for the data from the region Catalonia), Olivier Mestre (for the French data), Theo Brandsma (for the Dutch data), Lars Andresen (for the Norwegian data) and Sorin Cheval (for the Romanian data). José A. Guijarro of the Agencia Estatal de Meteorología and Petr Štěpánek of the Czech Hydrometeorological Institute are acknowledged for their efforts in checking the quality of preliminary versions of this benchmark.

References

- Auer, I., R. Böhm, A. Jurkovic, W. Lipa, A. Orlik, R. Potzmann, W. Schöner, M. Ungersböck, C. Matulla, K. Briffa, P.D. Jones, D. Efthymiadis, M. Brunetti, T. Nanni, M. Maugeri, L. Mercalli, O. Mestre, J.M. Moisselin, M. Begert, G. Müller-Westermeier, V. Kveton, O. Bochnicek, P. Stastny, M. Lapin, S. Szalai, T. Szentimrey, T. Cegnar, M. Dolinar, M. Gajic-Capka, K. Zaninovic, Z. Majstorovic, and E. Nieplova. HISTALP – Historical instrumental climatological surface time series of the greater alpine region. *Int. J. Climatol.*, **27**, pp. 17-46, doi: 10.1002/joc.1377, 2007.
- Auer, I., A. Jurkovic, R. Böhm, W. Schöner, and W. Lipa. *Endbericht StartClim2007A. Erweiterung und Vervollständigung des StartClim Datensatzes für das Element tägliche Schneehöhe. Aktualisierung des existierenden StartClim Datensatzes (Lufttemperatur, Niederschlag und Dampfdruck) bis 2007 04*. CD in print, 2008.
- Davis, A., A. Marshak, W. Wiscombe, and R. Cahalan. Multifractal characterizations of intermittency in nonstationary geophysical signals and fields. In *Current topics in*

- nonstationary analysis*, ed. G. Trevino, J. Hardin, B. Douglas, and E. Andreas, pp. 97-158, World Scientific, Singapore, 1996.
- Caussinus, H., and O. Mestre. Detection and correction of artificial shifts in climate series. *Appl. Statist.*, **53**, part 3, pp. 1-21, 2004.
- Schreiber, T., and A. Schmitz. Improved surrogate data for nonlinearity tests. *Phys. Rev. Lett.*, **77**, pp. 635-638, 1996.
- Schreiber, T., and A. Schmitz. Surrogate time series. *Physica D*, **142**(3-4), pp. 346-382, 2000.
- Venema, V., S. Bachner, H. Rust, and C. Simmer. Statistical characteristics of surrogate data based on geophysical measurements. *Nonlinear Processes in Geophysics*, **13**, pp. 449-466, 2006a.
- Venema, V., F. Ament, and C. Simmer. A stochastic iterative amplitude adjusted Fourier transform algorithm with improved accuracy. *Nonlinear Processes in Geophysics*, **13**, no. 3, pp. 247-363, 2006b.

Submitted homogenized contribution

3.1 Introduction

This section describes the main characteristics of the tested methods. This report will only list the features used to homogenize the benchmark; many tools have additional possibilities. Most of the algorithms test for relative homogeneity, which implies that a candidate series is compared to some estimation of the regional climate (“comparison phase”). Comparison may be performed using one composite reference series assumed homogeneous (e.g. SNHT), several ones, not assumed homogeneous (MASH), or via direct pairwise comparison (USHCN, PRODIGE); see Table 4. The comparison series are computed as the difference (temperature) or ratio (precipitation) of the candidate and the reference. The time step of comparisons may be annual, seasonal or monthly (parallel or serial). When several comparisons are made because multiple references are used or monthly data are analyzed in parallel, a synthesis phase is necessary, that may be automatic, semi-automatic, or manual.

The comparison series are tested for changes. Detection implies a statistical criterion to assess significance of changes, which may be based on a statistical test – Student, Fisher, Maximum Likelihood Ratio (MLR) test, etc. – or criteria derived from information theory (Penalized likelihood). Detection requires an optimisation scheme, to find the most probable positions of the changes among all possibilities. Such a searching scheme may be exhaustive (MASH), based on semi-hierarchic binary splitting (HBS), stepwise, moving windows (AnClim) or based on dynamic programming (DP).

Table 4. Comparison and detection methods.

<i>Method</i>	Comparison		Detection		References
	<i>Comparison</i>	<i>Time step</i>	<i>Search</i>	<i>Criterion</i>	
<i>MASH</i>	Multiple references	Annual Parallel monthly	Exhaustive	Statistical test (MLR)	Szentimrey, 2007, 2008
<i>PRODIGE</i>	Pairwise Human synthesis	Annual Parallel monthly	DP	Penalized Likelihood	Caussinus and Mestre, 2004
<i>USHCN</i>	Pairwise Automatic synthesis	Serial monthly	HBS	Statistical test (MLR)	Menne et. al., 2009
<i>AnClim</i>	Reference series	Annual, parallel monthly	HBS, moving window	Statistical test	Štěpánek et. al., 2009
<i>Craddock</i>	Pairwise Human synthesis	Serial monthly	Visual	Visual	Brunetti et. al., 2006
<i>RhtestV2</i>	Reference series or Absolute	Serial Monthly	Stepwise	Statistical test (modified Fisher)	Wang, 2008b
<i>SNHT</i>	Reference series	Annual	HBS	Statistical test (MLR)	Alexandersson and Moberg, 1997
<i>Climatol</i>	Reference series	Parallel monthly	HBS, moving window	Statistical test	Guijarro, 2011
<i>ACMANT</i>	Reference series	Annual Joint seasonal	DP	Penalized Likelihood	Domonkos et. al., 2011

In the methods we investigate, the homogenization corrections may be estimated directly on the comparison series (SNHT). When several references or pairwise estimates are available, a combination of those estimates is used. PRODIGE uses a decomposition of the signal into three parts: a common signal for all stations, a station dependent step function to model the

inhomogeneities and random white noise. In some methods, raw monthly estimates are smoothed according to a seasonal variation.

Once a first correction has been performed, most methods perform a review. If inhomogeneities are still detected, corrections with additional breaks are implemented in the raw series (examination: raw data), except in MASH where the corrected series receive additional corrections, until no break is found (called cumulative in table D2).

Table 5. Correction methods.

Method	Estimation	Review	Monthly correction
<i>MASH</i>	Smallest estimate from multiple comparisons	Examination: cumulative	Raw
<i>PRODIGE</i>	ANOVA	Examination: raw data	Raw
<i>USHCN</i>	Median of multiple comparisons	No review	Annual coefficients
<i>AnClim</i>	Estimated from comparison	Examination: raw data	Smoothed
<i>Craddock</i>	Mean of multiple comparisons	Examination: raw data	Smoothed
<i>RhtestV2</i>	Estimated on comparison	No review	Annual coefficients
<i>SNHT</i>	Estimated on comparison	Examination: raw data	Raw, ?
<i>ACMANT</i>	Estimated from comparison	No review	Smoothed
<i>Climatol</i>	Estimated from comparison	No review	Raw

The 25 submitted contributions are listed in Table 6. In case there are multiple contributions with the same algorithms, the ones denoted by “main” are the ones where the developer of the algorithm deployed it with typical settings.

Table 6. Names of contributions, contributors, purpose and references.

Contribution	Operator	Main purpose
<i>MASH main</i>	Szentimrey, Lakatos	Main submission
<i>MASH Marinova</i>	Marinova	Experienced user
<i>MASH Kolokythas</i>	Kolokythas	First-time user
<i>MASH Basic, Light, Strict and No meta</i>	Cheval	Experimental ¹
<i>PRODIGE main</i>	Mestre, Rasol, Rustemeier	Main submission ²
<i>PRODIGE monthly</i>	Idem	Monthly detection
<i>PRODIGE trendy</i>	Idem	Local trends corrected
<i>PRODIGE Acquaotta</i>	Acquaotta, Fratianni	First-time users
<i>USHCN main</i>	Williams, Menne	Produced USHCNv2 dataset
<i>USHCN 52x, cx8</i>	Idem	Alternatives for small networks
<i>AnClim main</i>	Stepanek	Main submission
<i>AnClim SNHT</i>	Andresen	SNHT alternative
<i>AnClim Bivariate</i>	Likso	Bivariate test in AnClim
<i>iCraddock Vertacnik, Klancar</i>	Vertacnik, Klancar	Two first-time users
<i>PMTred rel</i>	Viarre, Aguilar	PMTred test of RhTestV2
<i>PMFred abs</i>	Viarre, Aguilar	PMFred test, absolute method
<i>C3SNHT</i>	Aguilar	SNHT alternative
<i>SNHT DWD</i>	Müller-Westermeier	SNHT alternative
<i>Climatol</i>	Guijarro	Main submission
<i>ACMANT</i>	Domonkos	Main submission

¹ Experimental version that performs the four rules to combine yearly and monthly data separately, instead of the standard consecutive way.

² Detection: yearly; Correction: temperature monthly, precipitation yearly.

The data can be found on the ftp-server: ftp://ftp.meteo.uni-bonn.de/pub/victor/costhome/homogenized_monthly_benchmark/. In this directory there is a list with the subdirectories and the contributions they contain in the file

algorithms.txt. The current list is shown in Table 7. As it will be possible to submit new contributions for some time to make the comparison with the existing ones easier, further contributions may be added in future. In every subdirectory is a machine-readable file algorithm.txt with four lines: the short name of algorithm (for tables), the long name of the algorithm, the operator and a comment line.

Table 7. Names the directories and the contributions they contain.

Directory	Names blind contributions	Directory	Names late contributions
h001	MASH Kolokythas	h101	PMF OutMonthMean
h002	PRODIGE main	h102	PMF OutMonthQM
h003	USHCN 52x	h103	PMF OutAnnualMean
h004	USHCN main	h104	PMF OutAnnualQM
h005	USHCN cx8	h105	PMT OutMonthMean
h006	C3SNHT	h106	PMT OutMonthQM
h007	PMTred rel	h107	PMT OutAnnualMean
h008	PMFred abs	h108	PMT OutAnnualQM
h009	MASH Marinova	h201	PMF MonthMean
h010	Climatol	h202	PMF MonthQM
h011	MASH main	h203	PMF AnnualMean
h012	SNHT DWD	h204	PMF AnnualQM
h013	PRODIGE trendy	h205	PMT MonthMean
h014	AnClim Bivariate	h206	PMT MonthQM
h015	ACMANT	h207	PMT AnnualMean
h016	iCraddock Vertacnik	h208	PMT AnnualQM
h017	iCraddock Klancar	h301	PRODIGE Automatic
h018	AnClim main	h302	AnClim late
h019	AnClim SNHT	h303	Climatol2.1a
h020	PRODIGE Acquaotta	h304	Climatol2.1b
h021	PRODIGE monthly	h305	Craddock late
h022	MASH Basic	h306	ACMANT late
h023	MASH Light		
h024	MASH Strict		
h025	MASH No meta		
inho	Inhomogeneous data		
orig	Original homogeneous data		

There are two special directories: inho and orig. The directory inho contains the inhomogeneous data, which was available to the participants for homogenisation. The directory orig contains the original homogeneous data, i.e. the data before the inhomogeneities have been inserted. This directory can also be seen as a contribution that found and corrected all breaks perfectly; it therefore also contains the corrected-files, which indicate the correct positions of all breaks, local trends and outliers.

In the directories named hnnn you can find the homogenised data returned by the participants. Directories with a number nnn below 100 are the official blind contributions that have been submitted before the deadline and thus before the truth was revealed. These contributions are detailed in Section 3.2. The other late contributions are described in Section 3.3 and should be interpreted with more care; see the article on the benchmark dataset.

3.2 Blind contributions

3.2.1. MASH

The original MASH (Multiple Analysis of Series for Homogenization) method has been developed for homogenization of climate monthly series based on hypothesis testing (Szentimrey, 1999). It was created with the intention to provide a user-friendly software package based on a mathematically well-established procedure.

MASH is a relative and iterative homogeneity test procedure that does not assume homogeneity of the reference series. Possible break points can be detected and corrected through comparisons with multiple reference series from subsets of stations within the same climatic area. This area can be limited for large networks based on the geographical distance, but this was not needed for the benchmark. From the set of available series, a candidate series is chosen and the remaining ones are treated as potential reference series. During one iterative step, several difference series are constructed from the candidate and various weighted reference series. The optimal weights are determined according to the kriging interpolation formula that minimizes the variance of the difference series. Provided that the candidate series is the only one with a break in common with all the difference series, these break points can be attributed to the candidate series.

Depending on the distribution of the examined meteorological element, an additive (e.g. temperature) or a multiplicative (e.g. precipitation) model can be applied, where the second case can be transformed into the first one.

MASH computes not only estimated break points and shift values, but the corresponding confidence intervals as well. The confidence intervals make it possible to use the metadata information - in particular the probable dates of break points – for automatic correction. For this, MASH uses three different decision rules, which are used consecutively and are described in more detail in Section 3.2.1.4.

We emphasize that such an iteration step including comparison, detection, attribution and correction is fully automatic. The outlier detection and missing value completion based on kriging interpolation are also part of the system. The latest version MASHv3.02 was used, which also handles daily data (Szentimrey, 2007, 2008).

3.2.1.1. MASH MAIN

For the benchmark homogenization, some modifications have been made to the MASH system. The first one is to start with a preliminary examination of the annual series. Normally MASH would start with the monthly series and later examine the seasonal and annual series. The new possibility is to begin with an examination of annual series and to use the detected breaks as preliminary information (metadata) for the standard application of MASH for monthly data. This information can be used as if it were metadata according to the three decision rules.

The second modification is implemented because of the large amount of missing data in the beginning period (1900-1924). In order to solve this problem, the homogenization procedure is divided into two parts. First, for the later period (1925-1999) mainly automatic homogenization and quality control procedures are applied. Second, the original (1900-1924) and the resulted homogenized (1925-1999) data series are merged and automatic missing data completion and cautious interactive homogenization procedures are applied to the full series.

3.2.1.2. MASH MARINOVA

This contribution uses the same version of MASHv3.02 described in Section 3.2.1.1, but has another operator. The final results are influenced by the subjective decisions taken during the execution of the homogenization procedure, especially for the period 1900 – 1924, and differ to some extent from those obtained in contribution 3.2.1.1.

3.2.1.3. MASH KOLOKYTHAS

This contribution has used the software MASHv3.02 (Szentimrey, 2007, 2008), without the two modifications described in section 3.2.1.1. The followed procedure consists of three steps. First, the missing values are filled for all examined series. Second, outliers and breaks-shifts are identified and corrected in the monthly series and later in the seasonal and annual series, without using any metadata. Third, a verification of the homogenization is performed; both on the actual and the final stage of homogenization.

3.2.1.4. MASH BASIC, LIGHT AND STRICT

Just as for the contribution 3.2.1.1 this contribution first homogenized the yearly series and then used these break points as metadata for the monthly homogenization. MASH combines metadata with the monthly breakpoints by taking into account their confidence intervals using three different rules. Normally these three rules are applied consecutively, but here they are applied on their own to examine the differences.

Using the strict rule, only breaks given by metadata are corrected. The basic rule computes the date of the break as a compromise between the one in the metadata and the break point in the difference series (Szentimrey, 2007, p. 10). The light rule picks the date in the metadata, if available, i.e. metadata always have priority if it is possible on the basis of confidence intervals.

3.2.1.5. MASH NO META

The homogenization was performed on surrogate temperature networks one to four, using the standard version of MASH v.3.02, i.e. without the two changes described in 3.2.1.1.. The homogenization was performed fully automatically on the datasets as provided, i.e. no yearly homogenization was performed to provide metadata for the monthly homogenization.

3.2.2. *PRODIGE*

The detection method of *PRODIGE* relies on pair-wise comparison of neighbouring series. These difference series are then tested for discontinuities by means of an adapted penalized likelihood criteria (Caussinus and Lyazrhi, 1997). If a detected change-point remains constant throughout the set of comparisons of a candidate station with its neighbours, it can be attributed to this candidate station. This attribution phase requires human intervention.

The corrections are computed by decomposing the signal into three parts: a common signal for all stations, a station dependent step function to model the inhomogeneities and random white noise (Caussinus and Mestre, 2004). The standard least squares procedure is used for estimating the various parameters. Estimated corrections within a series are the difference between its most recent station effect and the actual one.

3.2.2.1. *PRODIGE* MAIN

In this first contribution, the detection has been performed annually, without checking the month of changes. Furthermore, all series within each network constitute the neighbourhood and determine the regional climate part in the correction algorithm. Temperature corrections are estimated monthly, while precipitations are corrected using an annual coefficient.

3.2.2.2. *PRODIGE* MONTHLY

In the second contribution, the month of change is determined as well (in the first contribution the break is always at the end of the year). Moving neighbourhoods to compute the regional climate are used in the correction algorithm, using series with a correlation greater than 0.7, (sometimes it is lower when less than six series meet this threshold). The precipitation series are corrected using monthly coefficients.

3.2.2.3. PRODIGE TRENDY

In the third contribution is similar to the first one, but artificial local trends were corrected as well; their detection is currently still mainly manual.

3.2.2.4. PRODIGE ACQUAOTTA

In this contribution, the surrogate temperature networks have been homogenized.

The outliers have been identified at monthly scale. The lower threshold for outlier values is given by: $q_{25}-c*(q_{75}-q_{25})$, and the upper threshold by: $q_{75}+c*(q_{75}-q_{25})$; where q_{25} , q_{75} are the 25th and 75th quantiles and $c=1.5$. This procedure was applied by using the AnClim program (Štěpánek, 2005) and was chosen because the identification of the outliers by PRODIGE is too subjective.

Then the detection was performed on annual values and the break was identified always at the end of the year. For the choice of reference series two methods have been used. For the networks with only five stations all series were used as reference series while for the networks with more than five stations, only four series with a correlation greater than 0.80 were used. The reference series constitute the neighbourhood for both detection and correction purposes and determine the regional climate part in the correction algorithm. Temperature corrections are estimated monthly.

3.2.3. USHCNv2

In this approach (Menne and Williams, 2009), automatic comparisons are made between numerous combinations of temperature series in a region to identify cases in which there is an abrupt shift in one station series relative to many others in the region. The algorithm starts by forming a large number of pairwise difference series between serial monthly temperature values from a sub-network of stations. Each difference series is then statistically evaluated for abrupt shifts using the SNHT with a semi-hierarchical splitting algorithm. After all breaks are resolved within the combinations of difference series, the series responsible for each particular break is identified automatically through a process of elimination.

Adjustments are determined by calculating multiple estimates of each shift in the difference series using segments of a differences series formed with highly correlated, neighbouring series that are homogeneous for at least 24 months before and after the target change-point. When two change-points occur within 24 months in the target series, an adjustment is made for their combined effect. The range of pairwise shift estimates for a particular step change is treated as a measure of the confidence with which the magnitude of the discontinuity can be estimated. At least three estimates for the shift are required to determine statistical significance. If it is not possible to calculate three estimates of the magnitude of the shift, no adjustment is made; otherwise, the median adjustment is implemented.

3.2.3.1. USHCN MAIN

This version is the current standard and is used to produce the USHCN version 2 monthly temperature data (Menne et al., 2009).

3.2.3.2. USHCN 52X

52X is the same as the main one except that the process of elimination used to identify the series responsible for a particular break in the set of pairwise difference series has been altered to work more effectively in networks with the small number of stations found in the benchmark dataset. In 52d, at least two difference series formed between a target and its neighbours must share the exact same change-point date to allow for the shift to be appropriately identified as having been caused by the target series whereas in 52x the change-points in the target-neighbours differences are only required to be nearly coincident.

3.2.3.3. USHCN CX8

CX8 is identical to the main one except that an adjustment is made when all pairwise estimates of the target shift are of the same sign (i.e., there is no Tukey-type significance test used to evaluate the distribution of the pairwise shift estimates as described in Menne and Williams 2009).

3.2.4. AnClim

The homogenization is performed with a software package consisting of AnClim (Štěpánek, 2008), LoadData and ProClimDB (Štěpánek, 2009). The complete package offers database functionality, data quality control and homogenization, as well as time series analysis, extreme values evaluation, model outputs verification, etc.

Detection of inhomogeneities is carried out on a monthly scale applying AnClim while quality control, preparing reference series and correction of found breaks is carried by ProClimDB software. The software combines many statistical tests, types of reference series (calculated from distances or correlations, applying one reference series or pair-wise one) and time scales (monthly, seasonal and annual). All of these can be used to create an “ensemble” of solutions, which may be more reliable than a single method (if the detection results coincide).

To be able to test the large number of the networks in reasonable time, the program runs fully automatically, without need of any user intervention. Normally one would use the software results only as a first step for decision-making, and additionally consider metadata (if available), inspect the difference plots, etc. The obtained results for this study should thus be regarded as the worst possible and would be improved by human analysis. Furthermore, quality control is normally performed on daily (or sub-daily) data in which problems are more evident and correction of inhomogeneities is also applied in daily (sub-daily) data.

3.2.4.1. ANCLIM MAIN

Version 5.014 of AnClim and 8.358 of ProClimDB are used to homogenize this benchmark contribution. Outliers are rejected automatically only in the most evident cases if there is coincidence among several characteristics (settings are taken from daily data processing, details can be found in Štěpánek et al., 2009).

For the computation of the reference the limit for correlations, ρ , is set to 0.35 (median of monthly correlations computed using the first differenced series), while the limit for distances (altitude) is set to a maximum of 400 km (500 m), but maximally 5 neighboring stations are used. Two reference series are calculated as weighted average, one using correlations and one reciprocal distances, $1/d$. The weights are ρ^{pw} and d^{pw} , with the power of the weights pw set to 1 for temperature and 3 for precipitation. Neighbour stations values were standardized to average and standard deviation of the tested series prior weighted average calculation.

To homogenize the benchmark data set three tests are applied: SNHT (Alexandersson, 1986), the bivariate test of Maronna and Yohai (1978) and a t-test. The time scales of tested series are months, seasons and year.

A break in the ensemble is accepted, and thus corrected, if there are at least 15% hits in the same year of a given station. In this computation of the number of hits, breaks found in the monthly series have unity weight, breaks in the seasonal series a weight of 2 and in the annual series a weight of 5.

Correction has been performed using data from 10 years before and after a change (by comparison with the reference series). The monthly adjustments are smoothed with a low-pass filter with weights 0.479, 0.233, 0.027 to each side. Detection and correction is performed in

two iterations. Only breaks which are minimally 2 years from the edge of the series are corrected and the correction is only implemented if it improves the cross-correlation.

3.2.4.2. ANCLIM SNHT

In this contribution, the SNHT method is used as implemented in version 5.012 of AnClim and version 8.397 of ProClimDB is used. The criteria for the homogeneity test are the same as described in 3.2.4.1 with the following exceptions. The limit in the quality control for the distance between neighbours is 200 km and the difference in altitude is 800 m. All proposed outliers are accepted.

The reference time series for detection is based on neighbours with a correlation larger than 0.80. The reference is computed using the reciprocal distance as weights, with the power of weights set to 0.5. Only one detection test is used, namely SNHT.

For calculating the size of the breaks, the same reference series as for detecting the breaks are used. The procedure of the homogeneity analysis of the adjusted time is run only once, i.e. no iterations are performed. Only breaks with a minimum of 4 years from the edge of the series (or from the nearest break point) are adjusted. Furthermore, the correction is only applied if the correlation between the test series and the reference series has improved with more than 0.5 %.

3.2.4.3. ANCLIM BIVARIATE

In this contribution, the reference series for homogeneity testing is calculated by ProClimDB software v8.309 based on correlations. For this purpose, five neighbouring stations with the same settings described in Section 3.2.4.1 are used. The length of a reference series is 40 years and the overlapping period is 10 years as recommended by Štěpánek et al. (2009).

Using the version v5.012 of AnClim, each station is tested manually on monthly basis with the bivariate test. Adjustments are performed using data from 10 years before and after a change and are based on comparison with the reference series. The procedure is repeated several times (2 to 3 iterations are needed) in order to obtain series without significant break points, i.e. homogeneous series.

3.2.5. Craddock

The Craddock test (Craddock, 1979) is a relative test that accumulates normalized differences between two series according to the formula, which is based on Schönwiese and Malcher (1985),

$$s_i = s_{i-1} + \frac{\bar{t}}{\bar{r}} r_i - t_i,$$

where i is a running index, t_i denotes the series to be tested and r_i the reference series; \bar{t} and \bar{r} represent the mean over the whole period; temperature is expressed in Kelvin.

The ideal Craddock graph for a homogeneous series should be a straight line with $s_i = 0 \forall i$ and in realistic situations, a curve with no systematic deviations from zero. The typical signal of an inhomogeneity is a discontinuity in the first derivative of the Craddock function.

The homogenization approach is the same discussed in Brunetti et al. (2006). It consists of a revisited version of the HOCLIS procedure (Auer et al., 1999). HOCLIS rejects the a priori existence of homogeneous reference series. Instead of using one single reference series (e.g. a weighted average series of neighbouring stations), a multiple pairwise Craddock test is performed, where each series is tested against other series. The test is based on the hypothesis of the constancy of temperature differences and precipitation ratios. The break signals of one series against all others are then collected in a decision view graph and the breaks are assigned to the single series according to the estimated probability. This system avoids trend

imports and an inadmissible adjustment of all series to one or a few “homogeneous reference series”.

Once a break is identified, the adjustments are evaluated from those reference series that result homogeneous in a sufficiently long sub-period centred on the break year, and that well correlate well with the candidate one. The monthly adjustments from each reference series are then fitted with sine and cosine functions (which can be chosen to be between 1 and 4 harmonics) to smooth the noise and to extract only the physical signal.

Especially when discontinuities are not very high, the signal in the homogeneity test may not be very clear and the researcher’s choices will be influenced by his working philosophy and skill. Thus, the performance is strongly user dependent.

3.2.5.1. iCRADDOCK VERTACNIK

The software package iCraddock is based on the software written by Michelle Brunetti, but includes a graphical user interface. In this application the number of pair tested does not have an upper limit. Furthermore, in most cases sine and cosine functions are fitted to the monthly adjustments to account for an annual cycle in the inhomogeneity.

3.2.5.2. iCRADDOCK KLANCAR

This contribution uses the same software and settings as described in Section 3.2.5.1, but has another operator.

3.2.6. *RhTestV2*

The RHtestV2 software package provides a user friendly graphical interface for detecting and adjusting artificial mean-shifts in data series. Relative and absolute homogenization can be performed. Both approaches use different statistical tests, as described below.

The critical values take the effect of short-range correlations into account and correct for over-detection near the end-points of the time series as well. Both statistical tests are constructed using autoregressive Gaussian noise (Wang, 2008b). If metadata is available, expected dates of change-points can be indicated.

Result-files document the significance of the change-points. In addition to the date (month and year) of the change-point, it details the p-values and the PF_{\max} or PT_{\max} statistics (for absolute and relative detection respectively), followed by their 95th percentiles and their 95% uncertainty range. The user has to set a criterion to consider a change-point significant and can delete the least significant change-points. In this case, the algorithm has to be relaunched to assess the significance of remaining change-points. This operation has to be performed until all the retained breaks are significant. No outlier detection and removal was applied. To estimate the magnitude of the shifts, the algorithm uses a multi-phase regression model (see appendix in Wang (2008b) for more details).

3.2.6.1. PMTRED REL

Relative homogenization uses the penalized maximal t-test of the software package RHTESTV2 (Wang et al., 2007) applied over separately computed references, as the software package does not include a tool for their computation or recommendations on how to do so. RhTestyV2 assumes that there are no trends in the difference time series. If we expect an artificial shift in mean, the time series can be considered as two independent samples from two normal distributions with the same unknown variance, but different means (Wang et al., 2007). The position of the change-point is associated with the maximal value of the log likelihood ratio (see Csörgo and Horváth (1997) for details).

For each time series and for each iteration, a reference is computed by weighting the best-correlated series by their correlation (Pearson's method). The minimum correlation to consider

a series in the computation of the reference is set to 0.5. If less than three stations have a correlation higher than 0.5, the three best correlated series are used. This method can be robust for monthly series of temperature (where regional effects are less pronounced), but depends on the network density. For precipitation data, local climatic factors can affect the correlations between the nearest stations. Results have to be interpreted carefully. All time series are detrended before applying the algorithm.

3.2.6.2. PMF_{FRED} ABS

Absolute homogenization (without reference) is carried out using a penalized maximal F-test of the software package RHTESTV2 (Wang, 2008b). The time series can have a constant linear trend during the whole period. The most probable change-point is the one that maximizes the F criterion (Wang, 2003). Corrections for precipitation and temperature time series are similarly estimated.

3.2.7. SNHT

The Standard Normal Homogeneity Test was first described by Alexandersson (1986) as a likelihood ratio test. It is performed on a difference series (temperature, or in general additive variables) or ratio series (precipitation, or in general multiplicative variables) between the candidate station (i.e. the site to be homogenized) and a reference series. The reference is usually calculated as a weighted average of well-correlated neighbouring series sharing the same climate signal. First the difference or ratio series is normalized by subtracting the mean and dividing by the standard deviation. In its simplest form, a single break is detected at the maximum of:

$$T_v = v\bar{z}_1^{-2} + (n-v)\bar{z}_2^{-2}$$

where T_v is the test statistic for position v , and \bar{z}_1 is the mean for the series from the first data point to v and \bar{z}_2 is the mean of the series from $v+1$ to the end, i.e. index n .

The test was reformulated to hierarchically detected multiple break points by testing the sub-series defined after a break-point is found in Alexandersson and Moberg (1997); AM97 from now on. The procedure described there considers each time series as potentially containing multiple breaks and tests iteratively and several times the series to determinate significant change-points. Additional tests were provided in AM97 to allow for different variances before and after v and to try to identify artificial trends in the series, instead of sudden jumps. Critical values of the statistics are provided from Monte-Carlo simulations, although in recent years some authors have discussed their adequacy, providing new critical levels after more comprehensive simulations (i.e., Khaliq and Ouarda, 2007)

3.2.7.1. C3SNHT

The implementation of SNHT presented here is an automatic version of the AM97 procedure. For each station a specific network is created, meeting the requirements of containing the optimum reference stations for that particular site (main candidate station, MCS), according to a combination of data overlap, geographical distance and correlation. This network is homogenized through a full AM97 procedure, but only results for the MCS are retained. A separate AM97 procedure is run for each station in the network. Segmentation is done in each test by identifying the most recent break and retesting the remaining data until the segment is too small. For the benchmark, detection is performed on the annual series and all the identified significant breaks are corrected in to monthly values. In a real application, this quick procedure would be completed with the consideration of the seasonal series and different variables (i.e., if t_{\max} and t_{\min} are available, the test would be also applied to t_{mean} and DTR to derive a solid correction pattern, as described in Brunet et al., 2008). In addition, an unavoidable metadata analysis would be performed to prevent contamination by breaks in the

reference time series. For each network, the specific settings varied according to the number of stations and missing values – for most cases the outlier threshold was set to 4 times the standard deviation during detection (and 3.5 during correction), the significance level of breaks to 90% during detection (and all breaks are corrected in the correction phase) and the minimum segment length was set to 9 years in most cases.

3.2.7.2. SNHT DWD

The test is a simple version of SNHT, which can only detect breaks, i.e. it does not consider local trends nor does it detect outliers. It is applied independently to the time series of the different months with an error level of 10 %.

The reference time series is derived from the ten best neighbouring stations. For the selection of these stations, limits for the horizontal distance from the station and for the difference in elevation can be defined in dependence on the climatological parameter tested.

As the results in the individual months may differ in date and correction factor for the detected inhomogeneities as well as for the fact that any inhomogeneity is detected at all, the results have to be processed manually. For that purpose usually metadata with information concerning the station history are used. However, for this formal test of procedures, only gaps in the results from month to month have been filled and homogeneities detected in single months only have been deleted. Furthermore, checks to achieve a fixed date for the inhomogeneities in all months and a smooth seasonal variation of correction factors have been applied.

The procedure is normally used for a large network of stations that have already been checked for random errors and outliers, so that only climatologically reasonable data can be expected. Therefore, the results for the rather small station networks with extreme outliers in the benchmark date may probably be rather weak and sometimes detect “inhomogeneities” which are only caused by the inhomogeneous reference series.

3.2.8. *Climatol*

The homogenization function of the R package "Climatol" begins by normalising all monthly series in parallel; temperature is normalized by subtracting its mean and dividing by its standard deviation, precipitation is first transformed with a square root and then normalized by dividing by its mean. As series are often incomplete, the mean and standard deviation of the whole study period are unknown, and therefore they must be computed iteratively until getting stable values. In the beginning of every iteration, a reference time series is computed as a weighted average of the closest neighbours. The weights (w) are determined by a bell shaped function of the distance (d) by the equation $w = (1 + (d/b)^2)^{-1}$, where the constant b is 100 km for detection and 10 km for computing the corrections. For detection, the nearest six stations with data on a specific date are used; for correction, the nearest four stations with data.

The difference between the observed and reference series are used in search for outliers and breaks. Outliers greater than a prescribed threshold are deleted; the default threshold is four times the standard deviation, but in three networks this threshold needed to be increased to preserve a sufficient number of data points. Break detection is performed on the serial monthly time series based on a test of squared relative mean difference (SRMD) in a running window of 48 months (4 years) before and after the break. The threshold for SRMD was 53.9, which is equivalent to a significance level of about 10^{-7} . Part of the series is split off at the most significant break, forming a new series, which is treated as an independent new station located at the same position. The whole procedure is repeated until no more breaks are found. In a final stage, all missing data are filled. Although the function admits different parameterizations, only its default values were used here.

3.2.9. ACMANT

The Adapted Caussinus-Mestre Algorithm for Networks of Temperature series (ACMANT) has been developed for the automatic homogenization of temperature data with monthly resolution. In an iterative procedure, one series of the network is corrected in every step.

Reference series are built from the other time series of the same climatic network, weighting them with the squared spatial correlations among the series of consecutive differences (Peterson and Easterling, 1994). When different component series are available for sections of the candidate series, different reference series can be applied for different parts of the same candidate series. New reference series are built after each step of homogenization or outlier filtering.

The detection is performed using bivariate statistics on the difference time series of both the annual means as well as summer-winter differences, with step-functions to model the breaks. The common minimum of the sum of squared errors and the Caussinus-Lyazhri criterion (*Caussinus and Lyazrhi, 1997*) indicates the timing of the inhomogeneities (IHs). For the timing, first, the year is calculated and only later the month is determined, by examining 10 months moving windows around the December of the selected year. The IH-search is performed for each time series of the network, but the corrections are applied only to the one with the largest IHs. Monthly correction factors are modelled as an annual harmonic function and always applied relative to the pre-homogenized values.

When accumulated anomalies for subperiods of 5-10 months remain high after correction, a secondary univariate homogenization on monthly values is accomplished on subperiods of 60 months around the maximum of accumulated anomalies. Maximally two IHs are detected for a 60-month period. After this step the corrections for each break of the whole series are recalculated. The method interpolates data gaps before homogenization, and filters outliers both in the beginning as well as at the end of the homogenization procedure.

3.3 Late contributions.

3.3.1 PMF and PMT

After the deadline 16 surrogate temperature contributions similar to the blind contributions *PMTred rel* and *PMFred abs* were produced, but with the detection and correction functions from the new software package RhTestV3. In half of the contributions the, now known, outliers have been removed in advance, in the others not (marked as *OUT*). Furthermore, half of the contributions corrected monthly (*Month*) and the other half yearly values (*Annual*); half did so correcting the mean values (*Mean*), half with quantile matching (*QM*).

3.3.2. PRODIGE automatic

An automatic version of PRODIGE was developed by Elke Rustemeier. This late contribution is similar to *PRODIGE main*, but the interpretation of the change points, found by the comparison with a group of neighbor stations, is performed automatically. It computes a weighted mean number breaks per year, based on the cross-correlations between the stations. The decision to accept a break depends on thresholds, which were found by training on the first two precipitation networks. Because this contribution was trained on part of the benchmark dataset, these errors may not be representative.

3.3.3. AnClim late

The main differences between the blind *AnClim main* and the new contribution *AcClim late* are as follows. In *AnClim main*, distances for finding the neighbor stations were used and from the neighbors one reference series using weighted (reciprocal distance) average was created. In the new contribution, correlations for selecting the neighbors are used and the inhomogeneity testing was performed as pair-wise comparison, i.e. repeating for each neighbor individually.

There is also difference in finding outliers (before running homogenization). Originally the same settings as for daily data were applied (inhomogeneities correction and quality control using AnClim and ProClimDB software is tuned for daily data), now both quality control and inhomogeneities correction was tuned also for monthly data.

3.3.4.8. Climatol2.1

Climatol's blind contribution showed good results for detection, but strongly reduced the trends. After the deadline two new contributions, *Climatol2.1a* and *Climatol2.1b*, were submitted. The important changes are as follows. Precipitation data are first transformed by applying a cubic root to the data, instead of a square root.

The homogenization in *Climatol2.1a* begins by normalizing all series by subtracting (temperature) or dividing by the mean (precipitation) of serial monthly data, while in the blind contribution the means were computed for every month individually. In *Climatol2.1b* the normalization is performed by subtracting the mean and dividing by the standard deviation (full standardization), and then estimating all normalized data as an average of their closest neighbours (up to 10 by default, if available).

As series are often incomplete, the means of the whole period are unknown, and therefore they must be computed iteratively until getting stable values. The new stopping criterion for the iterations is stricter.

Furthermore, the test of the squared relative mean difference was replaced by the SNHT test. To avoid misleading results due to multiple change points, SNHT is applied first on 120 month windows moved forward in steps of 60 months, and in a second stage it is applied on the whole series (the traditional way). The whole process is repeated until no SNHT value is greater than a prescribed threshold, and in a third stage all missing data (original and resulting from the series splits) are computed again (this time using an inverse distance weighted average).

In addition the maximum number of series to compute the references series is increased to ten, the reference series, for shift detection, are computed as the simple mean of all these surrounding stations, and the homogenised series have been built on the longest homogenous subperiod, in stead of the first one, as in the blind contribution. Outliers are removed if they are greater than a preset number of standard deviations (5 by default, although 7 or 8 have been used for the precipitation series).

6.2.2. *Craddock late*

The two very good *iCraddock* contributions were submitted by users knowledgeable about homogenization, but new to the algorithm. After the deadline we have been able to obtain a contribution by Michele Brunetti, who is an experienced Craddock user. This contribution might still be considered blind as Brunetti declares not to have looked at the original homogeneous data.

Four noteworthy strategies were used. Firstly, the most relevant pairs of stations are selected not only based on correlation, but for climatological similarity, e.g. exposure. Secondly, often only a part of the homogeneous subperiod is used for correction. Thirdly, also breaks that are not clearly evident are corrected. Finally, depending on the strength of the seasonal cycle of the break, the operator selects annual or monthly corrections.

6.2.1 *ACMANT late*

ACMANT is a new method. The development of the method is continued after the blind test experiment. The main adjustments of this new contribution *ACMANT late* are as follows. Instead of the step-by-step correction of the network data, a pre-homogenization is applied. In this step series are not used to compute a reference that later are a candidate series to avoid biases. Furthermore, *ACMANT late* applies the decomposition model of PRODIGE (Sect. 3.2.2) for the final adjustment. Further small modification have been made relative to the early version. The full description of the method can be found at <http://www.c3.urv.cat/members/softpeter.html>.

ACMANT late is optimized based on the benchmark data itself. It is thus not clear how much of this performance is due to overfitting and may not be fully realized in an application to a real dataset. However, those results are promising.