

Climate Data and Monitoring  
WCDMP-No. 84



**Eighth Seminar for Homogenization and  
Quality Control in Climatological Databases  
and Third Conference on Spatial Interpolation  
Techniques in Climatology and Meteorology**

(Budapest, Hungary, 12-16 May 2014)



**World  
Meteorological  
Organization**  
Weather • Climate • Water

**© World Meteorological Organization, 2014**

The right of publication in print, electronic and any other form and in any language is reserved by WMO. Short extracts from WMO publications may be reproduced without authorization, provided that the complete source is clearly indicated. Editorial correspondence and requests to publish, reproduce or translate this publication in part or in whole should be addressed to:

Chair, Publications Board  
World Meteorological Organization (WMO)  
7 bis, avenue de la Paix  
P.O. Box 2300  
CH-1211 Geneva 2, Switzerland

Tel.: +41 (0) 22 730 84 03  
Fax: +41 (0) 22 730 80 40  
E-mail: [Publications@wmo.int](mailto:Publications@wmo.int)

NOTE

The designations employed in WMO publications and the presentation of material in this publication do not imply the expression of any opinion whatsoever on the part of the Secretariat of WMO concerning the legal status of any country, territory, city or area, or of its authorities, or concerning the delimitation of its frontiers or boundaries.

Opinions expressed in WMO publications are those of the authors and do not necessarily reflect those of WMO. The mention of specific companies or products does not imply that they are endorsed or recommended by WMO in preference to others of a similar nature which are not mentioned or advertised.

**EIGHTH SEMINAR  
FOR HOMOGENIZATION AND QUALITY CONTROL IN  
CLIMATOLOGICAL DATABASES**

**AND**

**THIRD CONFERENCE  
ON SPATIAL INTERPOLATION  
TECHNIQUES IN CLIMATOLOGY AND METEOROLOGY**

**Budapest, Hungary, 12 – 16 May 2014**

**Organized by the Hungarian Meteorological Service (OMSZ)**

**Supported by WMO and OMSZ**

**Edited by Mónika Lakatos, Tamás Szentimrey, Annamária Marton**

## PREFACE

Homogenization of climate data series and spatial interpolation of climate data play a growing role in the meteorology and climatology. The data series are usually affected by inhomogeneities due to changes in the measurement conditions (relocations, instrumentation) therefore a direct analysis of the raw data series can lead to wrong conclusions about climate change. Reconstruction of meteorological fields and gridded databases require spatial interpolation methods.

The first seven Seminars for Homogenization and Quality Control in Climatological Databases as well as the first two Conferences on Spatial Interpolation Techniques in Climatology and Meteorology were held in Budapest and hosted by the Hungarian Meteorological Service.

The 7th Seminar in 2011 was organized together with the final meeting of the COST Action ES0601: Advances, in Homogenization Methods of Climate Series: an integrated approach (HOME), while the 1st Conference on Spatial Interpolation was organized in 2004 in the frame of the COST Action 719: The Use of Geographic Information Systems in Climatology and Meteorology. Both the seminar and the conference series were supported by WMO.

The 8th Homogenization Seminar and the 3rd Conference on Spatial Interpolation were organized together considering certain theoretical and practical respects. Theoretically there is a strong connection between these topics since the homogenization and quality control procedures need spatial statistics and interpolation techniques for spatial comparison of data. On the other hand the spatial interpolation procedures (e.g. gridding) need homogeneous data series with high quality. Practically the CARPATCLIM project that was launched in 2010 and ended in 2013 is a good example for this problem. The main purpose of the project was to produce a gridded database for the Carpathian region based on homogenized data series. The experiences of this project may be useful for the implementation of gridded databases.

The Organizers

## CONTENTS

PREFACE .....	2
CONTENTS .....	3
MATHEMATICAL QUESTIONS OF HOMOGENIZATION AND QUALITY CONTROL.....	
<b>Tamás Szentimrey, Mónika Lakatos, Zita Bihari</b> .....	5
IRELAND WITH HOMER .....	
<b>John Coll, Mary Curley, Séamus Walsh, John Sweeney</b> .....	23
THE ACMANT2 SOFTWARE PACKAGE .....	
<b>Peter Domonkos</b> .....	46
HOMOGENIZATION OF MONTHLY TEMPERATURE SERIES IN ISRAEL - AN INTEGRATED APPROACH FOR OPTIMAL BREAK-POINTS DETECTION .....	
<b>Yizhak Yosef, Isabella Osetinsky-Tzidaki, Avner Furshpan</b> .....	73
THE WMO/MEDARE INITIATIVE: BRINGING AND DEVELOPING HIGH-QUALITY HISTORICAL MEDITERRANEAN CLIMATE DATASETS INTO THE 21 <sup>ST</sup> CENTURY ...	
<b>Khalid Elfadli and Manola Brunet</b> .....	86
HOMOGENIZATION OF SPANISH MEAN.....	
<b>José A. Guijarro</b> .....	98
MATHEMATICAL QUESTIONS OF SPATIAL INTERPOLATION .....	
<b>Tamás Szentimrey, Zita Bihari, Mónika Lakatos</b> .....	107
PRACTICAL ASPECTS OF RAW, HOMOGENIZED AND GRIDDED DAILY PRECIPITATION DATASETS.....	
<b>Predrag Petrović, Gordana Simić, Ivana Kordić</b> .....	115
HOMOGENIZATION OF MONTHLY AIR TEMPERATURE .....	
AND MONTHLY PRECIPITATION SUM DATA SETS .....	
COLLECTED IN UKRAINE .....	
<b>Skrynyk O., Savchenko V., Radchenko R. and Skrynyk O.</b> .....	128
HOMOGENIZATION PROCESS IN THE CLIMATE OF CARPATHIAN REGION PROJECT, VERIFICATION RESULTS.....	
<b>M. Lakatos, T. Szentimrey, Z. Bihari, and S. Szalai</b> .....	134
BIASES AND CORRECTIONS OF WIND SPEED TIME SERIES .....	
<b>Csilla Péliné Németh, Judit Bartholy, Rita Pongrácz,</b> .....	
<b>Tamás Szentimrey, Kornélia Radics</b> .....	151

PROGRAMME ..... 160  
LIST OF PARTICIPANTS ..... 165

# MATHEMATICAL QUESTIONS OF HOMOGENIZATION AND QUALITY CONTROL

**Tamás Szentimrey, Mónika Lakatos, Zita Bihari**

Hungarian Meteorological Service  
szentimrey.t@met.hu

## **Abstract**

There are several methods and software for the homogenization of climate data series but there is not any exact mathematical theory of the homogenization. At the examinations mainly the physical experiences are considered while the mathematical formulation of the problems is neglected in general. Moreover occasionally there are some mathematical statements at the description of the methods in the papers – e. g. capability to correct the higher order moments – but without any proof and this way is contrary to the mathematical conventions of course. As we see the basic problem of the homogenization is the unreasonable dominance of the practical procedures over the theory and it is the main obstacle of the progress. Therefore we try to formulate some questions of homogenization in accordance with the mathematical conventions. The planned topics to be discussed are as follows.

- The mathematical definition of the inhomogeneity and the aim of homogenization. It is necessary to clarify that the homogenization of climate data series is a distribution problem instead of a regression one.
- Relation of monthly and daily data series homogenization.
- Mathematical overview on the methodology of spatial comparison of series, inhomogeneity detection, correction of monthly series.
- Relation of theoretical evaluation and benchmark for methods, validation statistics.

## **1. INTRODUCTION**

First let us consider the abstract schema of the meteorological examinations. The initial stage is the meteorology that means the qualitative formulation of the given problem. The next stage is the mathematics in order to formulate the problem quantitatively. The third stage is to develop software on the basis of the mathematics. Finally the last stage is again the meteorology that is the application of the developed software and evaluation of the obtained results. In the practice however the mathematics is sometimes neglected.

Concerning our topic we have the following question. What is the mathematics of homogenization in meteorology? There are several methods and software for the homogenization of climate data series but unfortunately there does not exist any exact well elaborated mathematical theory of this problem. At the climatological examinations mainly the physical experiences are dominated while the mathematical formulation of the problems is neglected in general. We do not argue the importance of the physical aspects but the applied not too advanced mathematics is in contrast with the fact that the methods are declared to be based on the mathematical statistics. Moreover often there are some mathematical statements at the description of the methods in the papers – e. g. capability to correct the higher order moments – but without any proof and this way is contrary to the mathematical conventions of

course. As we see the basic problem of the homogenization is the unreasonable dominance of the practical procedures over the theory and it is the main obstacle of the progress. As a consequence of this practice the exact evaluation of the methods is also very problematic or properly speaking it is unrealistic and the progress of the homogenization research activity is doubtful. Therefore we try to provide a general approach for the mathematical formulation of homogenization in accordance with the mathematical conventions. We believe the correct mathematical principles can promote understanding and clarifying the questions of homogenization in climatology.

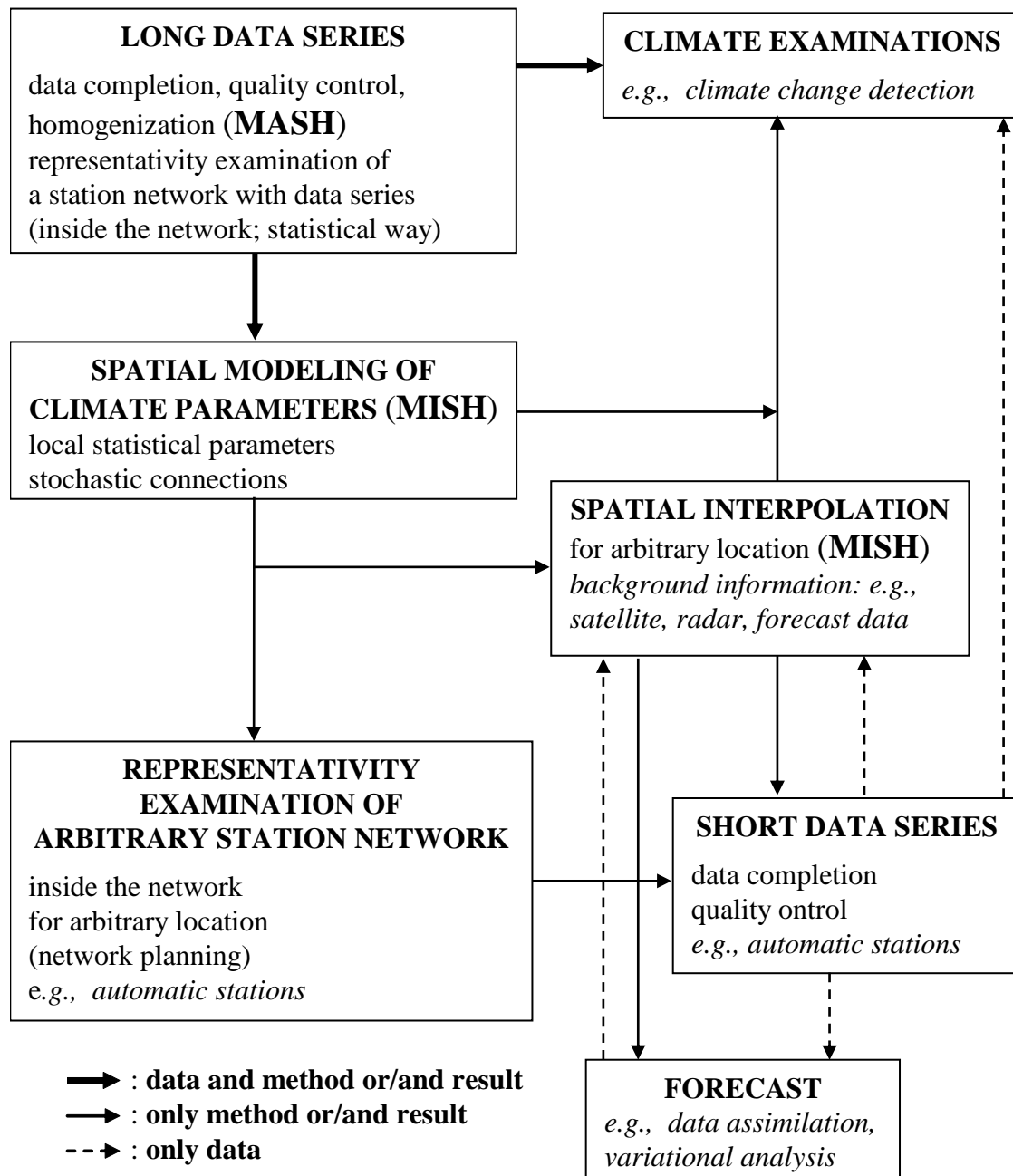


Fig. 1. Block diagram for the possible connection between various basic meteorological topics and systems



In our conception the meteorological questions and topics cannot be treated separately. Therefore we present a block diagram (*Fig. 1*) to illustrate the possible connection between various important meteorological topics. The software MASH (Multiple Analysis of Series for Homogenization; *Szentimrey*, 1999, 2014) and MISH (Meteorological Interpolation based on Surface Homogenized Data Basis; *Szentimrey* and *Bihari*, 2014) were developed by us. These software were applied also in CARPATCLIM project (<http://www.carpatclim-eu.org>).

## 2. MATHEMATICAL FORMULATION OF CLIMATE DATA HOMOGENIZATION

Unfortunately the exact theoretical, mathematical formulation of the problem of homogenization is neglected at the meteorological studies in general. Therefore we try to formulate this problem in accordance with the mathematical conventions. First of all it is emphasized that the homogenization is a distribution problem and not a regression one.

### 2.1 General mathematical formulation and some theorems

#### *Notation*

Let us assume we have daily or monthly climate data series:

$Y_1(t)$  ( $t = 1, 2, \dots, n$ ): candidate time series of the new observing system.

$Y_2(t)$  ( $t = 1, 2, \dots, n$ ): candidate time series of the old observing system.

$1 \leq T < n$ : change-point, series  $Y_2(t)$  ( $t = 1, 2, \dots, T$ ) can be used before  
and series  $Y_1(t)$  ( $t = T + 1, \dots, n$ ) can be used after the change-point.

The appropriate theoretical cumulative distribution (CDF) functions are:

$$F_{1,t}(y) = P(Y_1(t) < y) \quad , \quad F_{2,t}(y) = P(Y_2(t) < y) \quad y \in (-\infty, \infty) \quad , \quad t = 1, 2, \dots, n$$

It is very important to remark that as a consequence of some natural changes - e.g. annual cycle, climate change - the series of distribution functions  $F_{1,t}(y)$ ,  $F_{2,t}(y)$  ( $t = 1, 2, \dots, n$ ) may change in time! In the statistical climatology the climate change is equivalent with the changing probability of the meteorological events. The inhomogeneity of data series can be defined on the basis of the distribution functions.

#### *Definition 1*

The merged series  $Y_2(t)$  ( $t = 1, 2, \dots, T$ ),  $Y_1(t)$  ( $t = T + 1, \dots, n$ ) is inhomogeneous, if the identity of the distribution functions  $F_{2,t}(y) \equiv F_{1,t}(y)$  ( $t = 1, 2, \dots, T$ ) is not true.

### **Definition 2**

The aim of the homogenization is the adjustment or correction of values  $Y_2(t)$  ( $t = 1, 2, \dots, T$ ) in order to have the corrected values  $Y_{1,2h}(t)$  ( $t = 1, 2, \dots, T$ ) with the same distribution as the elements of series  $Y_1(t)$  ( $t = 1, 2, \dots, T$ ) have, i.e.:

$$\mathbf{P}(Y_{1,2h}(t) < y) = \mathbf{P}(Y_1(t) < y) = F_{1,t}(y) \quad y \in (-\infty, \infty), t = 1, 2, \dots, T \quad (1)$$

The formula (1) means the equality in distribution:  $Y_{1,2h}(t) \stackrel{d}{=} Y_1(t)$  ( $t = 1, 2, \dots, T$ )

### **Remark 1**

Within the same climate area, if the variables  $Y_1(t), Y_2(t)$  ( $t = 1, 2, \dots, T$ ) have identical distribution, i.e.  $Y_2(t) \stackrel{d}{=} Y_1(t)$  ( $t = 1, 2, \dots, T$ ), then the merged series  $Y_2(t)$  ( $t = 1, 2, \dots, T$ ),  $Y_1(t)$  ( $t = T + 1, \dots, n$ ) is homogeneous.

Some mathematical existence and unicity theorems can be proved in connection with the homogenization.

### **Theorem 1 (existence)**

Let us assume about the random variables  $Y_1, Y_2$  and their distribution functions  $F_1(y), F_2(y)$ , that  $\mathbf{P}(Y_j \in (a_j, b_j)) = 1$  and  $F_j(y)$  is a strictly increasing continuous function on the interval  $(a_j, b_j)$  ( $j = 1, 2$ ). Then applying the transfer function  $Y_{1,2h} = F_1^{-1}(F_2(Y_2))$  we obtain that the variable  $Y_{1,2h}$  has the same distribution like  $Y_1$  i.e.  $\mathbf{P}(Y_{1,2h} < y) = \mathbf{P}(Y_1 < y) = F_1(y)$ .

*Proof.*

The distribution function of  $Y_{1,2h}$  is as follows.

$$\mathbf{P}(Y_{1,2h} < y) = \mathbf{P}(Y_1 < y) = 0 \quad \text{if } y \leq a_1 \quad \text{and} \quad \mathbf{P}(Y_{1,2h} < y) = \mathbf{P}(Y_1 < y) = 1 \quad \text{if } y \geq b_1.$$

Furthermore if  $a_1 < y < b_1$  then,

$$\begin{aligned} \mathbf{P}(Y_{1,2h} < y) &= \mathbf{P}(F_1^{-1}(F_2(Y_2)) < y) = \mathbf{P}(F_2(Y_2) < F_1(y)) = \\ &= \mathbf{P}(Y_2 < F_2^{-1}(F_1(y))) = F_2(F_2^{-1}(F_1(y))) = F_1(y) \end{aligned}$$

### **Theorem 2 (unicity)**

Let us assume again about the random variables  $Y_1, Y_2$  and their distribution functions  $F_1(y), F_2(y)$ , that  $\mathbf{P}(Y_j \in (a_j, b_j)) = 1$  and  $F_j(y)$  is a strictly increasing continuous function on the interval  $(a_j, b_j)$  ( $j = 1, 2$ ). Let  $h(s)$  be also a strictly increasing continuous function on the interval  $(a_2, b_2)$ .

Then the distribution function  $P(h(Y_2) < y) = F_1(y)$  if and only if  $h(Y_2) = F_1^{-1}(F_2(Y_2))$ .

*Proof.*

According to the *Theorem 1*,  $P(F_1^{-1}(F_2(Y_2)) < y) = F_1(y)$ .

In the other direction - supposing  $P(h(Y_2) < y) = F_1(y)$  - we use the notations of the next *Lemma 1* that is,

$$g(y) = F_1^{-1}(F_2(h^{-1}(y))) \quad \text{and} \quad Y = h(Y_2).$$

Then  $g(y)$  is a strictly increasing continuous function on the interval  $(a_1, b_1)$  and the distribution functions of variables  $Y$ ,  $g(Y)$  are identical,  $P(g(Y) < y) = P(Y < y) = F_1(y)$ .

Consequently according to the *Lemma 1*,

$$g(Y) = Y \quad \text{i.e.} \quad h(Y_2) = F_1^{-1}(F_2(h^{-1}(h(Y_2)))) = F_1^{-1}(F_2(Y_2)).$$

### ***Lemma 1***

Let us assume about the random variable  $Y$  and its distribution function  $F(y)$ , that  $P(Y \in (a, b)) = 1$  and  $F(y)$  is a strictly increasing continuous function on the interval  $(a, b)$ .

Let  $g(y)$  be also a strictly increasing continuous function on the interval  $(a, b)$ .

Then the distribution function  $P(g(Y) < y) = F(y)$  if and only if  $g(Y) = Y$ .

*Proof.*

Let us suppose that  $P(g(Y) < y) = F(y)$  and  $a < y < b$ . Then,

$$F(y) = P(g(Y) < y) = P(Y < g^{-1}(y)) = F(g^{-1}(y))$$

Consequently if  $a < y < b$  then,

$$F(y) = F(g^{-1}(y)) \quad \text{that is} \quad y = g^{-1}(y) \quad \text{then} \quad g(y) = y.$$

## **2.2 Arising mathematical questions to be solved**

Let us suppose the merged series is given that is,

$$Y_2(t) \quad (t = 1, 2, \dots, T), \quad Y_1(t) \quad (t = T + 1, \dots, n)$$

In addition we suppose that the assumptions of the former theorems are fulfilled, consequently the theoretical correction or transfer formulas for the series elements are,

$$Y_{1,2h}(t) = F_{1,t}^{-1}(F_{2,t}(Y_2(t))) \quad (t = 1, 2, \dots, T) \quad (2)$$

However these transfer formulas are theoretical ones and if we want to apply them in the practice then a number of mathematical statistical estimation problems are arising. The most important problems are as follows.

- Estimation, detection of the change point(s)  $T$ .
- Estimation of the theoretical distribution functions  $F_{1,t}(y)$ ,  $F_{2,t}(y)$  ( $t = 1, 2, \dots, T$ ):
  - i,  $F_{1,t}(y)$ ,  $F_{2,t}(y)$  may change in time because of the climate change and the annual cycle, consequently the methodology of the use of the empirical distribution functions is very doubtful.
  - ii, There is no sample for  $F_{1,t}(y)$  ( $t = 1, 2, \dots, T$ ) and  $F_{2,t}(y)$  ( $t = T + 1, \dots, n$ ) usually.

These mathematical problems are insolvable generally! Therefore only relative methods can be used with some model assumptions. In addition some simplifications are also necessary.

### 2.3 Mathematical formulation for normal distribution

The homogenization problem is very complicated in general case however in case of normal distribution a much simpler mathematical formula can be obtained. We emphasize that the normal distribution is a special case but it is basic one in the mathematical statistics as well as in the meteorology. For example the normal distribution model can be accepted for the temperature variables in general.

#### *Theorem 3*

Let us assume the data series have normal distribution that is,

$$Y_1(t) \in N(E_1(t), D_1(t)), \quad Y_2(t) \in N(E_2(t), D_2(t)) \quad (t = 1, 2, \dots, n),$$

where  $E(Y_1(t)) = E_1(t)$ ,  $E(Y_2(t)) = E_2(t)$  are the means or expected values and

$D(Y_1(t)) = D_1(t)$ ,  $D(Y_2(t)) = D_2(t)$  are the standard deviations.

Then the transfer formula of homogenization:

$$Y_{1,2h}(t) = F_{1,t}^{-1}(F_{2,t}(Y_2(t))) = E_1(t) + \frac{D_1(t)}{D_2(t)}(Y_2(t) - E_2(t)) \quad (t = 1, 2, \dots, T)$$

#### *Proof.*

If the data series have normal distribution then the distribution functions can be written in the following form,

$$F_{1,t}(y) = \Phi\left(\frac{y - E_1(t)}{D_1(t)}\right), \quad F_{2,t}(y) = \Phi\left(\frac{y - E_2(t)}{D_2(t)}\right)$$

where  $\Phi(s)$  is the standard normal distribution function.

Consequently,

$$\begin{aligned}
Y_{1,2h}(t) &= F_{1,t}^{-1}(F_{2,t}(Y_2(t))) = F_{1,t}^{-1}\left(\Phi\left(\frac{Y_2(t) - E_2(t)}{D_2(t)}\right)\right) = \\
&= D_1(t) \cdot \Phi^{-1}\left(\Phi\left(\frac{Y_2(t) - E_2(t)}{D_2(t)}\right)\right) + E_1(t) = E_1(t) + \frac{D_1(t)}{D_2(t)}(Y_2(t) - E_2(t)) \quad (t = 1, 2, \dots, T)
\end{aligned}$$

since the inverse of function  $\Phi\left(\frac{y-m}{d}\right)$  is  $d \cdot \Phi^{-1}(y) + m$ .

## 2.4 Mathematical questions in case of normal distribution to be solved

In case of normal distribution according to the *Theorem 3* we have a much simpler transfer formula for correction than the general form (2), that is,

$$Y_{1,2h}(t) = E_1(t) + \frac{D_1(t)}{D_2(t)}(Y_2(t) - E_2(t)) \quad (t = 1, 2, \dots, T) \quad (3)$$

This formula is a simple linear one that means if the data series have normal distribution it is sufficient to homogenize the means and standard deviations only that is equivalent with the homogenization of the first two moments. We emphasize that the normal distribution is a basic model in the mathematical statistics as well as in the meteorology and there is no “tail distribution” problem at this important distribution according to the *Theorem 3*! At the normal distribution if the means and standard deviations are homogenous then the higher order moments are also homogeneous and there is not any inhomogeneity in the tails of the distributions. It is in contrast with the popular assumption based on parallel measurements as it is very likely the inhomogeneity in the tails of the distributions at the daily data series. As regards the parallel measurements a mathematical examination for them will be presented at Section 2.5.

Returning to the formula (3) although it is much simpler than (2), there are still a number of mathematical statistical estimation problems to be solved as follows.

- Estimation, detection of the change point(s)  $T$ .
- Estimation of the statistical parameters  $E_1(t), D_1(t), E_2(t), D_2(t)$  ( $t = 1, 2, \dots, T$ ):
  - i,  $E_1(t), D_1(t), E_2(t), D_2(t)$  may change in time because of the climate change and the annual cycle.
  - ii, There is no sample for  $E_1(t), D_1(t)$  ( $t = 1, \dots, T$ ) and  $E_2(t), D_2(t)$  ( $t = T + 1, \dots, n$ ) usually.

However these mathematical problems are still very complicated! Therefore only relative methods can be used with some model assumptions. In addition some simplifications are also necessary.

The most often applied transfer formula in the practice can be obtained from the formula (3) with the following simplifications,

$$D_2(t) = D_1(t) , \quad E_2(t) - E_1(t) = E \quad (t = 1, 2, \dots, T)$$

Then the transfer formula is,

$$Y_{1,2h}(t) = Y_2(t) - E \quad (t = 1, 2, \dots, T)$$

This is the homogenization in mean applied in the practice mostly (Section 3).

## 2.5 Mathematical examinations of parallel measurements

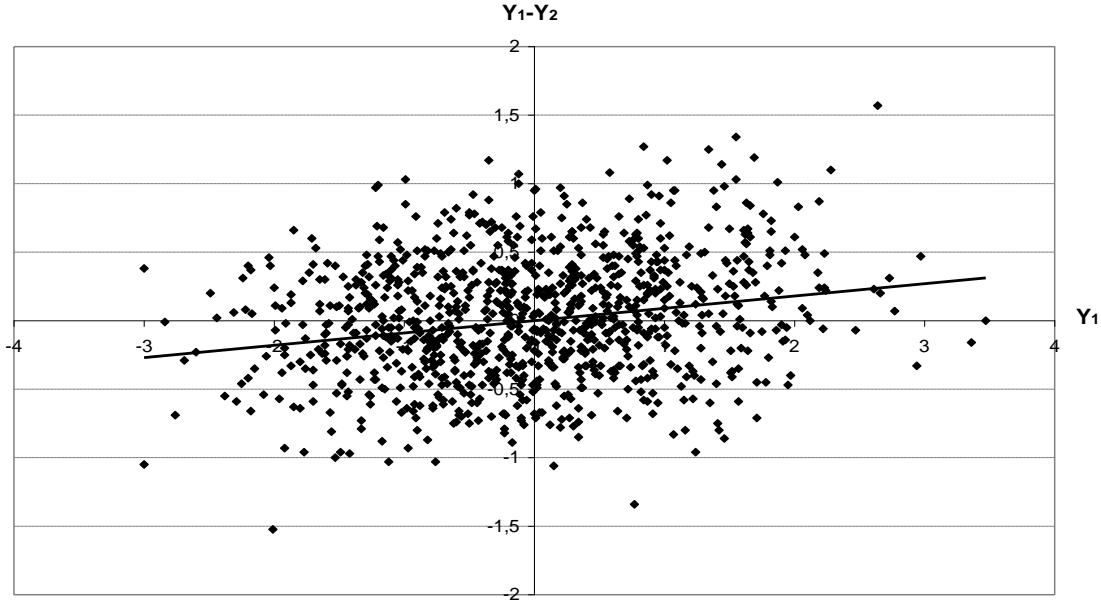
On the one hand the daily data series are very important for studying extremes. On the other hand there is a popular assumption based on parallel measurements and some physical considerations as it is very likely the inhomogeneity in the tails of the distributions at the daily data series. What is the reason of this assumption?

Essentially the reason is an observed phenomenon at the extremes, namely the differences of parallel measurements are larger in case of extremes. In our opinion, this observed phenomenon has a simple and logical reason, and it is superfluous to look for some complicated physical explanation for the inhomogeneity. The simple reason is that the extremes may be expected at different moments in case of parallel measurements, or in other words, there may be systematic biases in rank order! It is a natural phenomenon, and for illustration a trivial example is presented according to the probability theory.

*Example 2.5* Let  $Y_1(t) \in N(0,1)$ ,  $Y_2(t) \in N(0,1)$  ( $t = 1, 2, \dots, n$ ) be series of independent and standard normally distributed variables with expected values  $E(Y_1(t)) = E(Y_2(t)) = 0$ , with standard deviations  $D(Y_1(t)) = D(Y_2(t)) = 1$ , and with correlation between the series  $\text{corr}(Y_1(t), Y_2(t)) = \rho$  ( $t = 1, 2, \dots, n$ ).

Then the mean difference  $E(Y_1(t) - Y_2(t)) = 0$  of course, however, the difference  $Y_1(t) - Y_2(t)$  is not independent from the elements  $Y_1(t)$ ,  $Y_2(t)$  if  $\rho \neq 1$ , and, e.g., the conditional expectation of difference  $Y_1(t) - Y_2(t)$  given  $Y_1(t)$ , or equivalently the regression of difference  $Y_1(t) - Y_2(t)$  on  $Y_1(t)$  is  $E(Y_1(t) - Y_2(t) | Y_1(t)) = (1 - \rho) \cdot Y_1(t)$ .

Consequently, the difference  $Y_1(t) - Y_2(t)$  is an expectedly monotonous increasing function of  $Y_1(t)$  if  $\rho \neq 1$ . This is the theory, but it can be demonstrated in practice too. We generated such standard normal series by the Monte Carlo method with parameters  $\rho = 0.9$ ,  $n = 1000$ . In this case,  $E(Y_1(t) - Y_2(t) | Y_1(t)) = 0.1 \cdot Y_1(t)$  and the difference series  $Y_1(t) - Y_2(t)$  as a function of series  $Y_1(t)$  is plotted in Fig. 2.



**Fig. 2. Difference series  $Y_1(t) - Y_2(t)$  as a function of series  $Y_1(t)$**

It is evident that the conditional expectation of difference  $Y_1(t) - Y_2(t)$  is monotonous increasing function of  $Y_1(t)$ , consequently the difference may be larger mainly in the case of extreme values. There is no inhomogeneity it is a general phenomenon that can be observed also at the meteorological measurements.

## 2.6 Problem of inhomogeneity of the standard deviation

There is also a popular assumption applied in the practice that the correction in mean is sufficient for monthly and annual series, and that the correction of higher order moments is necessary only in the case of daily data series. In general, it is tacitly assumed that the averaging is capable to filter out the inhomogeneities in the higher order moments. However, this assumption is false, for example, if there is a common inhomogeneity in the standard deviation of daily data, we may have the same inhomogeneity in monthly data.

### ***Lemma 2***

Let us assume  $Y(t)$  ( $t = 1, \dots, 30$ ) are daily data and the monthly average is  $\bar{Y} = \frac{1}{30} \sum_{t=1}^{30} Y(t)$ .

Let us introduce some inhomogeneity of the mean and the standard deviation for the daily data by a linear function:

$$Y_{ih}(t) = \alpha \cdot (Y(t) - E(Y(t))) + E(Y(t)) + \beta \quad (t = 1, \dots, 30)$$

Then the expected values and the standard deviations are:

$$E(Y_{ih}(t)) = E(Y(t)) + \beta, \quad D(Y_{ih}(t)) = \alpha \cdot D(Y(t)) \quad (t = 1, \dots, 30)$$

Let us see the new monthly average:  $\bar{Y}_{ih} = \frac{1}{30} \sum_{t=1}^{30} Y_{ih}(t)$ .

Then the expected value and the standard deviation also changed with the same measure like the daily values:

$$E(\bar{Y}_{ih}) = E\left(\frac{1}{30} \sum_{t=1}^{30} Y_{ih}(t)\right) = E\left(\frac{1}{30} \sum_{t=1}^{30} E(Y(t)) + \beta\right) = E(\bar{Y}) + \beta$$

$$D(\bar{Y}_{ih}) = D\left(\frac{1}{30} \sum_{t=1}^{30} Y_{ih}(t)\right) = D\left(\frac{1}{30} \sum_{t=1}^{30} \alpha \cdot Y(t)\right) = \alpha \cdot D\left(\frac{1}{30} \sum_{t=1}^{30} Y(t)\right) = \alpha \cdot D(\bar{Y}).$$

## 2.7 Mathematical formulation of conditional homogenization

In our earlier papers we made some criticism about the so called variable correction methods especially about their underlying principles. Their common assumption is that in case of daily data series the corrections for inhomogeneity have to vary according to the meteorological situation of each day in order to represent the extremes. This assumption is based also on the parallel measurements and some physical considerations as it was analyzed at Section 2.6. We do not agree with this argumentation and we think that the transfer function (2) is appropriate for correction of inhomogeneity however we began to develop a general mathematical form of the conditional homogenization that can be applied if we have some supplementary information as condition. The following parts are in draft form since on the one hand this mathematical development is still at early stage and on the other hand the detailed description should need advanced mathematical tools.

### 2.7.1 Conditional homogenization based on given events

Let  $B = \{B_j : j = 1, 2, \dots, M\}$  be a complete system of events:

$$B_i \cap B_j = \emptyset, \quad \sum_{j=1}^M P(B_j) = 1 \quad (\text{e.g. macrosynoptic weather situations})$$

Conditional homogenization of  $Y_2(t)$  on given events  $B$ ,

$$Y_{1,2h}(t, B) = F_{1,t,B_j}^{-1}\left(F_{2,t,B_j}(Y_2(t))\right) \Leftrightarrow B_j \text{ occurs at } t \quad (t = 1, 2, \dots, T)$$

where  $F_{1,t,B_j}(y)$ ,  $F_{2,t,B_j}(y)$  are the conditional distribution functions

of  $Y_1(t)$ ,  $Y_2(t)$ , given  $B_j$ , that is

$$F_{1,t,B_j}(y) = P(Y_1(t) < y | B_j), \quad F_{2,t,B_j}(y) = P(Y_2(t) < y | B_j) \quad y \in (-\infty, \infty), \quad t = 1, 2, \dots, T$$

Then as a consequence of Bayes and total probability theorems:

$$P(Y_{1,2h}(t, B) < y) = F_{1,t}(y) \quad y \in (-\infty, \infty), \quad t = 1, 2, \dots, T$$



### 2.7.2 Conditional homogenization in more general form

Let  $\mathbf{Z}(t)$  ( $t = 1, 2, \dots, n$ ) be a homogeneous climate (vector) time series

Conditional homogenization of  $Y_2(t)$  on given  $\mathbf{Z}(t)$ ,

$$Y_{1,2h}(t, \mathbf{Z}(t)) = F_{1,t,\mathbf{z}}^{-1}(F_{2,t,\mathbf{z}}(Y_2(t))) \Leftrightarrow \mathbf{Z}(t) = \mathbf{z} \quad (t = 1, 2, \dots, T)$$

where  $F_{1,t,\mathbf{z}}(y)$ ,  $F_{2,t,\mathbf{z}}(y)$  are the conditional distribution functions

of  $Y_1(t)$ ,  $Y_2(t)$ , given  $\mathbf{Z}(t) = \mathbf{z}$ , that is

$$F_{1,t,\mathbf{z}}(y) = P(Y_1(t) < y \mid \mathbf{Z}(t) = \mathbf{z}), \quad F_{2,t,\mathbf{z}}(y) = P(Y_2(t) < y \mid \mathbf{Z}(t) = \mathbf{z})$$

$$y \in (-\infty, \infty), \quad t = 1, 2, \dots, T$$

Then as a consequence of Bayes and total probability theorems:

$$P(Y_{1,2h}(t, \mathbf{Z}(t)) < y) = F_{1,t}(y) \quad y \in (-\infty, \infty), \quad t = 1, 2, \dots, T$$

## 3. RELATION OF DAILY AND MONTHLY HOMOGENIZATION

The theme of homogenization can be divided into two subgroups, such as monthly and daily data series homogenization. These subjects are in strong connection with each other of course, for example the monthly results can be used for the homogenization of daily data.

### 3.1 The general structure of daily data homogenization

If we have daily data series the general way of homogenization is,

- calculation of monthly series,
- homogenization of monthly series taking advantage of the larger signal to noise ratio,
- homogenization of daily series using the detected monthly inhomogeneities.

So we have the question how can we use the valuable information of detected monthly inhomogeneities for the daily data homogenization?

### 3.2 A popular procedure e.g. the variable correction methods

The typical steps of the procedure are as follows.

1. Homogenization of monthly series:

Break points detection, correction in the first moment (mean).

Assumption: homogeneity in higher order moments (e.g. standard deviation).

## 2. Homogenization of daily series:

There is a trial to homogenize also in higher order moments.

The used monthly information are only the detected break points.

However the following questions are arising at this procedure:

- Is it adequate model that we have inhomogeneity in higher moments only at daily series but not at monthly ones? Can this model be accepted according to the probability theory? According to Section 2.6 the correct answer is that this model cannot be accepted.
- Why are not used the monthly correction factors for daily homogenization? It seems to lose some valuable information obtained during the monthly homogenization.

### 3.3 An alternative procedure

We suggest an alternative procedure to homogenize both the daily and the monthly series.

The steps of the procedure in case of quasi normal distribution (additive model, e.g. temperature) are as follows.

#### 1. Homogenization of monthly series:

Break points detection, correction of the first two moments that is equivalent with the homogenization of mean and standard deviation. The correction is based on the transfer formula (3).

Assumption: homogeneity in higher order moments. This assumption is always right in case of normal distribution according to *Theorem 3*.

#### 2. Homogenization of daily series:

Homogenization of mean and standard deviation on the basis of the monthly results. The used monthly information are the break points and the monthly corrections of the mean and standard deviation. The correction is based on the transfer formula (3) considering Lemma 2. If the daily data are normally distributed then there is no inhomogeneity in the higher order moments according to *Theorem 3*.

In case of quasi lognormal distribution (multiplicative model, e.g. precipitation) also the above procedure can be applied for the data obtained by certain transformation based on logarithmization.

## 4. OVERVIEW ON HOMOGENIZATION IN MEAN OF MONTHLY SERIES

This section considers some various theoretical aspects of monthly series homogenization. In the practice the monthly series are homogenized in the mean mostly. The aim of these homogenization procedures is to detect the inhomogeneities of mean and to correct the series. In connection with the such type of homogenization methods we have to give solutions for the following mathematical problems: relative models, statistical spatiotemporal modelling of the series, methodology for comparison of series, break point (change point) and outlier detection, methodology for correction of series, quality control procedures, missing data completion,

usage of metadata, relation of daily and monthly homogenization, manual versus automatic methods, evaluation of methods (theoretical, benchmark).

In practice there are absolute and relative methods applied for homogenization. However the main problem of the application of absolute methods is that the separation of climate change signal and the inhomogeneity is essentially impossible. Relative methods can be applied if there are more station series given, which can be compared mutually. In this case the statistical spatiotemporal modelling of the series is a fundamental question. The adequate comparison, break point detection and correction procedures are depending on the chosen statistical model.

#### 4.1 General structure of additive spatiotemporal models

If the data series are normally distributed (e.g. temperature) then the additive model can be used. In case of relative methods a general form of additive model for more monthly series belonging to the same month in a small climate region can be written as follows,

$$X_j(t) = \mu(t) + E_j + IH_j(t) + \varepsilon_j(t) \quad (j = 1, 2, \dots, N; t = 1, 2, \dots, n), \quad (4)$$

where  $\mu(t)$  is the common and unknown climate change signal,  $E_j$  are the spatial expected values,  $IH_j(t)$  are the inhomogeneity signals and  $\varepsilon_j(t)$  are normal white noise series. The type of inhomogeneity  $IH(t)$  is in general a 'step-like function' with unknown break points  $T$  and shifts  $IH(T) - IH(T + 1) \neq 0$ , and  $IH(n) = 0$  is assumed in general.

The normal distributed vector variables  $\boldsymbol{\varepsilon}(t) = [\varepsilon_1(t), \dots, \varepsilon_N(t)]^T \in N(\mathbf{0}, \mathbf{C})$  ( $t = 1, \dots, n$ ) are totally independent in time. The spatial covariance matrix  $\mathbf{C}$  describes the spatial structure of the series.

If the data series are quasi lognormal distributed (e.g. precipitation) then the multiplicative model can be used that can be transformed into the additive one by certain logarithmic procedure.

#### 4.2 Methodology for comparison of series

The problem of comparison of series is related to the following questions: reference series creation, difference series constitution, multiple comparisons of series etc. This topic is very important for detection as well as for correction, because the efficient series comparison can increase both the significance and the power. The development of efficient comparison methods can be based on the examination of the spatial covariance structure of data series. The examined series  $X_j(t)$  ( $j = 1, \dots, N$ ) have to be taken as candidate and reference series alike, furthermore the reference series are not assumed to be homogeneous at the correct examinations!

The main problem arises from the fact that the shape of climate change signal is unknown. Therefore so-called difference series are examined in order to filter out the climate change signal  $\mu(t)$ . The simple difference series between pairs are  $Z(t) = X_j(t) - X_i(t)$ . However the difference series constitution can be formulated in more general way as well. Assuming that

$X_j(t)$  is the candidate series and the other ones are the reference series, then the difference series belonging to the candidate series can be constituted as,

$$Z_j(t) = X_j(t) - \sum_{i \neq j} \lambda_{ji} X_i(t) = IH_j(t) - \sum_{i \neq j} \lambda_{ji} IH_i(t) + \varepsilon_{Z_j}(t) \quad (5)$$

with condition of  $\sum_{i \neq j} \lambda_{ji} = 1$  for the weighting factors. As a result of the last condition, the unknown climate change signal  $\mu(t)$  has been filtered out. Consequently the inhomogeneities can be detected by the examination of the above difference series. In addition if we want to increase the signal to noise ratio in order to increase the power of detection then we have to minimize the variance of noise term  $\varepsilon_{Z_j}(t)$ .

The covariance matrix  $\mathbf{C}$  uniquely determines the optimum weighting factors that minimize the variance, and the optimal difference series created in this manner can be applied efficiently for the detection and correction procedures (MASH, *Szentimrey*, 1999, 2014). We mention that in case of using the generalized-least-squares estimation for the unknown climate change signal  $\mu(t)$ , also the optimal difference series is obtained with minimal variance. We have to examine more difference series in order to separate the appropriate detected inhomogeneities for the candidate series. More difference series created without common reference series and with minimal variances can be defined as optimal difference series system (MASH).

### 4.3 Methodology for break point (changepoint) detection

One of the basic tasks of the homogenization is the examination of more difference series in order to detect the break points and to attribute them for the candidate series. The key question of the homogenization software is to develop automatic procedures for this attribution problem!!!

The scheme of the break point detection is as follows. Let  $Z(t)$  be a difference series according to the formula (5), that is

$$Z(t) = IH_Z(t) + \varepsilon_Z(t) \quad (t = 1, \dots, n), \quad (6)$$

where  $IH_Z(t)$  is a mixed inhomogeneity of difference series  $Z(t)$  with  $K$  break points  $T_1 < T_2 < \dots < T_K$ . In general the number  $K$  and the position of the multiple break points  $T_1 < T_2 < \dots < T_K$  are unknown, furthermore the noise variables  $\varepsilon_Z(t) \in N(E_Z, \sigma_Z^2)$  ( $t = 1, \dots, n$ ) are totally independent in time. The basic types of the detection procedures are the stepwise and the multiple break points detection. Let us have the following notation of the estimates:  $\hat{K}; \hat{T}_1 < \hat{T}_2 < \dots < \hat{T}_{\hat{K}}$ .

The more sophisticated multiple break points detection procedures were developed for joint estimation of the break points. There may be different principles of these methods that are classical ways in mathematical statistics.

#### 4.3.1 Break point detection based on Bayesian Approach

The methods based on Bayesian model selection are the penalized likelihood methods. These methods are different in the penalty terms or criterions e.g. Akaike criterion, Schwarz criterion, Caussinus-Lyazrhi criterion.

The PRODIGE procedure (Caussinus and Mestre, 2004) based on the Caussinus-Lyazrhi criterion is an example for the penalized likelihood methods.

#### 4.3.2 Break point detection based on Test of Hypothesis

Another possibility is to use hypothesis test methods for the detection of break points. At the MASH method (Szentimrey, 1999, 2014) a hypothesis test procedure has been developed, as we want to avoid the type one error that is the damage of data series. The essence of this multiple break points detection procedure based on test of hypothesis on a given significance level is as follows.

If the detected break points of  $Z(t)$  are  $\hat{K}; \hat{T}_1 < \hat{T}_2 < \dots < \hat{T}_{\hat{K}}$ , then on the given significance level  $p$  (e.g.:  $p=0.1$ ):

i,  $Z(t)$  is not homogeneous above the intervals  $(\hat{T}_{k-1}, \hat{T}_{k+1}]$  because,

$$P\left(\exists(\hat{T}_{k-1}, \hat{T}_{k+1}] \text{ above that : } Z(t) \text{ homogeneous}\right) = p$$

Consequently the detected break points  $\hat{T}_k$  are not superfluous.

This means there is no serious type one error.

ii,  $Z(t)$  can be accepted to be homogeneous above the intervals  $(\hat{T}_{k-1}, \hat{T}_k]$ .

This means there is no serious type two error.

#### Remark

Confidence intervals are also given for the break points beside the point estimation at the method MASH (Szentimrey, 1999, 2014).

#### 4.4 Methodology for correction of series

Beside the detection another basic task of the homogenization is the correction of series. Calculating of correction factors can be based on the examination of difference series for estimation of shifts  $IH(\hat{T}_k) - IH(\hat{T}_k + 1)$  ( $k = 1, \dots, \hat{K}$ ) at the detected break points.

Almost all the methods use point estimation for the correction factors at the detected break points. For example the PRODIGE method (*Caussinus and Mestre, 2004*) uses the standard least squares technique to estimate the correction factors. Probably the generalized least squares estimation technique based on spatial covariance structure would be more efficient.

The MASH procedure (*Szentimrey, 1999, 2014*) is an exception because the correction factors are estimated on the basis of confidence intervals. The confidence intervals given for the break points and shifts make possible also the automatic usage of metadata at MASH!

## **4.5 Automation of methods and software**

One of the fundamental problems of homogenization procedures is the relation of the manual versus interactive or automatic methods. In the practice the simple manual methods (e.g. Craddock method) are very popular however these ones are unusable for the real climate observation networks. We have to understand the fact that a lot of stations must be examined together in general! The reality for the number of stations per network is more than 100 instead of 10-15 used at COST HOME benchmark dataset.

Therefore the key questions of the homogenization methods and software are,

- firstly, the quality of homogenized data,
- secondly, the quantity of stations.

If we want to fulfill both respects it is necessary to develop automatic procedures.

Further necessary conditions required for automation of methods and software are,

- ability for automatic attribution of break points for the candidate series,
- automatic usage of metadata.

To solve the above problems without advanced mathematics is impossible!!!

## **4.6 Possibilities for evaluation of the methods**

### **4.6.1 Theoretical evaluation**

If want to obtain a real image of the methods, then the theoretical evaluation of their mathematical basis is indispensable.

### **4.6.2 Benchmark**

The COST Action ES0601 (HOME) has executed a blind intercomparison and validation study for monthly homogenization methods. The methods were tested on a realistic benchmark dataset. The benchmark contained simulated data with inserted inhomogeneities (*Venema et al., 2012*). Testing the methods on a generated benchmark dataset seems to be an objective validation procedure however we have to know also the limits of such type of examinations.

The interpretation of benchmark results is not a trivial problem, since these are depending on different factors, such as:

- tested methods (quality, manual or automatic),

- testing benchmark dataset (quality, adequacy),
- testers (skilled or unskilled),
- methodology of evaluation (validation statistics).

The creation of adequate benchmark dataset and the development of appropriate validation statistics are critical points and they need also strong theoretical mathematical background.

We remark that the question of the comparison of manual methods to automatic ones seems similar to the comparison of handmade and factory products. Or how can we compare the results of a manual time consuming method with a skilled tester versus the results of an interactive method with an unskilled tester. The method or the user is tested if we evaluate the test results?

## 5. SOFTWARE MASH

### The most important properties of MASHv3.03

(Multiple Analysis of Series for Homogenization; *Szentimrey* 1999, 2008, 2014)

#### Homogenization of monthly series:

- Relative homogeneity test procedure.
- Step by step iteration procedure: the role of series (candidate, reference)
- changes step by step in the course of the procedure.
- Additive (e.g. temperature) or multiplicative (e.g. precipitation) model
- can be used depending on the climate elements.
- Including quality control and missing data completion.
- Providing the homogeneity of the seasonal and annual series as well.
- Metadata (probable dates of break points) can be used automatically.
- The homogenization results and the metadata can be verified.

#### Homogenization of daily series:

- Based on the detected monthly inhomogeneities.
- Including quality control and missing data completion for daily data.

#### Some MASH specialty

- use of spatial covariance for comparison of series
- automatic attribution of break points for the candidate series
- automatic use of metadata

Our MISH-MASH software can be downloaded from:

[http://www.met.hu/en/omsz/rendezvenyek/homogenizationand\\_interpolation/software/](http://www.met.hu/en/omsz/rendezvenyek/homogenizationand_interpolation/software/)

## References

- Caussinus, H, Mestre, O. 2004: Detection and correction of artificial shifts in climate series, *Appl. Statist.*, 53, Part 3, pp. 405-425.
- Szentimrey, T., 1999: Multiple Analysis of Series for Homogenization (MASH), *Proceedings of the Second Seminar for Homogenization of Surface Climatological Data*, Budapest, Hungary; WMO, WCDMP-No. 41, pp. 27-46.
- Szentimrey, T., 2008: Development of MASH homogenization procedure for daily data, *Proceedings of the Fifth Seminar for Homogenization and Quality Control in Climatological Databases*, Budapest, Hungary, 2006; WCDMP-No. 68, WMO-TD NO. 1434, 2008, pp. 116-125.
- Szentimrey T. and CARPATCLIM Homogenization-Interpolation Team, 2012: Final report on quality control and data homogenization measures applied per country, including QC protocols and measures to determine the achieved increase in data quality. *Deliverable D1.12 of CARPATCLIM*, homepage: <http://www.carpatclim-eu.org/pages/deliverables/>
- Szentimrey, T. 2013: Theoretical questions of daily data homogenization, *Időjárás* Vol. 117. No. 1, January-March 2013. pp. 113-122.
- Szentimrey, T., Bihari, Z., 2014: Manual of interpolation software MISHv1.03, Hungarian Meteorological Service, p. 60.
- Szentimrey, T., 2014: Manual of homogenization software MASHv3.03, Hungarian Meteorological Service, p. 69.
- Venema, V. K. C. et al., 2012: Benchmarking homogenization algorithms for monthly data, *Climate of the Past*, 8, 89-115



# IRELAND WITH HOMER

**John Coll<sup>1</sup>, Mary Curley<sup>2</sup>, Séamus Walsh<sup>2</sup>, John Sweeney<sup>1</sup>.**

<sup>1</sup> Irish Climate Analysis and Research Units, Department of Geography, NUI Maynooth,  
Maynooth, Co Kildare.

<sup>2</sup> Met Éireann, Glasnevin Hill, Dublin.

## **Abstract**

Our instrumental knowledge of climate change prior to the mid-19th century is heavily reliant on a few long meteorological series, mostly from Europe, and even here good instrumental series longer than 150 years are rare. However, climate change studies based only on raw long-term data are potentially flawed due to the many breaks introduced from non-climatic sources. Consequently quality controlled and homogenised climate data is desirable for basing climate related decision making on. HOMER was applied to eighty eight monthly precipitation station identified from the *Met Éireann database as stations with longer contiguous records* ranging from ~40 to 70 years between 1941 and 2010. Results on the reference networks and their associated geographical distances identified by the HOMER algorithm are compared with those derived via first difference correlations in a parallel statistical computing approach. Although only shown in the context of case study examples in the results here, results across the analysis for all 88 station records and their potential neighbour series indicate that both first difference correlations and HOMER geographical distance selections yield often corresponding neighbour series which are largely statistically and spatially coherent.

Keywords: Monthly precipitation, Ireland, homogenisation, HOMER

## **1. INTRODUCTION**

### **1.1. A policy context for homogenised data**

Quality control and homogenisation of climate data are becoming increasingly important as European Union (EU) Member States examine methods and put in place mechanisms for delivering integrated climate services. However, in reality climate services rarely go beyond the scope of meteorological variables and the gap between the supply of climate services and the needs of users has been identified by the World Meteorological Organization (WMO, 2011). Therefore the provision of fully quality controlled and homogenised series will become one part of an envisaged ‘end-to-end’ delivery chain for endpoint data provision to the impacts and policy communities.

In line with this, there is already a recognised need for Ireland to engage with the current EU Joint Programming Initiative (JPI) Climate and Horizon 2020 (H2020) initiatives in order that an entity or consortium to provide climate services can be provided. Importantly however, if useful climate information is to be delivered, this must be tailored to meet the needs of users. This is rarely the case however, since the level of interaction between providers and users of

climate services is inadequate, therefore users need access to expert advice and support to help them select and properly apply climate information. It is therefore incumbent upon the scientific community to address how best to meet the needs of users of climate information by attempting to bridge the gap between the climate modelling community and the end users of climate information. To facilitate this the Irish Environmental Protection Agency (EPA) have already formulated a plan to develop climate services which, among other goals, seeks to improve the translation of scientific information in clear terms to decision makers (EPA, 2014).

## **1.2. Scientific motivation for homogenised series**

The increasing interest in climate modelling has spurred the development and testing of a variety of homogenisation techniques aimed at identifying and sometimes correcting inhomogeneities in data series which do not reflect real variations in climate. A homogeneous climate time series is defined as one where variability is only caused by changes in weather or to the climate (Freitas et al. 2013). Most of the homogenisation techniques are addressed in classical or Bayesian statistical frameworks, supported by parametric or nonparametric models. Long instrumental records are rarely if ever homogeneous and most decade- to century-scale time series of atmospheric data have been adversely impacted by inhomogeneities caused by, for example; changes in instrumentation, station moves, changes in the local environment such as urbanisation, or the introduction of different observing practices like a new formula for calculating mean daily temperature or different observation times. Our instrumental knowledge of climate change prior to the mid-19th century is heavily reliant on a few long meteorological series, mostly from Europe, and even here good instrumental series longer than 150 years are rare. However, climate change studies based only on raw long-term data are potentially flawed due to the many breaks introduced from non-climatic sources, consequently accurate climate data is an essential prerequisite for basing climate related decision making on.

If inhomogeneities are not accounted for properly, the results of climate analyses using these data can be erroneous (Peterson et al. 1998). Also, there is a difference between an inhomogeneous series and a series that is non-stationary, as a series can be non-stationary and homogeneous at the same time, i.e. the change points are only caused by real climatic variations (Beaulieu et al. 2009). Since, the user of climatic data series is often unaware of the presence or absence of inhomogeneities in the series, the inhomogeneities can interfere with the real climate change signal and lead to poor climatic or impact model calibration or biased studies of climate trends and variability (Beaulieu et al. 2009). Consequently, the detection and correction of these inhomogeneities is important before undertaking any kind of climate analysis. Therefore in recent times, and building on earlier work (e.g. Alexandersson 1986; Jones et al. 1986) several techniques have been developed for the detection and adjustment of non-climatic inhomogeneities (Cao and Wan 2012; Toreti et al. 2012; Freitas et al. 2013; Mestre et al. 2013). Arising from this new work, more recent techniques have been developed to detect and correct multiple change points using reference series (Peterson et al. 1998; Menne and Williams, 2005; Toreti et al. 2012). More recently a comprehensive analysis to assess different homogenisation techniques of climate series was included in scientific programme of the COST Action HOME ES 0601: Advances in Homogenisation Methods of Climate series: an integrated approach (HOME). The HOME objective was to develop a general homogenisation method for homogenising climate and environmental datasets. This task commenced in 2007 and was finalised in 2011 with the release of two new software

packages, HOMER for the homogenisation of monthly data and HOM/SPLIDHOM for daily data homogenisation (HOME, 2013).

The primary aim of this short paper is to summarise results using the HOMER software to homogenise monthly mean precipitation (Precip) totals for Ireland for the 1941 – 2010 period for an initial set of 88 stations from the Met *Éireann* monthly station data repository. A secondary aim is to compare the reference station networks identified within HOMER with a complementary approach using correlation and other statistical measures for the series in combination with spatial scrutiny in GIS.

## **2. MATERIALS AND METHODS**

### **2.1. Study area**

The study area is the whole island of Ireland, that covers ~84 421 km<sup>2</sup> on the Atlantic margin of northwest Europe, between ~51° and 56° N. Elevations reach up to 1038 m above sea level (a.s.l.) (Corrán Tuathail, Co. Kerry). Much of the island is lowland, partly surrounded by mountains, with a characteristic temperate oceanic climate. Mean annual temperature (averaged over 1961 to 1990) is highest on the south-west coast (10.4°C) and lowest inland (8.8°C). On average, annual precipitation ranges from 750 to 1000 mm in the drier eastern half of the country and >3000 mm yr<sup>-1</sup> in parts of the western mountains (Rohan 1986).

### **2.2. Stations and data**

Rainfall has been measured in Ireland since the early nineteenth century with a peak of over 800 rainfall stations in the late 1950s, and currently rainfall is recorded at synoptic and climatological weather stations; in addition, there is a wide network of voluntary rainfall observers (Walsh, 2012). The selected stations for this initial phase of work are distributed across the country, but more spatial clustering of the available series is apparent in the east (Figure 1). Based on an earlier audit and initial set of quality control procedures, the contiguous intact monthly records for this initial group of 88 stations ranged from ~40 to 71 years. Stations elevations were within the range of 5 – 404 m above sea level (a.s.l.) with a mean elevation of ~84 m.

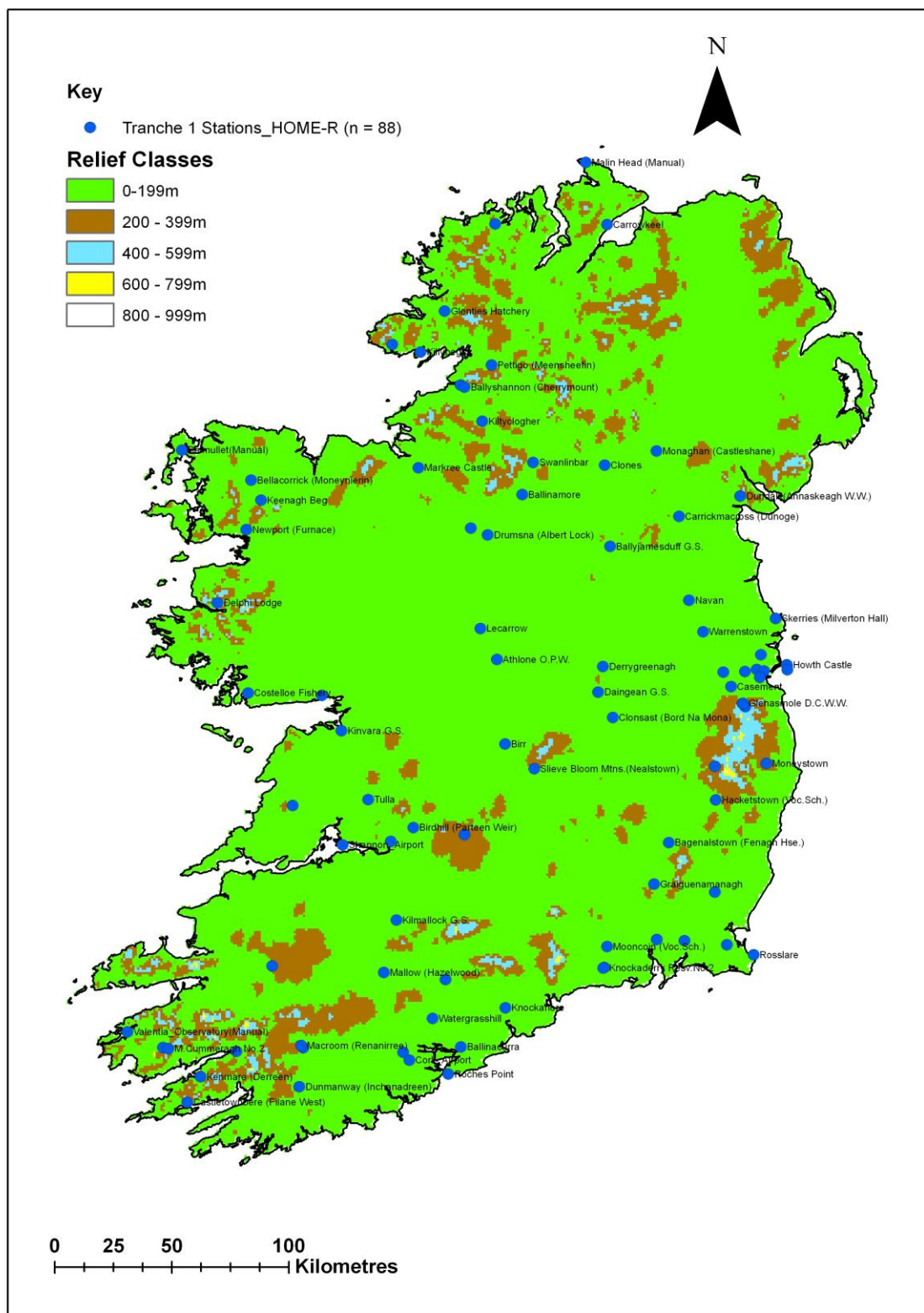
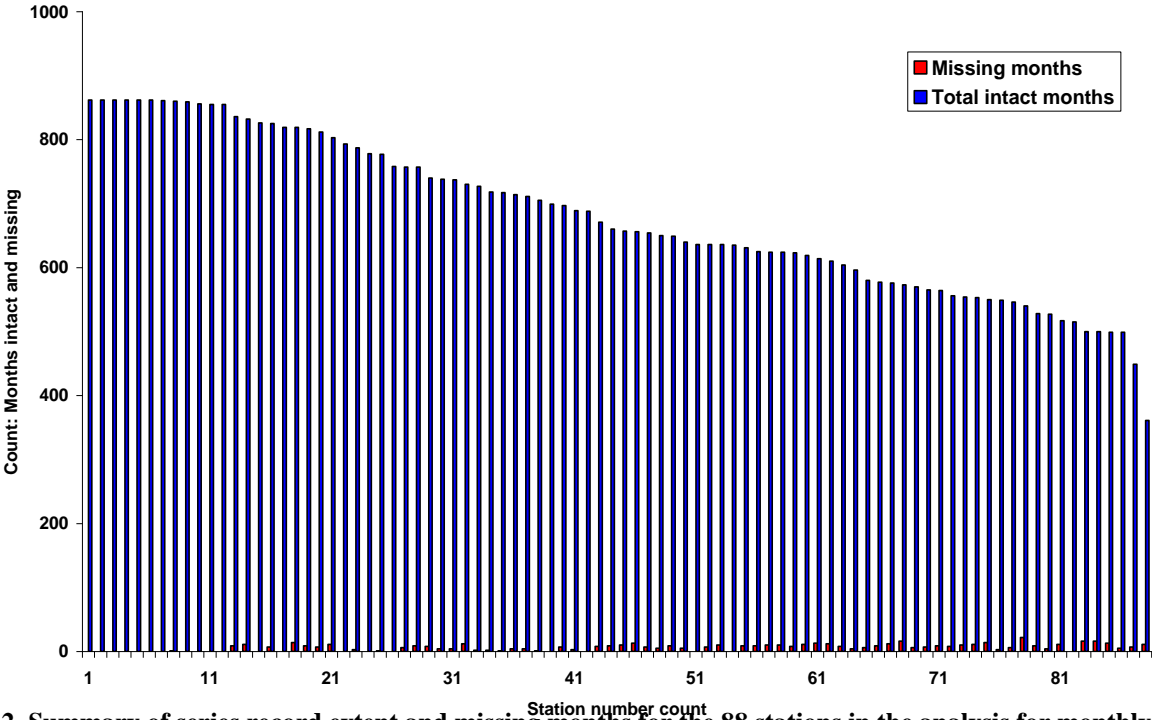


Fig. 1. Annotated map of the island of Ireland showing the selected Met *Éireann* monthly station locations. The locations of the Precipitation stations for the series which have been homogenised using the HOMER algorithm are denoted by blue circles. Upland areas are represented by graded brown, blue and yellow shading.

Following initial quality control procedures for the datasets, further routines involved an exploratory statistical analysis of the series to characterise the potential and limitations of the datasets, as well as to identify and correct missing values and outliers. Figure 2 provides an outline summary of the extent of the number of intact months and missing records for the station series. Issues with missing data in climate time series can be tackled with temporal interpolation using data from the same series before and after the gap, or with spatial interpolation using data from nearby stations (WMO, 2011). Recent work on establishing new 1981-2010 Long Term Averages (LTAs) for Ireland involved the implementation of comprehensive quality control procedures on all Met Éireann’s digital temperature and rainfall data (Walsh, 2012a, 2013); and the extension of this work involved backfilling the available precipitation records to 1941 (S. Walsh, pers. comm.). Complex estimation methods such as weighted averages, spline functions, linear regression and kriging which take into account the correlations with other elements can also be used to complete time series (Frei et al. 2013); and elements of each were used to deal with missing values for individual stations in the construction of the new LTAs (Walsh, 2012).



**Fig. 2. Summary of series record extent and missing months for the 88 stations in the analysis for monthly records spanning the period 1941 – 2010. The blue columns show the overall number of intact monthly records for the station series (n = 88); the red columns show the corresponding distribution of missing records across the station series.**

**2.3. Derivation of reference time series**

Two groups of homogeneity testing techniques can be distinguished and are usually referred to as ‘absolute’ and ‘relative’ methods; in the first set of procedures, the statistical tests are applied to each station data separately, whereas in the second set, the testing procedures use records from neighbouring reference stations which are assumed to be homogeneous (Costa and Soares, 2009). Reference series or reference sections are used in detection procedures in

many homogenisation methods (WMO, 2011), as well as being used to assess the quality of the homogenisation (Kuglitsch et al. 2009). These reference series do not need to be homogeneous (Szentimrey, 1999; Zhang et al. 2001; Causinus and Mestre, 2004), but must encompass the same climatic signal as the candidate series (Della-Marta and Wanner, 2006). Therefore a reference time series has ideally experienced all of the broad climatic influences of the candidate, but none of its artificial biases (WMO, 2011). In practise however, the two fundamental problems of homogenisation are that the nearby stations are also inhomogeneous and that typically more than one break is present (Lindau and Venema, 2013). In addition, the time series data collected at all sites within the same climatic region should be highly correlated, have similar variability, and differ only by scaling factors and random sampling variability (Costa and Soares, 2009).

Selection procedures for the surrounding stations to produce the reference series can be based on the distance between stations or on the correlation coefficients between candidate and potential time series, although there are advantages and disadvantages associated with both methods. Distance based methods will preserve geographic proximity, but time series from nearby stations with different climatic signals (for example due to a difference in elevation) can be selected. Whereas, when using more highly correlated neighbour time series, both the candidate and reference series will present similar variability, but station series with similar or coincident inhomogeneities with the candidate can be selected (Stepanek and Mikulova, 2008). Problems arise when the inhomogeneities in the climate data series are caused by simultaneous changes in the observational network, such as simultaneous changes in the measuring technique, as relative tests become insensitive since all series are affected at the same time (Tuomenvirta 2001; Wijngaard et al. 2003). Furthermore, ambiguous conclusions are possible when several neighbouring stations have inhomogeneities themselves (Reeves et al. 2007; Tayanç et al. 1998).

The most common approach for selecting reference stations is applying Pearson correlation matrices to establish the relationship between the candidate site and potential neighbour station data, and to take as reference the most closely correlated series (Boissonnade et al. 2002; Tayanç et al. 1998). A Pearson cross-correlation was done as a standard exercise for the 88 station series here, although the results are not reported. However, and as a more refined exercise prior to the application of HOMER, reference series were also produced using the first difference correlation coefficients for the series, this followed suggestions in, e.g.; Alexandersson and Moberg, 1997; Peterson et al. 1998; Stepanek and Mikulova, 2008 and Domonkos et al. 2012. A first difference series is made by subtracting year 1's observation from year 2, year 2 from year 3, etc. The correlation then is a measure of the similarity in year-to-year changes, and an inhomogeneity only impacts one observation rather than making all observations after the inhomogeneity artificially warmer or colder (WMO, 2011). By contrast, for the analysis using HOMER, the geographical distance option was selected to allow for comparison of the results obtained via the two approaches. A results comparison between both approaches for four regional case study base and neighbour series are provided in Section 3.3.

## 2.4. Application of HOMER

Features of the HOMER software were then implemented to detect and correct the inhomogeneities in the monthly Precip datasets for the 88 stations for the period 1941-2010. The software is one of the most recently developed for homogenisation, and was made available following a comparative analysis of the best available homogenisation algorithms performed within the COST Action ES0601 (HOME) (Venema et al. 2012). HOMER incorporates additional functions to perform fast quality control of the data, including functions of the CLIMATOL R package which allows the user to estimate the station density, correlogram, histograms, box plots, and cluster analysis (Guijarro, 2011). For the detection of heterogeneities in the datasets HOMER combines three detection algorithms: pairwise-univariate detection, joint detection and ACMANT-bivariate detection, and corrects the datasets using ANOVA (Mestre et al. 2013). However, the ACMANT detection functions are not applied for the homogenisation of Precipitation series within HOMER. If the precise month of change is not known, the default is to validate the break at the end of the year, since detection is mainly performed on annual indices (Mestre et al. 2013).

The models used in HOMER for imputation of missing data and for outlier correction are presented in Mestre et al. 2013. In these models missing datasets are corrected using ANOVA and outliers are detected by pairwise comparison of different time series between candidate and best neighbour time series. HOMER is an interactive semi-automatic method and in applying HOMER, users may choose between the fully automatic cghseg detection results, or a partly subjective pairwise comparison technique that is adapted from PRODIGE (Mestre et al. 2013). Therefore there is freedom for users to add subjective decisions based on metadata or research experience. Basic quality control and network analysis are adapted from CLIMATOL (Guijarro, 2011), and overall HOMER incorporates the best features of some other state-of-the-art methods (Mestre et al. 2013).

For the homogenisation of the 88 records reported here the pairwise comparison option was implemented alongside scrutiny of the station metadata, and thus represents the comparison technique adapted from PRODIGE. This choice also reflects, that in general, a combination of statistical methods and methods relying on metadata information is considered to be the most effective in detecting inhomogeneities (Wijngaard et al. 2003). As a final result HOMER provides corrected series with respect to inhomogeneities and missing values using multiple comparison and ANOVA respectively (Freitas et al. 2013).

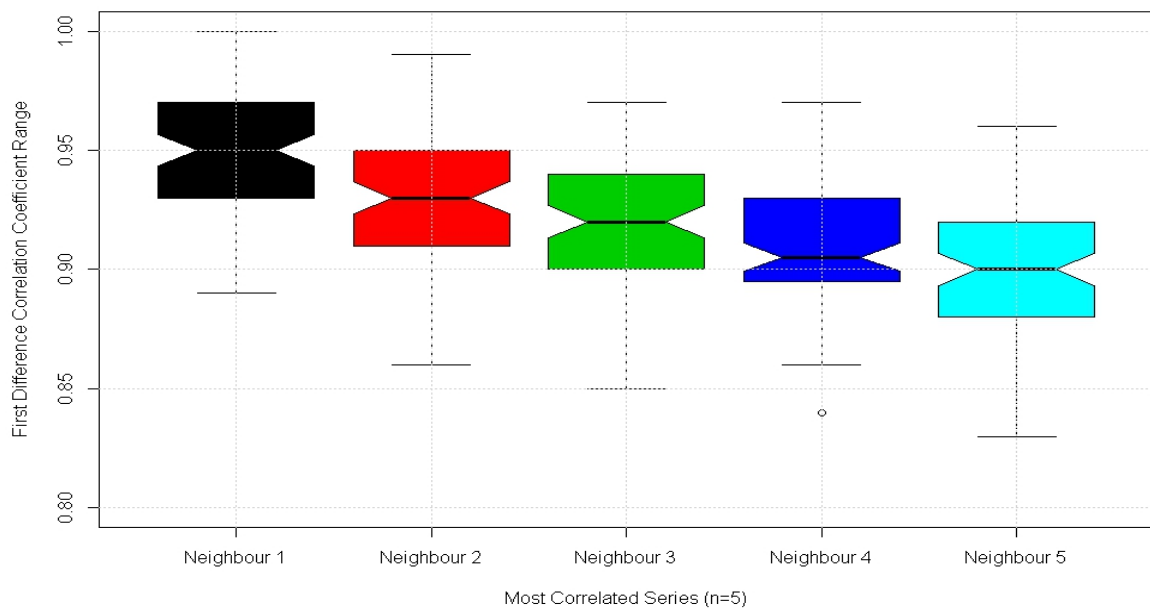
## 3. RESULTS

### 3.1. First difference correlation derived neighbour series

The identification of appropriate neighbour series allow searching for a change which occurs only in the base series. However, when the neighbour series are spatially distant from the base series or have a high elevation difference, the inhomogeneities in the base series may be hidden in the noise between the series due to the large spatial variability between the series. This is especially true for a noisy space-time variable such as precipitation. The first

difference correlations were computed using the infilled station time series (Walsh 2012) for 1941 – 2010, the common period of analysis for all the series.

The disparities between correlations computed from first difference series can be interpreted as a likely indicator of the presence of potential inhomogeneities in the base or in the reference series (Vincent, 1998). High correlation coefficients derived from the first difference series are likely to indicate a strong relationship between the reference series and base series without inhomogeneities. Whereas lower first difference correlation coefficients are more likely to indicate that the base or the neighbour series may contain inhomogeneities. Consequently, the variation in correlation between stations is likely to be a good indicator for the performance of the homogenisation methods, which are expected to perform better when the base and neighbour series are highly correlated. Therefore for deriving neighbour networks based on correlation, the five most closely correlated series were selected as calculated from the first difference series. If other statistical properties of the series, and in particular the variances and the data range are similar, these provide a further indicator on the likely performance for homogenisation. Hence box plots (as a useful summary of statistical properties) are the favoured means for summarising and communicating series characteristics in much of this report.



**Fig. 3. Waisted box plots summarising the first difference correlation coefficient range for the base series and their five most closely correlated neighbours. The plots describe a summary analysis for 88 station and 440 potentially available neighbour series. Boxes: interquartile range; whiskers 5th and 95th percentiles.**

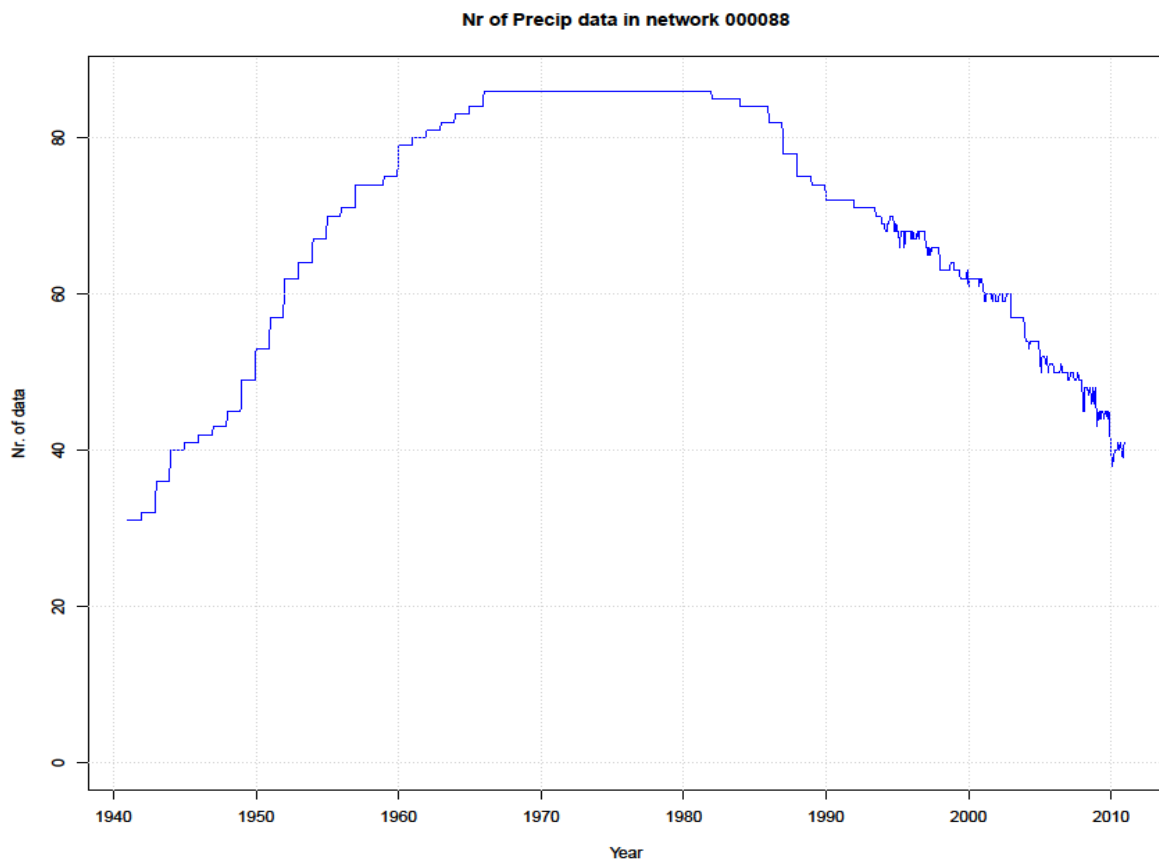
For the correlation exercise no limit for geographical distance or elevation difference between the stations was applied, although these relationships were checked in a subsequent GIS-based scrutiny of the candidate reference networks. The first difference correlation coefficient range for all 88 base station series and their five most closely correlated neighbours are summarised for all the station series used in the analysis in Figure 3. With a correlation coefficient range from 0.83 (least correlated neighbour series) to 1.0 (most



correlated neighbour series) and a mean correlation coefficient of 0.92 for all 88 stations and 440 potentially available correlated series from corresponding neighbour series, the indication is that major inhomogeneities among the series are unlikely. For a small and predominantly maritime-influenced country such as Ireland this is largely to be expected, at least for data at a monthly temporal scale where much of the daily and hourly synoptic and local terrain-induced variations will be smoothed out.

### 3.2. HOMER geographical distance derived neighbour series

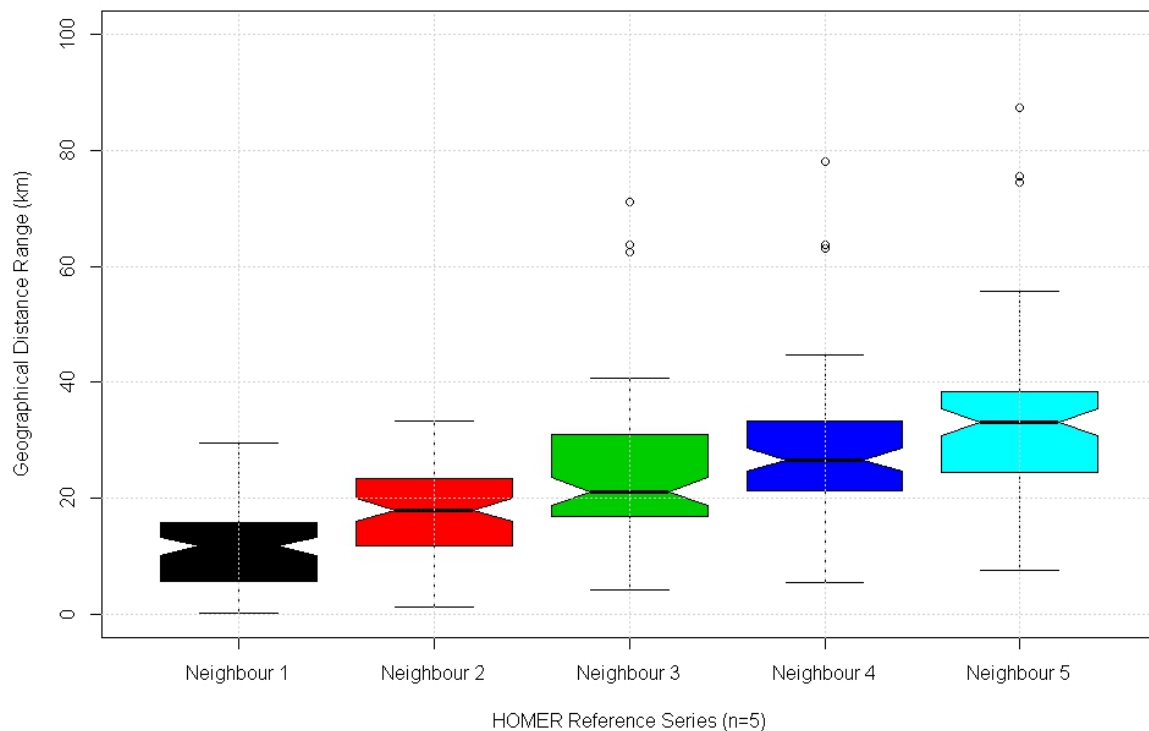
In contrast to the infilled series records used for the first difference correlation exercise, for the derivation of neighbour series in HOMER all the missing monthly values (1941-2010) were left in for the 88 station series with the widely used default entry value of -999.9. This reflects the inbuilt functionality of the software to deal with missing values depending on their distribution within the base and neighbour series. The distribution of years where most intact records were included in the HOMER analysis may be deduced from Figure 4 and clearly show the peak in available records between the late 1960s and early 1980s.



**Fig. 4. HOMER diagnostic summary illustrating the overall number of intact and absent months for the 88 station series records included in the 1941-2010 analysis. The blue line plot indicates the greater number of available records for nearly all station series in the 1960s-1980s compared to the earlier and later periods.**

The HOMER-computed geographical distance range for all 88 station base series and their five closest neighbours with similar series characteristics are summarised for all the station

series used in the analysis in Figure 5. As the algorithm iterates through, it is clear that the distance (kilometres) between base and neighbour series steadily increase as the algorithm searches for series with similar characteristics. The increasing distances associated with neighbour series 3 to 5 are apparent and some of the outlying geographical distance values are noteworthy; it may be the case that these increasing distances may have implications for homogeneity among some of the base and neighbour series identified within HOMER. As such this is an issue that will warrant further investigation as the work proceeds, particularly as new and different station series combinations are added for ongoing and future analysis.

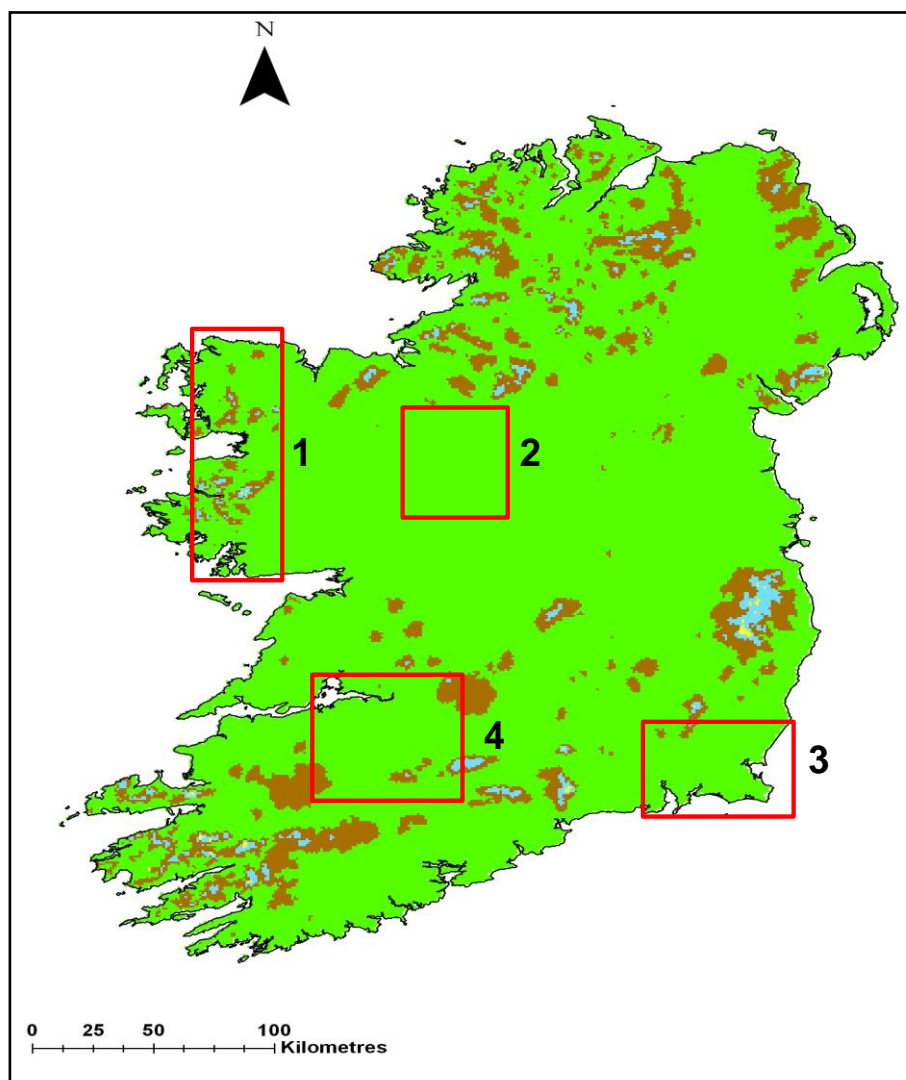


**Fig. 5. Waisted box plots summarising the HOMER-derived geographical distance range for the base series and their five most closely correlated neighbours. The plots describe a summary analysis for 88 station and 440 potentially available neighbour series. Boxes: interquartile range; whiskers 5th and 95th percentiles.**

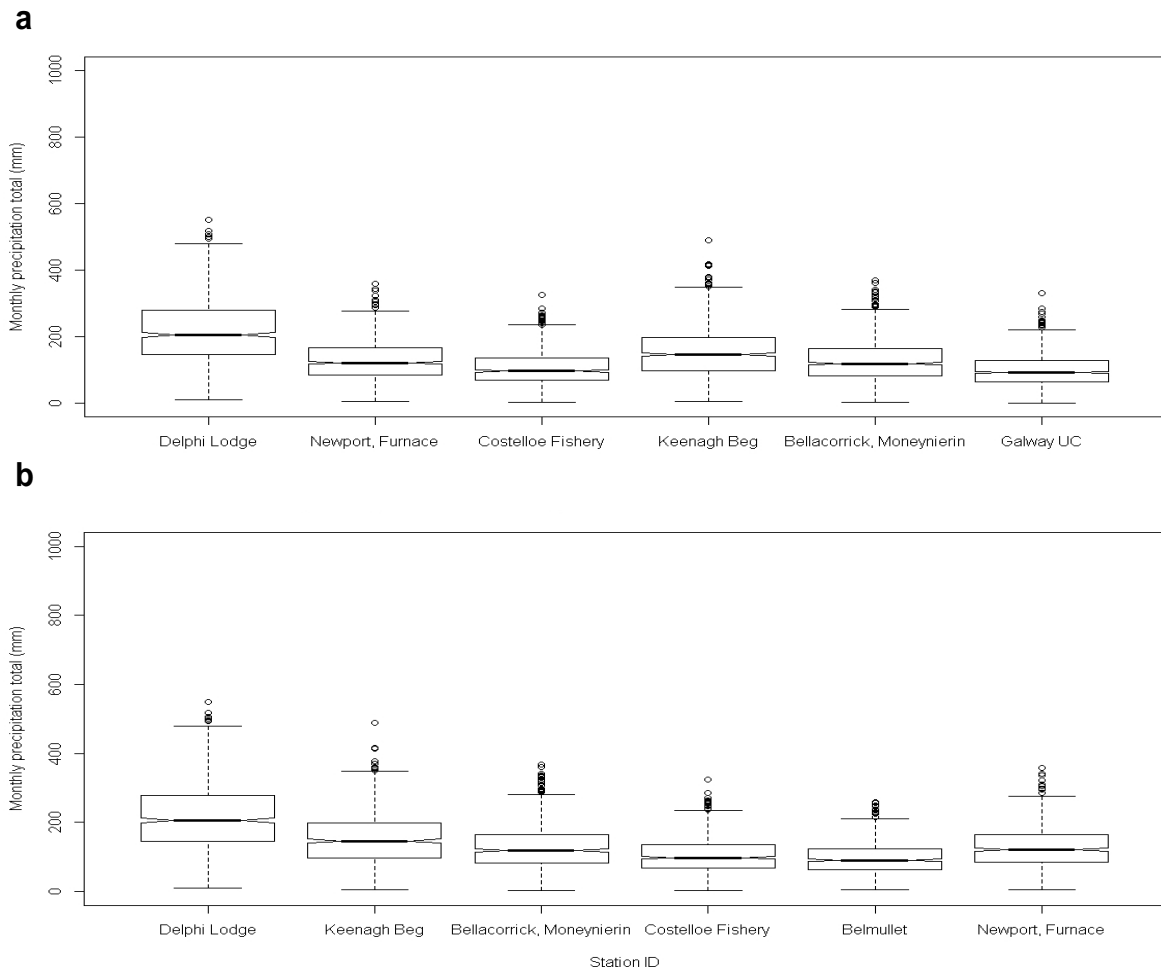
For example if some additional station records are added to the overall analysis (even if the available contiguous records are shorter) it will be the case that the geographical distances between base and new potential candidate neighbour series will change if a denser station series network can be constructed and used effectively. It is already recognised that if too many distant (or less correlated) neighbouring stations are used, the resulting reference may not reflect the true climatic signal of the candidate station properly (Boissonnade et al. 2002). These difficulties may increase dramatically with the increase in spatial variability of the data caused by the inherent variability of the element (e.g. precipitation) or the time series resolution (e.g. monthly data) (Costa and Soares, 2009), both conditions clearly pertain to the data analysed here.

### 3.3. Comparison of methods: First difference correlations and HOMER geographical distances for selected case study locations

Selective results are presented based on the comparative analysis between the two methods to arrive at neighbour series which are both spatially and statistically coherent. Figure 6 outlines the four case study regions for the selected locations, broadly these describe regional study locations for the west, west Midlands, south-east and south-west of the country. For the western case study example both approaches identified four reference series in common; Keenagh Beg, Bellacorrick, Costelloe Fishery and Newport (Figure 7). For both sets of results, the statistical properties of the base and neighbour series are largely coherent, although Delphi Lodge is a wetter station than the neighbours. First difference correlation coefficients ranged from  $r = 0.88 - 0.92$ , whereas geographical distances were 27.57 km – 48.35 km.

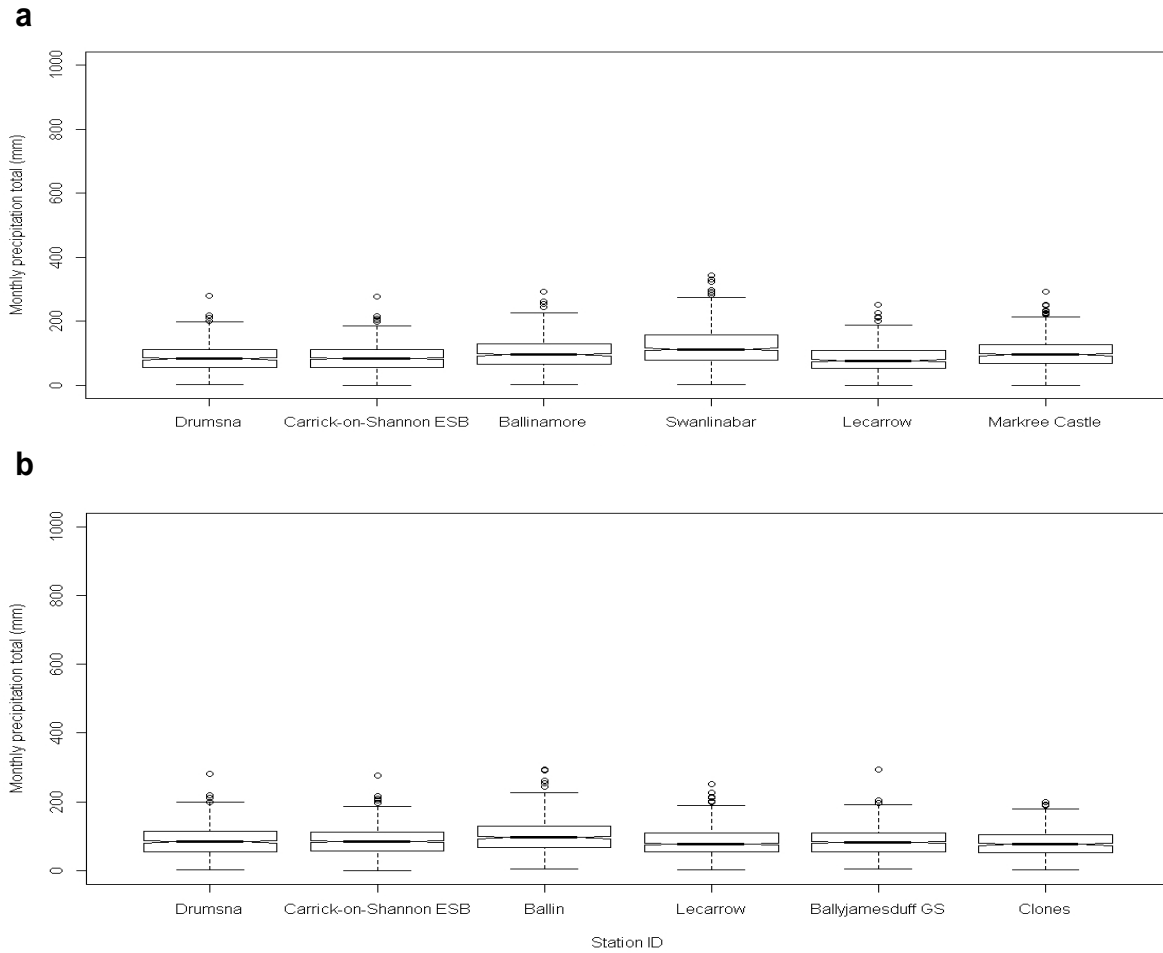


**Fig. 6. Annotated map of the island of Ireland showing the case study locations for selective results comparison. Upland areas are represented by graded brown, blue and yellow shading. The red squares denote the case study areas describing the derivation of neighbour for base series using; a) first difference correlation coefficients and b) geographical distances within the HOMER software. Box 1 outlines the western, Box 2 the west Midlands, Box 3 the south-east and Box 4 the south-west study regions respectively.**

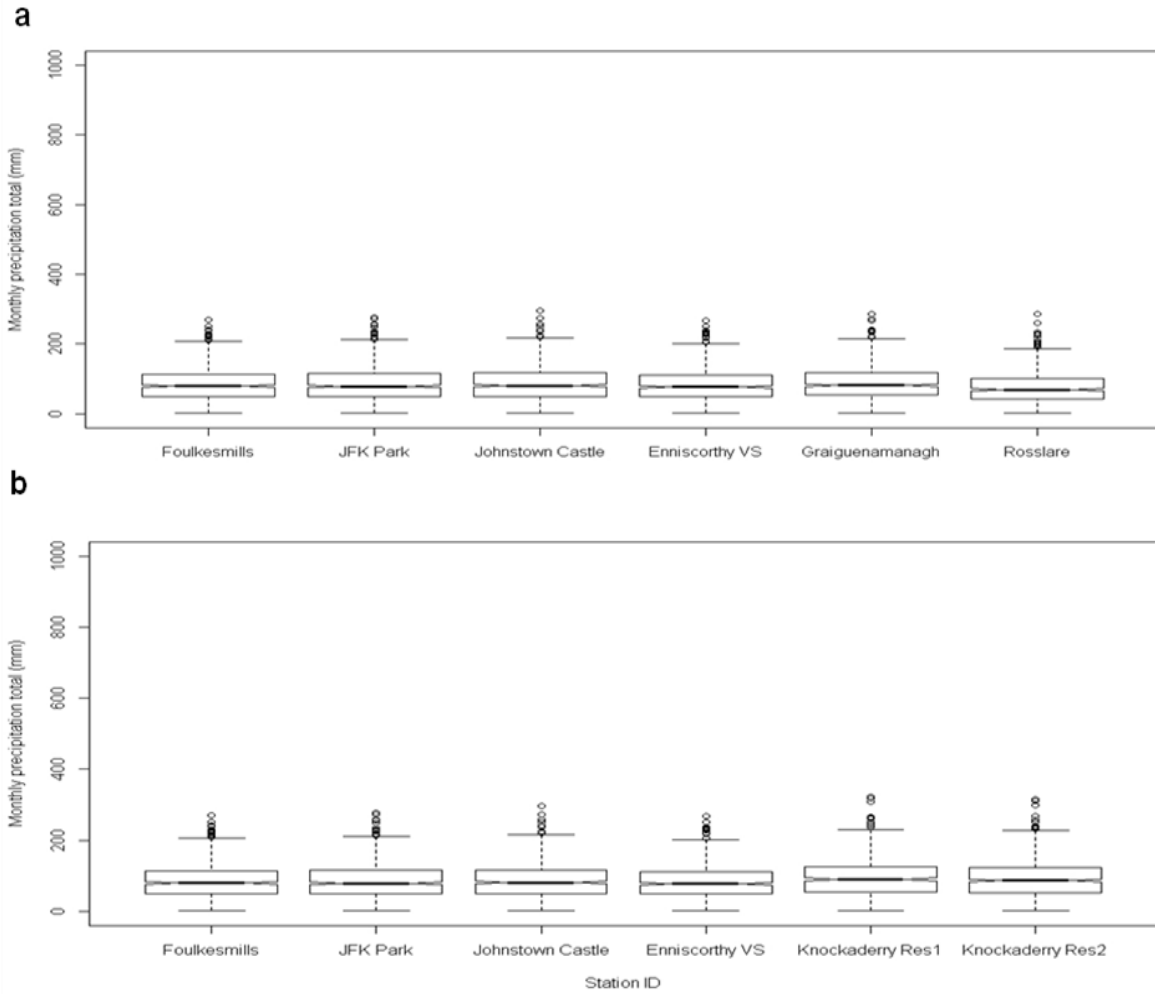


**Fig. 7. Box plot comparison of the base Delphi Lodge and neighbour series derived via; a) geographical distance in HOMER, and b) first difference correlation for the nearest five neighbours. Boxes: interquartile range; whiskers 5th and 95th percentiles.**

In the west Midlands case study example both approaches identified three reference series in common; Carrick-on-Shannon, Ballinamore, and Lecarrow (Figure 8). For both sets of results, the statistical properties of the base and neighbour series are largely coherent, although Swanlinbar as identified via HOMER is a wetter station than the neighbours. First difference correlation coefficients ranged from  $r = 0.88 - 0.95$ , whereas geographical distances were 6.36 km – 33.62 km. Similarly, for the south-east case study example both approaches identified three reference series in common; JFK Park, Johnstown Castle and Enniscorthy (Figure 9). For both sets of results, the statistical properties of the base and neighbour series are very similar. First difference correlation coefficients ranged from  $r = 0.93 - 0.97$ , whereas geographical distances were 9.50 km – 24.29 km.

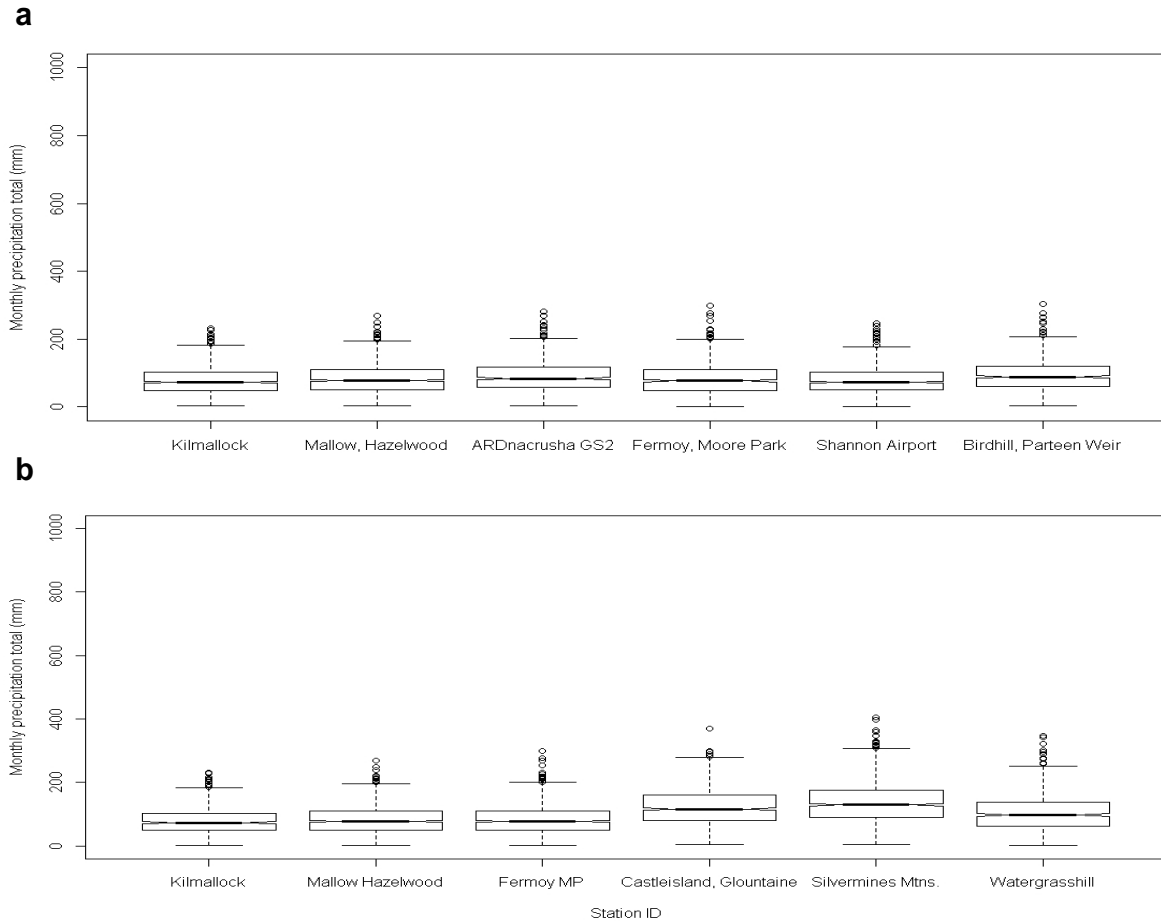


**Fig. 8. Box plot comparison of the base Drumsna (Albert Lock) and neighbour series derived via; a) geographical distance in HOMER, and b) first difference correlation for the nearest five neighbours. Boxes: interquartile range; whiskers 5th and 95th percentiles.**



**Fig. 9. Box plot comparison of the base Foulkesmills and neighbour series derived via; a) geographical distance in HOMER, and b) first difference correlation for the nearest five neighbours. Boxes: interquartile range; whiskers 5th and 95th percentiles.**

Whereas for the south-west case study example both approaches identified two reference series in common; Mallow and Fermoy (Figure 10). For both sets of results, the statistical properties of the base and neighbour series are largely coherent, although Silvermines Mountains as identified by the first difference correlations is a slightly wetter station than the neighbours. First difference correlation coefficients ranged from  $r = 0.86 - 0.94$ , whereas geographical distances were 19.07 km – 33.04 km. Although not shown here, a GIS-based check confirms results obtained via both approaches which are both spatially and climatologically coherent for each of the study regions.



**Fig. 10. Box plot comparison of the base Kilmallock and neighbour series derived via; a) geographical distance in HOMER, and b) first difference correlation for the nearest five neighbours. Boxes: interquartile range; whiskers 5th and 95th percentiles.**

### 3.4. Break detection and homogenisation of the base series in HOMER

Although only shown in the context of some regional case study examples in the previous section, for all 88 of the candidate base series, results obtained between both methods identified largely similar reference stations and series. In addition, both the statistical properties of the reference series and the geographical location of stations were coherent climatologically.

On this basis the full pairwise detection option using geographical distance as the specified metric for the reference series was implemented in HOMER, and the breaks detected were interpreted alongside the available metadata for the stations. Breaks were detected in 11 from the 88 base series analysed, with multiple breaks detected in three of the series, results are summarised in Table 1.

*Table 1.* HOMER break detection summary for the 11 stations with consistent breaks in the base and neighbour series. Break years are provided with the months in paranthesis.

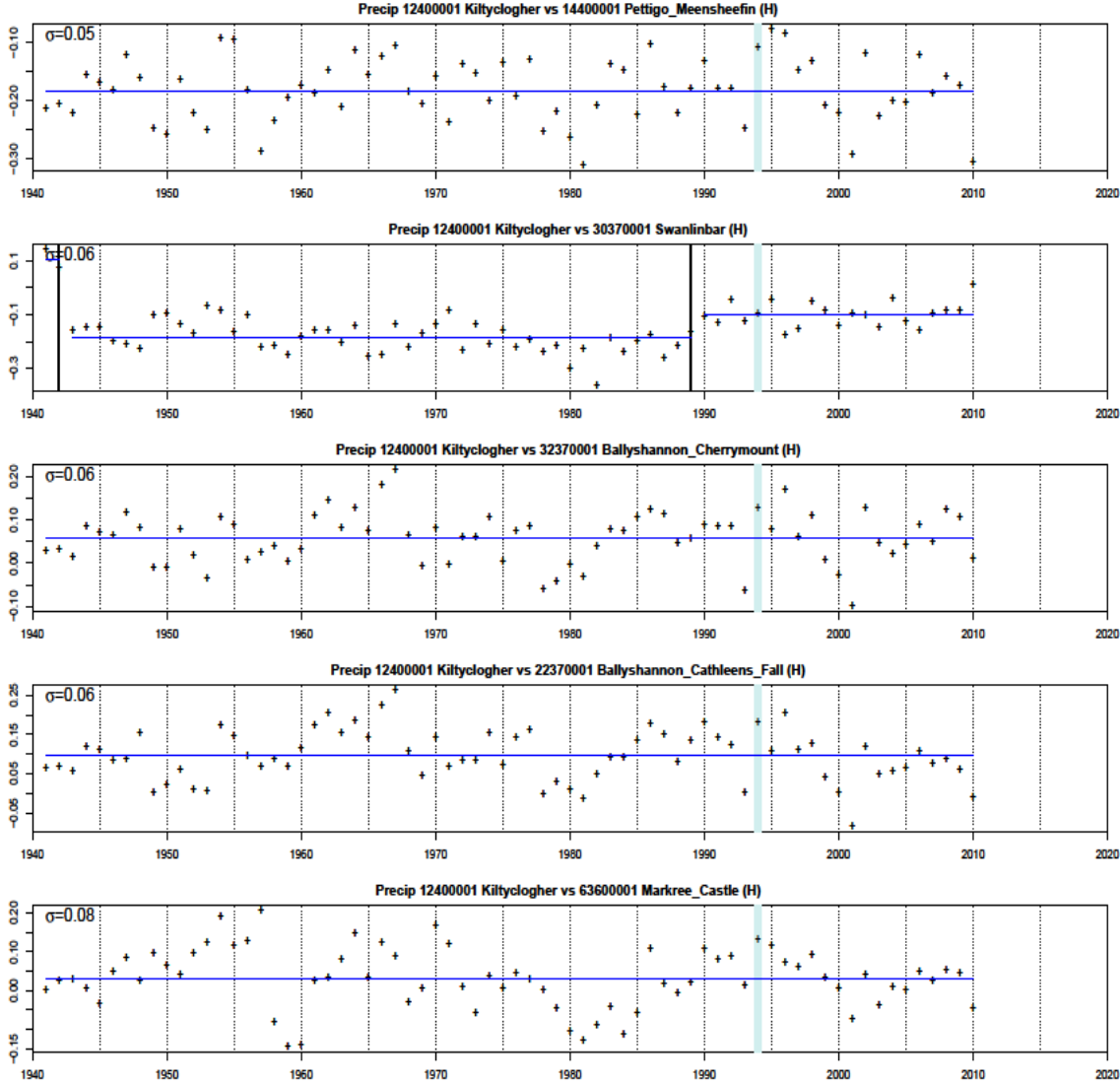
<b>Station Name and ID</b>	<b>Break Years (Month)</b>	<b>Metadata verified</b>
Delphi Lodge (626)	1976 (12), 1980 (12)	No
Creelough (944)	2000 (12), 2007 (12)	No
Roches Point (1004)	1955 (11)	Yes
Kiltyclogher (1240)	1994 (5)	Yes
Drumsna, Albert Lock (1529)	1954 (7), 1965 (8), 1969 (1)	Yes
Knockaderry Reservoir (1) (1712)	1974 (12)	No
Watergrasshill (2404)	1974 (11)	Yes
Howth, Danesforth (3023)	1966 (12)	No
Slieve Bloom Mountains (3513)	1971 (3)	Yes
Macroom, Renanirree (3804)	1978 (11)	Yes
Silvermines Mountains (4819)	2005 (12)	No

Station metadata checks revealed the cause of breaks for 6 from the 11 series as summarised in Table 1. These were:

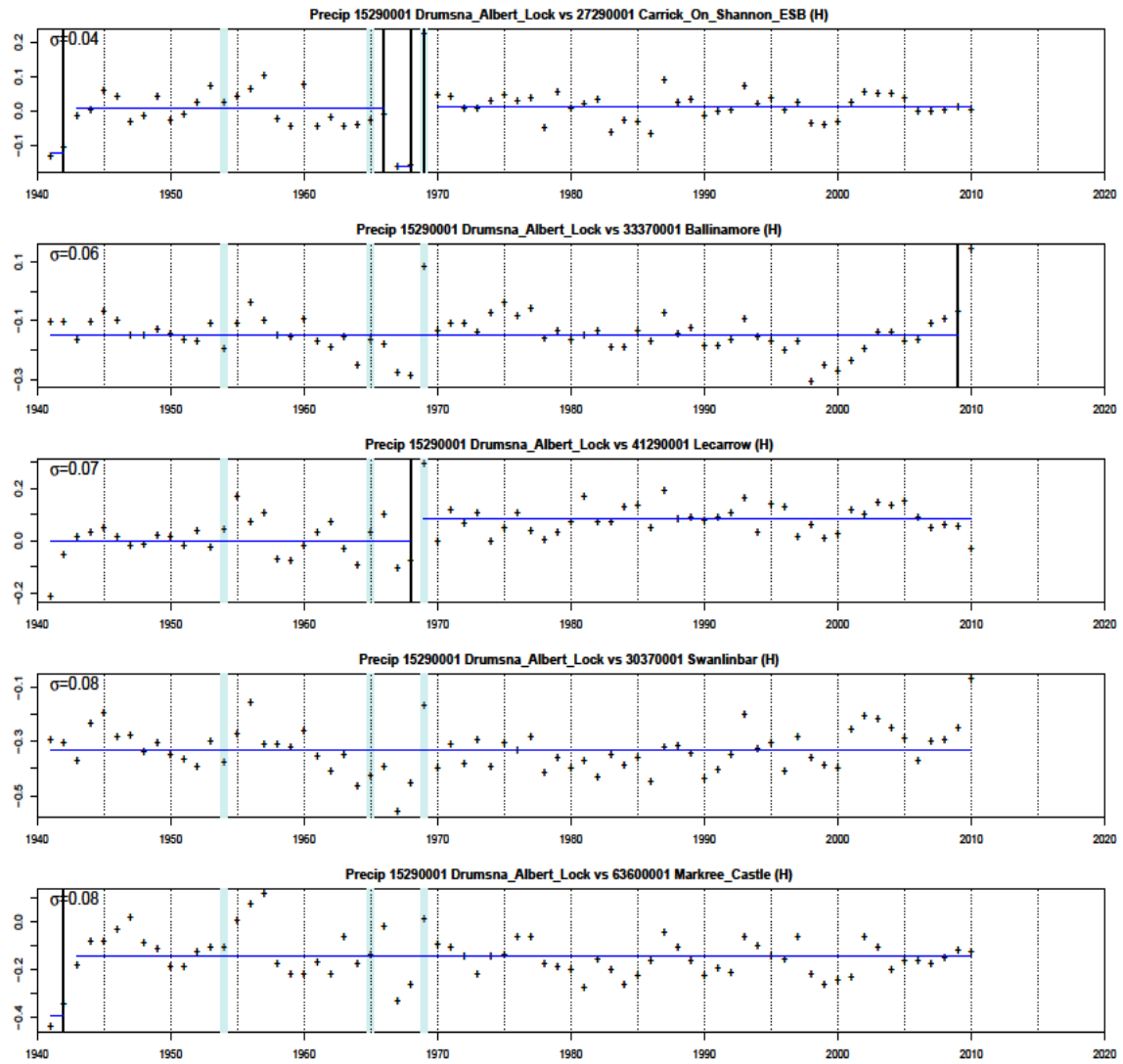
- Roches Point, November 1955 – station moved to a new elevation;
- Kiltyclogher, May 1994 – gauge replacement;
- Drumsna, July 1954 – new gauge after 7 year gap; August 1965 – new mm measure introduced; January 1969 – defective gauge replaced;
- Watergrasshill, November 1974 – change of observer;
- Slieve Bloom Mountains, March 1971 – station site move;
- Macroom, November 1978 – defective gauge replaced.



In the interests of brevity, sample output results are presented for only two of the base and neighbour series pairwise detection plots for the HOMER homogenised data. These illustrate the metadata verified single and multiple breaks in the Kiltyclogher and Drumsna series respectively (Figures 11 and 12).



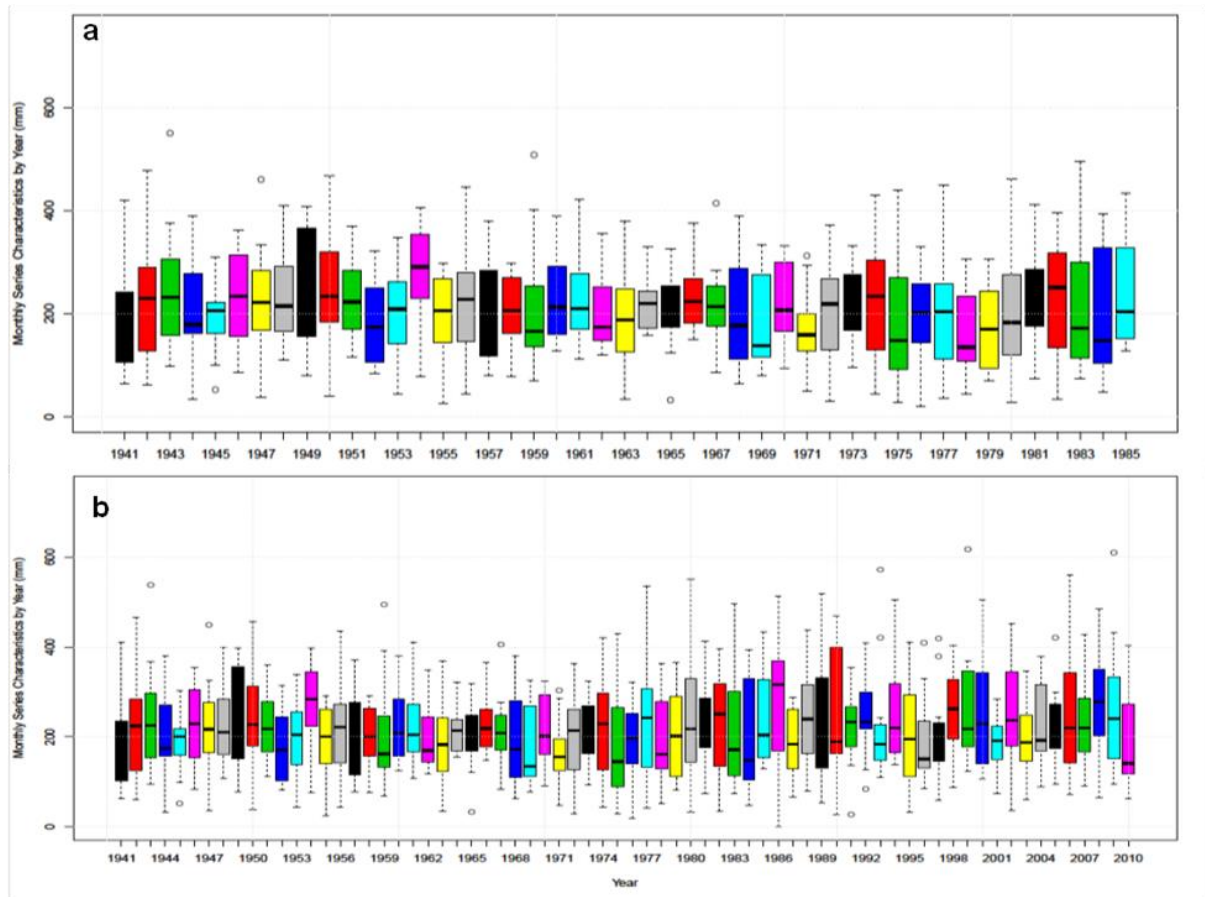
**Fig. 11. HOMER screen capture of the metadata verified break in the homogenised Kiltyclogher annual base series and the corresponding neighbour series. The sky blue lines denote the consistent single break detection across the base and neighbour series for May 1994.**



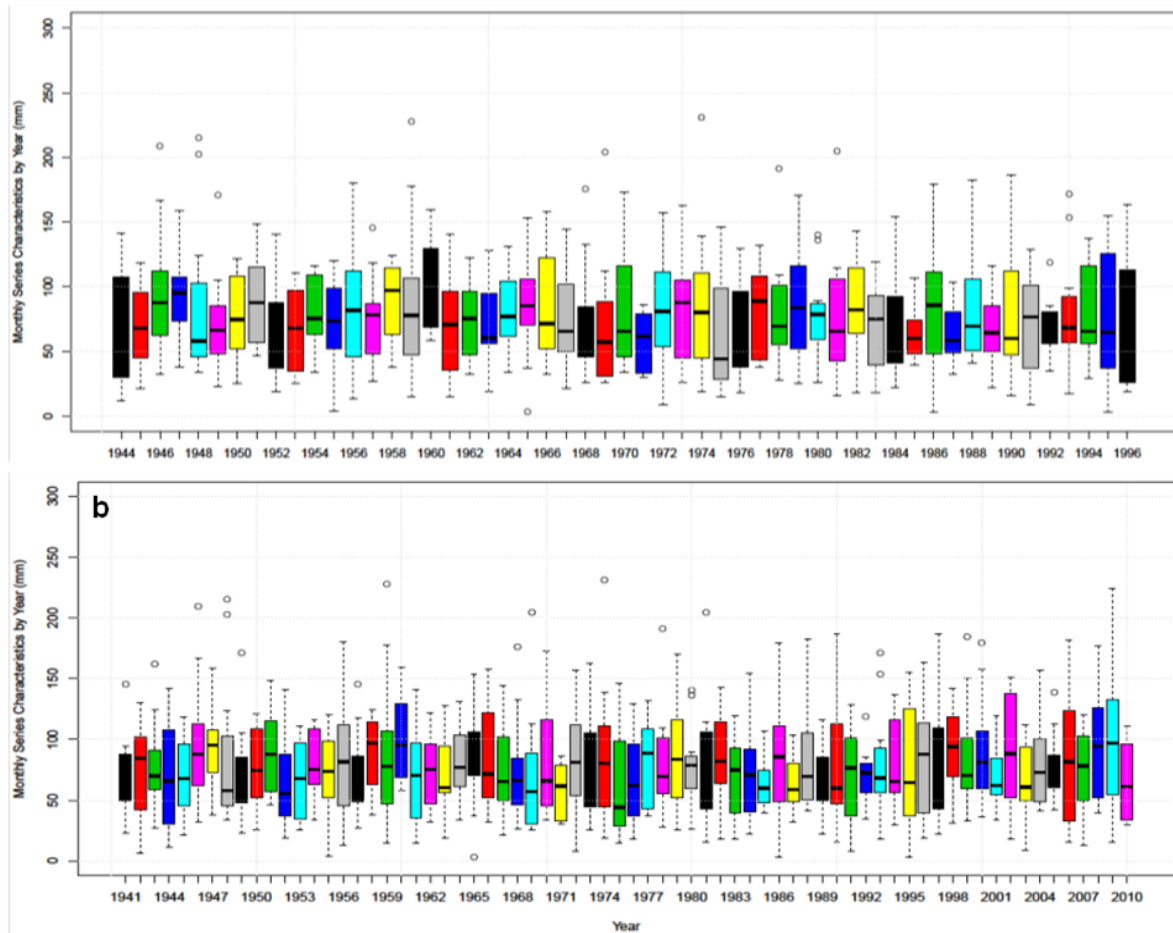
**Fig. 12.** HOMER screen capture of the metadata verified break in the homogenised Drumsna annual base series and the corresponding neighbour series. The sky blue lines denote the consistent multiple break detection across the base and neighbour series in 1954, 1965 and 1969.

### 3.5. Sample results: HOMER homogenised series

Based on the neighbour series selected on the basis of geographical distance, and after the metadata checks and user interaction with the software, HOMER provided a new and homogenised output for the focal base series based on the PRODIGE comparison technique. Examples are provided here for two of the case study series where the original series have been extended using HOMER-derived values based on the neighbour series, Delphi Lodge (Figure 13) and Kilmallock (Figure 14).



**Fig. 13.** Time series box plots of a) the available Delphi Lodge annual series (1941-1985), and b) of the HOMER homogenised data for the same series extended with reference to the neighbour data (1941-2010). Boxes: interquartile range; whiskers 5th and 95th percentiles.



**Fig. 14. Time series box plots of the a) available Kilmallock annual series (1944-1996), and b) of the HOMER homogenised data for the same series extended with reference to the neighbour data (1941-2010). Boxes: interquartile range; whiskers 5th and 95th percentiles.**

Delphi Lodge and Kilmallock are presented as sample output in favour of Foulkesmills and Drumsna as both these latter series comprised more or less fully intact data (1941-2010 and 1941-2008 respectively) prior to the homogenisation exercise. There was therefore less scope to illustrate the imputation of missing values by HOMER compared to the greater number of missing years in the Delphi Lodge and Kilmallock series. The data are presented as annual time series box plots to highlight that the HOMER-derived missing values are statistically coherent with the data range and other characteristics of the original series. In the examples here, for the 1986-2010 segment for Delphi Lodge and the 1997-2010 segment for Kilmallock, the additions are HOMER-derived extensions based on neighbour station data. The different Y-axis scales between the Figures are to account for different total precipitation receipts at both station locations.

However, for Delphi Lodge, it is noted that there are more outlying values associated with the neighbour-derived infill segment compared to the original series. During the homogenisation work, the Delphi Lodge series was noted as a case where further work is required and a series which may benefit from future work which identifies more statistically similar neighbour series as it is a wetter station than the neighbour series used here.

## 4. DISCUSSION AND CONCLUSIONS

As much of the preceding content reports on and discusses results to date, as befits an analysis using case study examples to report on interim results on a first application of HOMER to Irish precipitation data, discussion here is kept brief. First difference correlations were used to critically evaluate HOMER results output in a parallel statistical computing framework. In turn this builds on other approaches to exploratory analysis of the station series and their inter-relationships both statistically and spatially not reported on here. HOMER has accurately detected breaks in 11 from the 88 monthly precipitation series so far analysed in this approach, of which 6 are confirmed from the metadata. The results indicate that both approaches yield valid and statistically similar corresponding neighbour series, and although only shown in the context of case study examples in the results presented here, results across the analysis for all 88 station records and their potential neighbour series indicate that both first difference correlations and HOMER geographical distance selections yield overlapping neighbour series which are largely statistically and spatially coherent.

Ongoing analysis which is extending some of the functionality the HOMER algorithm offers indicate very good prospects for the work going forward. To date there are relatively few openly available comparative results as HOMER is still new and not in wide use, and any published work to date relates to applying the algorithm for the homogenisation of temperature (i.e. Freitas et al. 2013; Mestre et al. 2013; Luhunga et al. 2014). Therefore, and insofar as we are aware, this remains one of the most substantial applications and tests of the algorithm in relation to precipitation data, and hence the ongoing work remains novel and is certainly the first such application for Irish monthly precipitation series.

### Acknowledgements

We thank the Irish Environmental Protection Agency (EPA) for the project funding and support. This research was supported under grant 2012-CCRP-FS.11.

### References

- Alexandersson H (1986). A Homogeneity Test Applied to Precipitation Data. *International Journal of Climatology*, 6 (6); 661-675.
- Alexandersson H and Moberg A (1997). Homogenisation of Swedish temperature data. 1. Homogeneity test for linear trends. *International Journal of Climatology*, 17; 25-34.
- Beaulieu C, Seidou O, Ouarda, TMBJ and Zhang, X (2009). Intercomparison of homogenization techniques for precipitation data continued: Comparison of two recent Bayesian change point models. *Water Resources Research*, 45:W08410.
- Boissonnade AC, Heitkemper LJ, Whitehead D (2002). Weather data: cleaning and enhancement. In: Dischel RS (ed) *Climate risk and the weather market: financial risk management with weather hedges. Risk Waters*, pp 73–93.

- Cao LJ, and Yan ZW (2012). Progress in Research on Homogenization of Climate Data. *Advances in Climate Change Research*, 3 (2); 59-67.
- Caussinus, H and Mestre O (2004). Detection and correction of artificial shifts in climate series. *Applied Statistics*, 53: 405 – 425.
- Costa, AC Soares, A (2009). Homogenization of Climate Data: Review and New Perspectives Using Geostatistics. *Mathematical Geosciences* (2009) 41; 291–305
- Della-Marta, PM and Warner, H (2006). A method for homogenising the extremes and mean of daily temperature measurements. *Journal of Climate*, 19: 4179-4197.
- Domonkos, P (2011). Adapted Caussinus-Mestre algorithm for networks of temperature series (ACMANT). *International Journal of Geoscience*, 2: 293-309.
- Environmental Protection Agency (2014). Climate Services Note. Dublin. 4pp.
- Freitas L, Pereira MG , Caramelo L, Mendes M, and Nunes LF (2013). Homogeneity of Monthly Air Temperature in Portugal with HOMER and MASH. *Idojaras*, 117 (1); 69-90.
- Guijarro, JA (2011). Climatol Version 2.0, an R contributed package for homogenisation of climatological series. State Meteorological Agency, Balearic Islands Office, Spain.
- HOME (2013). Homepage of the COST Action ES0601 - Advances in Homogenisation Methods of Climate Series: An Integrated Approach (HOME).
- Jones PD, Raper SCB, Bradley RS, Diaz HF and others (1986). Northern Hemisphere Surface Air Temperature Variations: 1851-1984. *Journal of Climate and Applied Meteorology*, 25 (2) 2; 161-179.
- Kuglitsch, FG, Toreti, A, Xoplaki, E, Della-Marta, PM and others (2009). Homogenisation of daily maximum temperature series in the Mediterranean. *Journal of Geophysical Research*, 114; 1-6.
- Lindau, R and Venema, V (2013). On the multiple breakpoint problem and the number of significant breaks in homogenization of climate records. *Idojaras*, 117 (1); 1-34.
- Luhunga, PM, Mutayoba, E, Ng'ongolo, HK (2014). Homogeneity of Monthly Mean Air Temperature of the United Republic of Tanzania with HOMER. *Atmospheric and Climate Sciences*, 4: 70-77.
- Menne, MJ and Williams, CN (2005). Detection of Undocumented Change-points Using Multiple Test Statistics and Composite Reference Series. *Journal of Climate*, 18 (20); 4271-4286.
- Mestre O, Domonkos P, Picard F, Auer I and others (2013). HOMER: A Homogenization Software - Methods and Applications. *Idojaras*, 117 (1); 47-67.
- Peterson TC, Easterling DR, Karl TR, Groisman P and others (1998). Homogeneity Adjustments of *in Situ* Atmospheric Climate Data: A Review," *International Journal of Climatology*, 18 (13), 1998, 1493-1517.
- Reeves J, Chen J, Wang XL, Lund R, Lu Q (2007). A review and comparison of changepoint detection techniques for climate data. *Journal of Applied Meteorology and Climatology* 46:900–915
- Rohan PK (1986). The climate of Ireland. The Stationery Office, Dublin
- Štěpánek P, and Mikulová K (2008). Homogenisation of air temperature and relative humidity monthly means of individual observation hours in the area of the Czech and Slovak Republic. In: *5<sup>th</sup> Seminar for Homogenisation and Quality Control in Climatological Databases*. Hungarian Met. Service, Budapest; 147-163.
- Tayanç M, Dalfes HN, Karaca M, Yenigün O (1998). A comparative assessment of different methods for detecting inhomogeneities in Turkish temperature data set. *International Journal of Climatology* 18(5); 561–578.
- Toreti, FG, Kuglitsch, A, Xoplaki, E and Luterbacher, J. (2012). A Novel Approach for the Detection of Inhomogeneities Affecting Climate Time Series. *Journal of Applied Meteorology and Climatology*, 51 (2); 317-326.
- Tuomenvirta H (2001). Homogeneity adjustments of temperature and precipitation series – Finnish and Nordic data. *International Journal of Climatology* 21(4); 495–506.
- Venema V, Mestre O, Aguilar E, Auer I and others (2012). Benchmarking Homogenization Algorithms for Monthly Data. *Climate of the Past*, (8): 89-115.

- Vincent, LA (1998). A technique for the identification of inhomogeneities in Canadian temperature series. *Journal of Climate*, 11; 1094 – 1105.
- Walsh, S (2012). *A summary of climate averages for Ireland 1981-2010*. Met Éireann Climatological Note 14. Dublin. 16pp.
- Walsh, S (2013). *Setting the scene: the climate of Ireland 1900-2012*. In: Ireland's climate: the road ahead (Eds: Gleeson, E., McGrath R. and Treanor M.). Met Éireann, Dublin. 102pp.
- Wijngaard J, Klein Tank AMG, Können GP (2003). Homogeneity of 20th century European daily temperature and precipitation series. *International Journal of Climatology* 23(6): 679–692.
- World Meteorological Organisation (2011). *Guide to Climatological Practices*. WMO/No 100, Geneva.
- Zhang, XB, Vincent, LA, Hogg, WD and Niitsoo, A (2001). Temperature and precipitation trends in Canada during the 20<sup>th</sup> century. *Atmosphere-Ocean*, 38: 395-429.

# THE ACMANT2 SOFTWARE PACKAGE

**Peter Domonkos**

Centre for Climate Change (C3), University Rovira i Virgili, Tortosa, Spain  
(peter.domonkos@urv.cat)

## 1. INTRODUCTION

Author has developed the statistical software package ACMANT2, which includes computer programs for the automatic homogenization of mean temperature (Tmean), daily maximum temperature (Tmax), daily minimum temperature (Tmin) and precipitation amounts (PP). The software treats either daily or monthly input, but the detection of inhomogeneities (IH) and the calculation of adjustment terms are always done on the annual or monthly scale, then daily data are adjusted with downscaling the monthly adjustment terms. This study will present the structure and the most important segments of ACMANT2 and will discuss why ACMANT was one of the best performing homogenization method in the international tests of the European project COST ES0601 (its popular name, HOME, will be referred hereafter). In the study, the latest version of ACMANT is referred to as ACMANT2, its previous version as ACMANT1, while its constant properties are often assigned to “ACMANT” without index.

Temperature and precipitation time series can be homogenized with ACMANT2 and the homogenization of these variables is done with very similar algorithms. In this study the description of temperature homogenization is provided only in detail, but the important differences between the algorithms for temperature homogenization and precipitation homogenization will be mentioned.

The organization of the study is as follows: In the next section, the motivation and the brief history of the development of ACMANT is presented. In the third section, the most important theoretical properties of ACMANT2 are shown. In section 4 some efficiency results are shown, while in section 5 the computer programs of the ACMANT2 package and their use are briefly described. The study is supplied with appendixes (AP) with the detailed descriptions of I) the differences between ACMANT1 and ACMANT2 in the homogenization of Tmean and Tmax of mid- or high latitudes; II) the specific rules of Tmin homogenization and homogenization of any temperature variable in tropical regions.

## 2. MOTIVATION AND BRIEF HISTORY OF THE DEVELOPMENT OF ACMANT

During HOME (2007-2011), the new method ACMANT for homogenizing monthly temperature series appeared and was found to be one of the most effective methods just in its first version, leaving relatively small residual errors in the homogenized series.



I have been dealing with testing the efficiency of homogenization methods since 2003 (Domonkos, 2008, 2011a, 2013a, etc.) and I often found large differences between the efficiencies of different methods. Moreover, homogenization sometimes might worsen the quality of observational series (Domonkos, 2013a), therefore a new approach is needed in our general view around the task of time series homogenization. While the review study of Peterson et al. (1998) suggested that time series homogenization is generally recommendable anything is the method applied, the enhanced need for more reliable and more accurate observational data for climate change and climate variability studies forces us to select and use the best performing methods. The international tests with the HOME benchmark dataset (HBM) confirmed that the differences between method efficiencies are large (Venema *et al.*, 2012) and based on these tests only five methods can be recommended for homogenizing monthly temperature and precipitation series, namely MASH (Szentimrey, 1999), PRODIGE (Caussinus and Mestre, 2004), ACMANT (Domonkos, 2011b, referred here D2011), USHCN (Menne and Williams, 2009) and the Craddock-test (Craddock, 1979). In the last stage of HOME, the HOMER method was created (Mestre *et al.*, 2013) from the best performing segments of PRODIGE and ACMANT and incorporating the network wide joint segmentation method (Picard *et al.*, 2011). After HOME, the climatologist community still has important tasks in continuing test experiments (Domonkos, 2013b), since the efficiencies measured by HOME are based on the use of a not very large benchmark dataset, i.e. 15 networks for testing each of temperature and precipitation homogenization methods (Venema *et al.*, 2012). On the other hand, HOME recognized several weak points of the tested homogenization methods, fostering a new stage of the methodological developments.

In my opinion, the creation of an effective homogenization method must be based on three principles, namely: i) Consideration of the statistical properties of observational data (including its IHs), for which the homogenization method will be used; ii) Relying on the best results of earlier achievements; iii) Creating additional value with innovation and automation.

i) *Consideration of the statistical properties of observed temperature and precipitation series* - The mean frequency of detected IHs in European and North American climate records is around 5-6 per 100yr and per station (Auer *et al.*, 2005; Menne *et al.* 2009; Venema *et al.*, 2012), although it depends on network density and the examined climatic variable (Menne *et al.*, 2009) and on the homogenization method applied (Domonkos, 2011a). Note that IHs are modelled as a sudden shift in the mean and referred to as “break” throughout this study when no other specification is given. The true frequency of breaks is likely higher than their detected frequency, because small-size shifts and short-term biases often cannot be detected (Brohan *et al.*, 2006; Menne *et al.*, 2009; Domonkos, 2011a, 2013a). Therefore, the true frequency of breaks in observational time series is expected to be at least equal but likely higher than 5 breaks per station and per 100yr, whilst other kinds of IHs (e.g. trend-like biases) may additionally occur in the series. The optimal way of homogenizing datasets with such IHs is the use of multiple break methods, i.e. methods in which the joint structure of IHs are searched and corrected directly, taking into account the mutual effects of individual IHs. The development of multiple break methods began in the last decade of the 20<sup>th</sup> century (apart from some basically subjective methods not considered in here) and now we have four methods of this kind: MASH, PRODIGE, HOMER and ACMANT. Considering that at present, climatologists apply approximately 40 methods for homogenizing temperature and precipitation series (Domonkos and Efthymiadis, 2013), multiple break methods compose a

small cluster of the existing methods. It is striking that the cluster of the best performing methods in HOME tests is almost identical with the cluster of multiple break methods, showing that HOME tests justified the advantage of using multiple break methods. Consequently, the incorporation of multiple break techniques in ACMANT was a good decision.

Another positive feature of ACMANT related to this point is that the semi-sinusoid annual cycle of Tmean and Tmax biases is taken into account in the homogenization procedure. In mid- and high latitudes the annual cycle of insolation results in the annual cycle of biases caused by various IHs (Drogue *et al.*, 2005; Domonkos and Štěpánek, 2009; Brunet *et al.*, 2011), since bias-sizes are often related directly or indirectly to the duration and intensity of insolation. In ACMANT, the homogenization of Tmean and Tmax observed in mid- or high latitudes is performed with a bivariate detection where the two variables are the annual mean and the amplitude of summer – winter difference. The calculation of adjustment terms and some other routines of the computer program also consider the semi-sinusoid cycle of biases and I believe that these properties significantly contribute to the high efficiency of ACMANT.

ii. *Relying on the best results of earlier achievements* - ACMANT is based on the detection and correction method of PRODIGE (the name “ACMANT” came from “Adapted Caussinus - Mestre Algorithm for homogenising Networks of Temperature series). In testing detection parts of homogenization methods, PRODIGE showed the highest efficiency (Domonkos, 2011a, 2013a), while ANOVA is a correction method, with which even the results of other homogenization methods could be improved (Domonkos *et al.*, 2011). The adaptation of routines, which once were effective in other homogenization methods, sets a good basis for the creation of new methods that could outperform the earlier methods.

iii. *Creating additional value with innovation and automation* - Three kinds of added value will be discussed here: a) Separating time scales in the detection of IHs; b) Exploitation of the partial regularity in the seasonal changes of biases caused by IHs; c) Automation of the homogenization procedure.

a) Although the main goal of time series homogenization is not the break detection (but the minimization of the residual non-climatic biases), break detection is included in every homogenization method, since the identification of break positions helps to eliminate the biases. On annual or multi-annual timescale, the spatial correlations and hence the signal to noise ratio is higher than on monthly scale. In addition, monthly data is often affected by annual cycle of bias, which is obviously absent in annual data. On the other hand, statistical samples are larger when time series are examined on monthly resolution, thus the advantages and drawbacks could seemingly be compensated by each-other. However, one main difficulty of the homogenization task is that both the number of breaks and their positions must be estimated from the sample, and the uncertainty can be reduced if breaks are searched at the time scale in which they manifest themselves best.

b) Tmean and Tmax data are often affected by seasonally varying biases and such variations can be modelled by sinusoid cycles. ACMANT and HOMER are the only methods, which exploit this feature of temperature data in their break detection algorithms.

c) ACMANT and HOMER have been developed from PRODIGE, but while PRODIGE and HOMER are semi-objective methods, ACMANT is a fully objective and fully automatic homogenization method. Note that we use the term “objective method” in the sense that the results do not depend on homogenizers. The objectivity and automation has four advantages:

i) The use of automatic methods is advisable for large datasets and practically the only option for homogenizing datasets including 50 or more time series. ii) Testing of automatic methods is straightforward and tests are easily manageable even with huge test datasets. Owing to these tests the performance of automatic methods is more transparent than that of other methods. iii) A homogenization product of a fully objective method can be reconstructed at any time. iv) Automatic methods are easy-to-use for climatologists.

Note, however, that automatic methods are not competitors of subjective or semi-objective methods, since for the homogenization of small networks with the help of metadata, the use of subjective or semi-objective methods is preferred.

The ACMANT1 software was only for homogenizing monthly means of Tmean and Tmax observed in mid- or high latitudes. The full description of ACMANT1 is published in D2011. However, I have made new computer programs for homogenizing other variables and even the structure and content of the algorithm for homogenizing Tmean and Tmax have changed significantly since then. The two most important novelties in the general structure of ACMANT2 relative to its earlier version are as follows: i) In ACMANT2 the adjustment terms are always calculated by ANOVA, also in the phase of Pre-homogenization; ii) A new subroutine has been included, “Filtering of outlier period”, that is always applied just after the common outlier filtering. The purpose of filtering of outlier periods is to remove large, short-term biases before homogenization, similarly as individual outlier values are removed. The decrease of the residual root mean square error in homogenized series can be expected from these changes.

The parameterization of ACMANT2 is based on tests with various large test datasets similar in climatic characteristics to the HBM, but varied in the properties of IHs included in them.

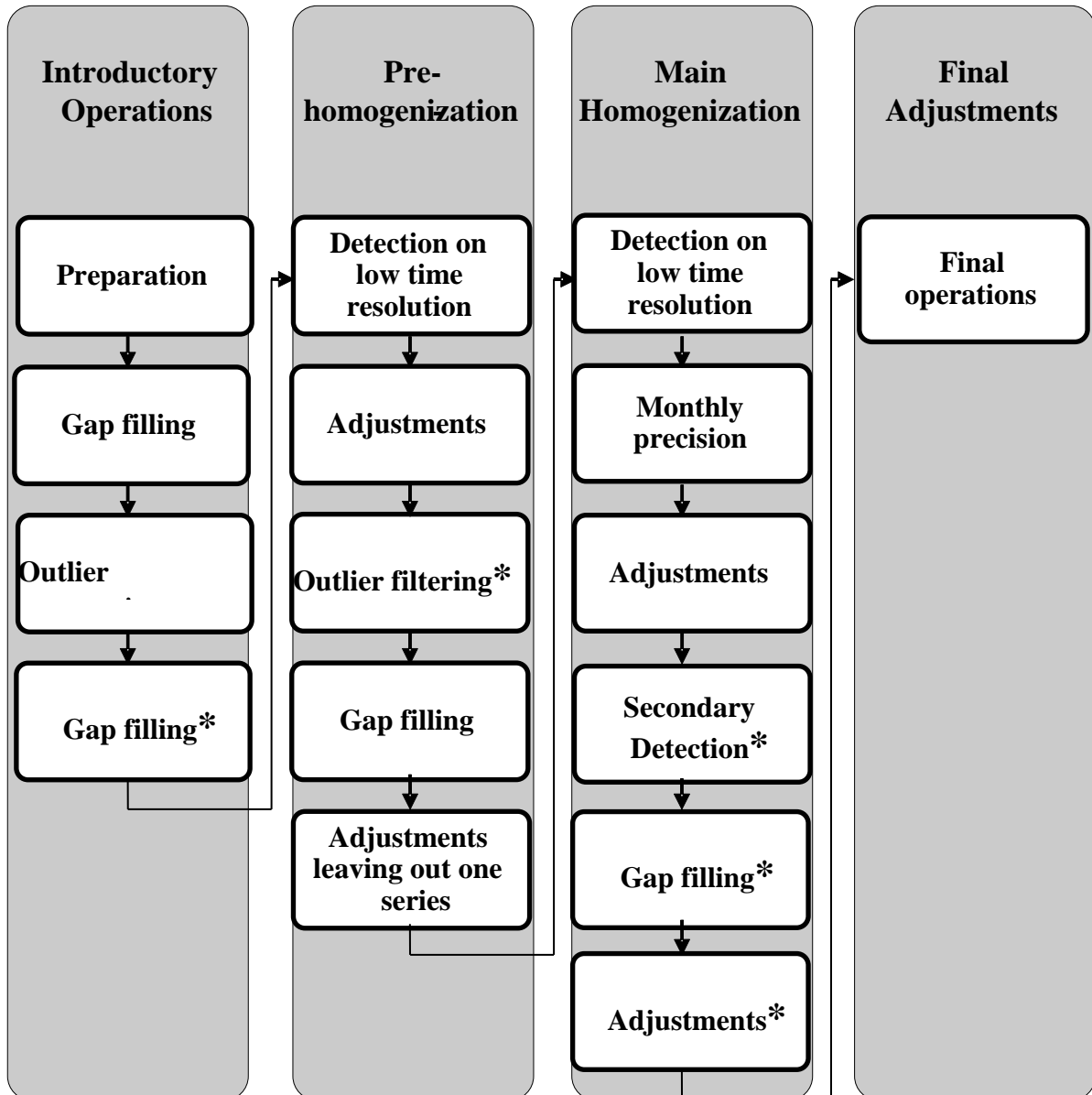
### **3. ACMANT2: STRUCTURE AND KEY THEORETICAL PROPERTIES**

ACMANT2 is composed of four main segments, namely Introductory Operations, Pre-homogenization, Main Homogenization and Final Adjustments, and each main segment includes various routines (e.g. for break detection, outlier filtering, bias correction). Some routines are common for more than one main segment (Fig. 1), since the accuracy of certain operations increases with the improving homogeneity of the data during the procedure, and thus the repeated application of such routines improves the accuracy of the final results.

Fig. 1 shows the most important segments only, i.e. operations, like calculation of spatial correlations, building reference series, exclusion of detected breaks of insignificant size, etc. are not shown. On the other hand, routines marked with asterisk indicate that they are not included in precipitation homogenization. For homogenizing PP, most routines are applied in the same way as for temperature homogenization, after the row PP values are converted by a quasi-logarithmic transformation. However, the work on monthly scale is strongly reduced in PP homogenization due to the often large spatial and temporal irregularity in monthly PP totals (particularly in semiarid regions).

### 3.1. Relative homogenization based on reference series

ACMANT2 is a relative homogenization method, which means that the detection of IHs is performed by examining the differences between a candidate series and a reference series, and then any detected IH is assigned to the candidate series. The method of creating reference series from composite series mostly follows the rules of Peterson and Easterling (1994) with some differences in the details. Composite series are weighted according to the squared spatial correlations of monthly temperature anomalies and the first difference series



**Fig. 1.** Scheme of ACMANT2. The shown segments are common for all computer programs included in ACMANT2, except that the ones marked with (\*) are included in temperature homogenization only.

(increment series) are used for calculating the correlations, in order to reduce the impact of IHs on the empirical correlations. Possible effects of IHs in the reference composites are not considered during the calculations of the spatial differences. Note, however, that in the phase of Main Homogenization the reference composites have been pre-homogenized (while the

candidate series remains the raw, outlier filtered series), and in this way the IHs of the reference composites have markedly reduced impact on the final results. The parameterization for the calculation of reference series is shown in AP I-2.1.

A speciality of ACMANT2 is that the detection of IHs on low time resolution (i.e. detection of biases for at least 3 year long sections of the candidate series) has two main phases. The goal of the first phase within the Pre-homogenization segment is to remove or reduce the large-size biases of the time series, those that will be reference composites in the second phase, within the Main Homogenization stage. The first and second phases are performed in almost the same way, with a small change only in the parameterization. In the Pre-homogenization phase, the future candidate series is excluded from the calculation of adjustment terms, and thus the multiple use of the same spatial relationship (including error term) is excluded.

The use of reference composites flexibly changes if the number of available reference composites is different for diverse sections of the candidate series (AP I-2.2), and all the reference composites of at least 0.4 spatial correlation with the candidate series are utilized, disregarding possible differences in the starting and ending dates of the series. Note that in ACMANT2, reference series are never used for calculating adjustment-terms. It is an important detail in which multiple break methods differ from more traditional homogenization methods (e.g. Standard Normal Homogeneity Test [SNHT] by Alexandersson and Moberg, 1997; RHTest by Wang, 2008, etc.).

### **3.2. Inhomogeneity detection on low time resolution**

ACMANT2 includes two markedly different types of break detection. One is for identifying long-standing biases whose characteristic time is longer than 2 years (referred as detection on low time resolution), while the other is for identifying temporarily existing, short-lived but large size biases lasting from 3 months to 24 months.

#### **3.2.1. Fitting optimal step function (univariate detection)**

Fitting optimal step function is a known technic for the detection of multiple breaks in time series (Hawkins, 1972; Caussinus and Mestre, 2004). Presuming that a time series contains  $K$  IHs and all of them are sudden shifts of the mean values (i.e. breaks), the time series is modelled by a step function of  $K + 1$  steps. The optimal step function can be found with the variance minimization of the data relative to the step function model. As true IHs are often sudden shifts (e.g. due to station relocation, instrumental change), this model is realistic. When gradually increasing bias (trend-like) IHs occur, the step function approach still provides fair (although slightly less accurate) results, transforming the trend into two or more steps.

Let the annual mean ( $E$ ) and the section mean (upper stroke) for step  $k$  of variable  $x$  be defined by (1) and (2),

$$E(x_j) = \frac{\sum_{m=1}^{12} x_{j,m}}{12} \quad (1)$$

$$\overline{\mathbf{X}}_{\mathbf{k}} = \frac{1}{j_k - j_{k-1}} \sum_{i=j_{k-1}+1}^{j_k} x_i \quad (2)$$

then the optimal step function for time series  $\mathbf{Q}$  of length  $L$ , including  $K$  breaks is given by (3) and (4).

$$\min_{[j_1, j_2, \dots, j_K]} \left\{ \sum_{k=0}^K \sum_{i=j_k+1}^{j_{k+1}} (E(q)_i - \overline{\mathbf{E}(q)}_{\mathbf{k}})^2 \right\} \quad (3)$$

$$j_0 = 0, \quad j_{K+1} = L \quad (4)$$

Note that when step function fitting is applied in ACMANT2, the minimum length of a step is at least 3 time units, i.e. 3 years or 3 months depending on the time step between two adjacent values (5).

$$j_{k+1} - j_k \geq 3 \quad \forall k \in \{0 \leq k \leq K\} \quad (5)$$

### 3.2.2. Bivariate detection of breaks with fitting optimal step function

As it has been mentioned, biases in Tmean and Tmax series often have sinusoid annual cycle linked to the annual cycle of insolation. Therefore when a break occurs for station relocation or change in the instrumentation, etc., both the annual means and the amplitude of seasonal cycle can be affected. In the bivariate detection of Tmean and Tmax homogenization, breaks with common timings are searched for the annual mean ( $E$ ) and for the amplitude of summer-winter difference ( $Z$ ).  $Z$  is defined as:

$$Z(x_j) = \frac{1}{3.5} \sum_{m=1}^{12} c_m x_m \quad (6)$$

where  $m$  denotes calendar month and the monthly coefficients ( $c_m$ ) are negative in winter and positive in summer:

$$c_1 = c_{11} = c_{12} = -1,$$

$$c_2 = -0.5,$$

$$c_3 = c_4 = c_9 = c_{10} = 0,$$

$$c_5 = c_6 = c_7 = 1,$$

$$c_8 = 0.5$$

Then the best fitting step function to series  $\mathbf{Q}$  including  $K$  breaks is given by (7).

$$\min_{[j_1, j_2, \dots, j_K]} \left\{ \sum_{k=0}^K \sum_{i=j_k+1}^{j_{k+1}} (E(q)_i - \overline{\mathbf{E}(q)}_k)^2 + c_0^2 (Z(q)_i - \overline{\mathbf{Z}(q)}_k)^2 \right\} \quad (7)$$

$c_0 = 5^{-0.5}$  (empirical constant). Note that

$$\overline{\mathbf{Q}}_k \equiv \overline{\mathbf{E}(q)}_k \quad (8)$$

by definition, the longer form in (3) and (7) is included only for showing the identity of the treatment for  $E$  and  $Z$ . Note also that for presenting the average of whole time series the index ( $k$ ) will be omitted.

### 3.2.3. Assessment of the number of breaks

The critical point of step-function fitting methods is the determination of  $K$ . In ACMANT2 a parameterized version of the Caussin - Lyazrhi criterion (Caussin and Lyazrhi, 1997) is applied. This criterion takes into account the reduction of variance due to the inclusion of breaks, but with balancing that with a penalty term due to the increasing number of steps, because the residual variance tends to decrease with the rising number of steps either the breaks between steps are significant or not. As a consequence, the rise of the number of steps will give better score only if the reduction of standard deviation overbalances the increase of the penalty term (9), (10).

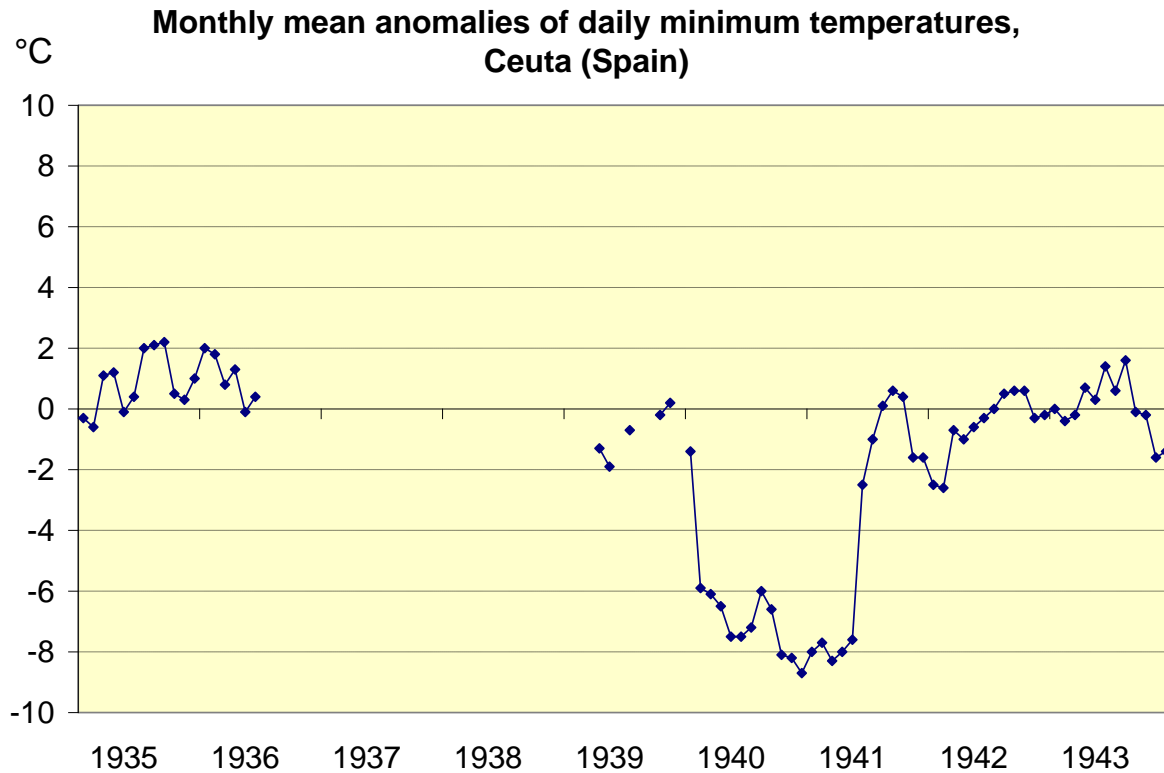
$$\ln \left\{ 1 - \frac{\sum_{k=0}^K (j_{k+1} - j_k) \cdot [(\overline{\mathbf{Q}}_k - \overline{\mathbf{Q}})^2 + c_0^2 (\overline{\mathbf{Z}(q)}_k - \overline{\mathbf{Z}(q)})^2]}{\sum_{i=1}^L (E(q)_i - \overline{\mathbf{Q}})^2 + c_0^2 (Z(q)_i - \overline{\mathbf{Z}(q)})^2} \right\} + S \quad (9)$$

$$S = p \frac{2K}{L-1} \ln(L) \quad (10)$$

The shown formula differs from the original one in one detail, i.e. the penalty term ( $S$ ) here includes an empirical coefficient ( $p$ ). The value of  $p$  is different for univariate and bivariate detection, as well as different in the Pre-homogenization than in the Main Homogenization. In the Main Homogenization  $p = 1.4$  in univariate detection and  $p = 1.0$  in bivariate detection. See more details in AP I-3.1 and AP II-1.

### 3.3. Detection of short-term biases

Short-term IHs can be modelled by a platform-shape bias from the true climatic values where the platform is the composition of a pair of breaks of the same shift-size but to the opposite directions (Fig. 2). Short-term IHs can be caused by temporal changes in the conditions of the observation. The frequency of the short-term biases can be much higher than their detected frequency, because the signal-to-noise ratio is relatively low for short sections of time series due to the limited sample size (Domonkos, 2013a). Experimental results indicate that the true frequency is really significantly higher than the detected frequency (Domonkos, 2011a; Rienzner and Gandolfi, 2011).



**Fig. 2. Large, short-term, platform-shaped bias between 02-1940 and 05-1941 in the Tmin of Ceuta (Spain).**



In ACMANT2, large-size, platform-shaped biases are detected in moving windows of monthly temperature anomalies. This new routine of ACMANT (such step is not included in ACMANT1) is named Filtering of outlier periods (AP I-4.1). This procedure has characteristics similar to filtering out outlier values, as well as to the minimization of standard deviation relative to platform-shape step functions. It is always performed just after the common Outlier filtering.

The detection of long-term biases might be affected by the existence of large-size short-term biases and vice versa. Therefore the detection and elimination of short-term biases is performed three times in ACMANT2, approaching step-by-step to the final solution. These operations are applied first in the Introductory Operations, then within the Pre-homogenization phase after the adjustment of long-term biases, and finally in the Main Homogenisation phase, after the adjustment of long-term biases. The way of the detection and correction of short term biases in the Main Homogenization differ from the Filtering of outlier periods, i.e. the routine Secondary Detection of ACMANT1 (D2011) has been kept with some little changes in the parameterization only (AP I-4.2 and AP II-4).

Note that in PP homogenization neither outlier filtering nor any kind of operation for filtering out short-term biases is performed.

### 3.4. Data adjustment

In ACMANT2, adjustment-terms are generally calculated with variance analysis (ANOVA, 3.4.1 – 3.4.2.). However, ANOVA determines temporal differences only, therefore one fix value (time series average, or a reference value of the time series which is considered to be unbiased) must be defined (3.4.3). For periods of very short-term biases, interpolation technique is applied instead of ANOVA (3.4.4).

#### 3.4.1. The ANOVA model for the assessment of adjustment terms

The ANOVA procedure determines the minimum variance of anomalies relative to the climate signal of an examined region, relying on the timings of detected breaks in all the examined time series of the region. It is proven that ANOVA provides the optimum estimation of adjustment-terms when the spatial gradients of climate are temporally constant and the list of detected breaks is correct (Mestre, 2004; Caussinus and Mestre, 2004). Moreover, experiments showed that ANOVA performs better than conservative correction methods even when the list of detected breaks is partially correct only (Domonkos *et al.*, 2011).

For applying ANOVA in time series homogenization, a model is set up. In this model, the observed values are considered to be the sum of the climate signal ( $u$ ), station effect ( $v$ ) and noise ( $\varepsilon$ ) for each time series and each time point (11).

$$\mathbf{X} = \mathbf{U} + \mathbf{V} + \boldsymbol{\varepsilon} \quad (11)$$

The spatial gradients of climate are temporally constant, which is approximately true for observational datasets when data for a specific climatic zone is examined. No other constraint is included for climate. Station effect means the sum of site effect (i.e. temporally constant difference relative to the climate signal) and the biases caused by IHs. In the model, all IHs are breaks, and their timings are known. This variance minimization can be solved with the construction of an equation system following the relationships in the model.

ANOVA searches the solution of (11) with the minimum variance of  $\varepsilon$  for the entire dataset. The minimum variance can be obtained by (12, 13).

$$Nu'_i + \sum_{s=1}^N v'_s = \sum_{s=1}^N x_{s,i} \quad \text{for every } i: i \in [j_{\min}, j_{\max}] \quad (12)$$

$$\sum_{i=j_k+1}^{j_{k+1}} u'_i + (j_{k+1} - j_k)v'_{s,k} = \sum_{i=j_k+1}^{j_{k+1}} x_{s,k} \quad \text{for every } s \text{ and } k \quad (13)$$

In (12) and (13)  $N$  denotes the number of station series,  $s$  the serial number of station series,  $j_{\min}$  and  $j_{\max}$  stand for the first and last years of the period, respectively, for which homogenisation is performed, while apostrophe denotes estimated variable.

When the period between adjacent breaks is very short or when simultaneous breaks occur in various station series, the efficiency of ANOVA is reduced. Moreover, if all time series have a detected break at the same time, then the equation system is undetermined. For these reasons, breaks of relatively small estimated size are sometimes deleted from the break list (AP I-5).

### 3.4.2. The use of ANOVA in ACMANT2

ANOVA is always applied using data without bias corrections, because the recursive application of ANOVA could multiply the errors of the estimated spatial relationships.

ANOVA can be applied separately for the variables under examinations, thus in the homogenization of Tmean and Tmax, ANOVA is applied separately for annual means ( $E(x)$ ) and summer-winter differences ( $Z(x)$ ), then the monthly adjustment-terms are derived from them. In the Pre-homogenisation, ANOVA is applied on data of annual resolution, while in the Main Homogenization the input data is monthly. In monthly resolution,  $\mathbf{X}$  is examined directly, instead of  $\mathbf{E}(x)$ . However,  $\mathbf{Z}(x)$  is a variable whose interpretation on monthly scale is not straightforward. Monthly values of  $\mathbf{Z}(x)$  are defined for each month ( $h$ ) of the series by using the data of the 12-month symmetric window around  $h$  (14).

$$Z(x)_{j,h} = \sum_{h'=h-5}^{h+5} c_{m(h')} \cdot x_{h'} + 0.5(c_{m(h-6)}x_{h-6} + c_{m(h+6)})x_{h+6} \quad (14)$$

(Note: close to the endpoints of the series the extent of window is limited by the data availability.) Coefficients  $c_m$  are the same as in (6).

If break  $k$  has the timing  $H(k)$  in monthly scale and  $\alpha_k$  and  $\beta_k$  denote the estimated station effects for the homogeneous section of  $[H(k-1)+1, H(k)]$  for  $\mathbf{X}$  and  $\mathbf{Z}(x)$ , respectively, then the estimated station effect ( $v'$ ) is given by (15) for each month of the section.

$$v'_{i,m} = \alpha_k + c^*_m \beta_k \quad (15)$$

Coefficients  $c^*_m$  differ from  $c_m$  in a way that they provide the same summer – winter difference as  $c_m$ , but with a harmonic annual cycle of the coefficients.

When only one variable is examined as in the homogenization of Tmin, the determination of adjustment term is simplified to (16).

$$v'_{i,m} = \alpha_k \quad (16)$$

(16) shows that in the Tmin homogenization of ACMANT2 the monthly adjustment terms are independent from the season of the year.

### 3.4.3. Selection of reference period

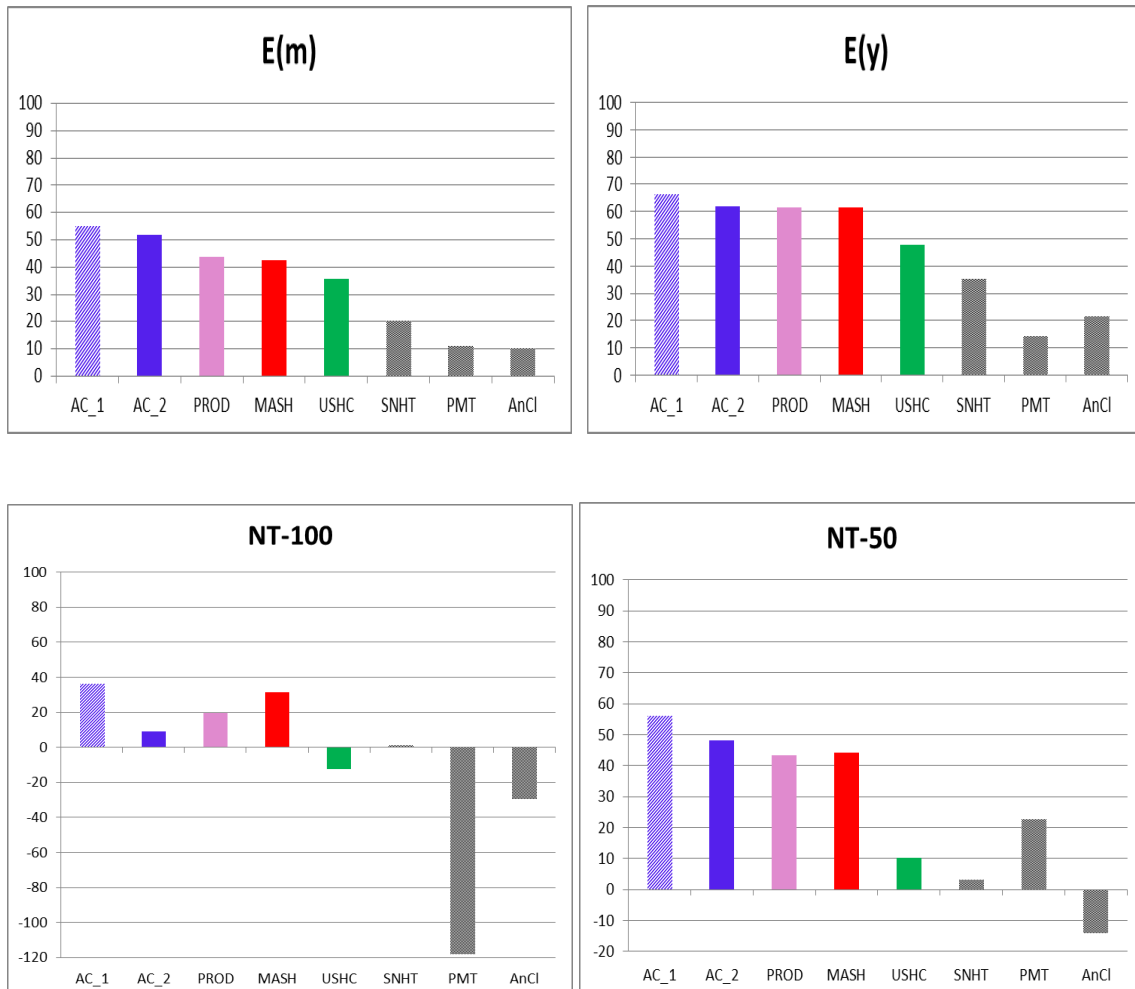
In ACMANT2, the values of the last homogeneous section of the series are considered to be unbiased, and it is considered to be reference period in the adjustment of the other sections of the series. This assumption is rather general in time series homogenization, but it might have unfavourable consequences when the last homogeneous section is too short for acquiring the accurate estimates of its statistical properties or when it has characteristics atypical for the site due to instrument error or for any other reasons. If the statistical properties of the last homogeneous section do not reflect well the true climate, its use as reference period may cause biases in the mean climatic characteristics of the site, as well as in the spatial climatic gradients. Notwithstanding, this problem has no effect on the reliability of the temporal variability. Note that the possible application of adjustment terms varying according to the probability distribution function value (percentile) of the raw data (as for instance in Della-Marta and Wanner, 2006) would make the homogenization results more sensible to the choice of the reference period, since any bias of the empirical probability distribution from the true climate in the reference period would be exported to all the other sections of the time series. Therefore as long as ACMANT remains fully automatic, percentile dependent adjustments will not be included in it.

### 3.4.4. Adjustment of short-term biases

In the Introductory Operations and Pre-homogenization, the values within the period of detected short-term biases are always substituted with interpolated values. By contrast, in the Main Homogenization, biases shorter than 6 months are adjusted by interpolation only, while the longer ones are corrected by ANOVA.

## 4. EFFICIENCY OF ACMANT2 IN HOMOGENIZING THE HOME BENCHMARK

If we would like to compare the efficiencies of different homogenization methods, it is still the HBM is the best for this purpose, since the characteristics of the HBM are rather close to the characteristics of observational time series and the efficiencies with this dataset are known for several homogenization methods (Venema et al., 2012). As the biases of the HBM series have quasi sinusoid annual cycle, the program for homogenizing Tmean and Tmax can be tested with this dataset.



**Fig. 3. Efficiency (%) in reducing RMSE of HBM with various homogenization methods. E(m) – RMSE of monthly values, E(Y) – RMSE of annual values, NT-100 – RMSE of network mean trends for the whole period (100 years) examined, NT-50 – RMSE of network mean trends over the last 50 years, AC\_1 – ACMANT1, AC\_2 – ACMANT2, PROD – PRODIGE, PMT – penalized maximum t-test of RHTest, AnCl – AnClim (Štěpánek et al., 2009).**

Fig. 3 shows the efficiency of ACMANT2 in comparison with the efficiencies of several other methods. It can be seen that the efficiency of ACMANT2 is slightly lower than that of ACMANT1. The slight decrease may have 3 reasons:

- i) The parameterization of ACMANT1 can be overfitted to HBM, since in the development of ACMANT1 I used the HBM.
- ii) In the HBM the number of short-term biases is unrealistically low, and thus the positive effect of the inclusion of filtering of outlier periods in ACMANT2 does not appear in the tests with this dataset. Note here that short-term biases are not inserted to HBM, thus short-term biases in HBM are present only in the rare cases of their accidental formation from randomly placed breaks.
- iii) Random fluctuation of the results due to the small sample size. Note here that the sample size of HBM (15 networks) is obviously very small for the assessment of the efficiency in reducing network mean trend errors.

Fig. 3. shows that even with the few percentages drop relative to its earlier version, ACMANT is one of the most effective homogenization methods, and considering the fully automatic methods tested by HOME, still ACMANT shows the highest efficiency. Note however, that several other homogenization methods have also been developed since the HOME tests, thus new comparative tests with a large new benchmark of realistic time series properties are needed to see more clearly the rank order of the efficiencies.

## **5. USE OF ACMANT2 SOFTWARE**

### **5.1. Some notes on the use of the software**

The software package has a manual available in web ([www.c3.urv.cat/data.html](http://www.c3.urv.cat/data.html)) together with the software. Therefore only some important points of the use are described here.

The software package includes 6 different computer programs which can be chosen according to the characteristics of the input raw data. Three programs are for homogenizing daily data, while the other three programs treat monthly data only. The three programs differ according to the variable treated: One program is for the homogenization of Tmean or Tmax, another one is for Tmin and the third one is for PP. The software package contains also some auxiliary files, their function and use is described in the Manual.

Anything is the input variable, some rules are common for the homogenization with ACMANT2. At least 4 time series with adequate spatial correlations (AP-I-2.1) are needed. The length of the input series may vary between 10 years and 200 years. The rules of input data preparation shown in the Manual must be followed accurately, otherwise the selected program will not run or will stop with some error message.

## 5.2. Selection of the appropriate program

The selection of the appropriate program seems to be straightforward, since the kind of the input variable (i.e.  $T_{min}$ ,  $T_{max}$ ,  $T_{mean}$  or PP) and its time resolution (daily or monthly) determine which program matches best. Yet there is a gap in this simple matching between variables and programs: The program including the bivariate detection for annual mean and summer – winter difference is proposed to use for  $T_{mean}$  and  $T_{max}$  from the mid- or high latitudes only. As quasi sinusoid annual cycle of biases is not expected in temperature data of the tropical belt and in monsoon regions, the program with bivariate detection is not recommended to use there. The recommended matching between variables and programs is shown in Table 1.

Table 1. Recommended matching between input data type and programs of ACMANT2 software

Input variable and region	Program		
	$T_{mean}\&T_{max}$	$T_{min}$	PP
$T_{mean}$ or $T_{max}$ in mid or high latitudes	X		
$T_{mean}$ or $T_{max}$ in tropical or monsoonal regions		X	
$T_{min}$ anywhere		X	
PP anywhere			X

## 5.3. Options offered for users

Although ACMANT2 is fully automatic, there some options beyond the choice of the appropriate program are offered for the users at the initiation of the homogenization procedure. Four kinds of options are asked from the users: i) in temperature homogenization: the programs can be run with or without outlier filtering; ii) in precipitation homogenization: the program can be run with dividing the year into snowy and rainy seasons or without such division; iii) in homogenization of daily data: inputting both daily and monthly data or monthly raw data is constructed from daily data by the program; iv) output data format.

i) The programs of temperature homogenization can be run with or without outlier filtering. The proposed mode is the inclusion of outlier filtering, but there is one exception: Users may check manually the detected outliers by ACMANT, and this check might find that some detected outliers should be considered true extreme values instead of outliers. In such a case the proposed continuation of the homogenization procedure is as follows: a) Justified outliers must be shown with missing data code in the raw data; b) Accepted extreme values must be left unchanged in the raw data; c) Repeating the run of ACMANT2 with the use of the modified input data and without outlier filtering.

ii) If in a part of the year the dominant form of the precipitation is snow, the starting and ending months of the snowy season must be introduced at the initialization of the PP homogenization. It is because the IHs of snow data often markedly differ from the IHs of rain,

due to the technical problems of catching snow precipitation and converting it into water amount comparable with rain precipitation. To manage this problem, PP homogenization in ACMANT2 includes two different modes, one mode is similar to the univariate homogenization of Tmin (when there is no snowy season in the region), while the other mode is bivariate homogenization in which the two variables are the total amount of the PP for the rainy season and that is for the snowy season. The program selects the appropriate mode automatically, after the requested parameters are introduced by the user.

iii) The detection of IHS and the calculation of adjustment terms are always based on monthly data. If the input is daily data, ACMANT2 develops the monthly dataset, does the homogenization and finally adjusts both the monthly and daily data. As input data may include missing values or outliers, the characteristics of the developed monthly dataset might depend on the treatment of these quality problems of the initial dataset. The development of monthly dataset in ACMANT2 is automatic, but there is an option for the users to introduce their own developed monthly dataset (together with daily data). This option is recommended to use in the case when the user has a monthly dataset which has been developed with the meticulous check of the possible quality problems in the daily data.

iv) The program offers that the output will consist of a default output package, but the user may select his choices if he wants. The goal of leaving free options in the form of the output package is to provide the opportunities of a) having homogenized dataset that is immediately applicable as input in extreme index calculation softwares, b) having the dataset with or without infilling the missing values, c) providing supplementary information about the spatial correlations and other characteristics related to the homogenization procedure.

## **6. CONCLUSIONS**

ACMANT2 homogenization software has recently been developed. This software has been prepared for the automatic homogenization of observational temperature and precipitation datasets. This paper together with D2011, provides the whole description of the temperature homogenization with ACMANT2. The high efficiency of temperature homogenization with ACMANT2 is illustrated with the homogenization of the HOME benchmark dataset. Due to its high efficiency, author recommends the use of the software in each case when the size of the dataset and the spatial correlations allow the use of automatic homogenization method. The use of ACMANT2 is particularly recommended for very large datasets when it is hard to use non-automatic methods.

## **Acknowledgements**

The research was supported by the European project UERRA FP7-SPACE-2013-1.

## References

- Alexandersson H, Moberg A. 1997. Homogenization of Swedish temperature data.1, Homogeneity test for linear trends. *Int. J. Climatol.*, **17**: 25–34.
- Auer I, Böhm R, Jurkovic A, Orlik A, Potzmann R, Schöner W, Ungersböck M, Brunetti M, Nanni T, Maugeri M, Briffa K, Jones PD, Efthymiadis D, Mestre O, Moisselin J-M, Begert M, Brazdil R, Bochnicek O, Cegnar T, Gajic-Capka M, Zaninovic K, Majstorovic Z, Szalai S, Szentimrey T, Mercalli L. 2005. A new instrumental precipitation dataset for the greater Alpine region for the period 1800–2002. *Int. J. Climatol.* **25**: 139–166.
- Brohan P, Kennedy JJ, Harris I, Tett SFB, Jones PD. 2006. Uncertainty estimates in regional and global observed temperature changes: a new data set from 1850. *Journal of Geophysical Research* **111**: D12106, doi: 10.1029/2005JD006548.
- Brunet M, Asin J, Sigró J, Bañón M, García F, Aguilar E, Palenzuela JE, Peterson TC, Jones P. 2011. The minimization of the screen bias from ancient Western Mediterranean air temperature records: an exploratory statistical analysis, *Int. J. Climatol.* **31**: 1879–1895.
- Caussinus H, Lyazrhi F. 1997. Choosing a linear model with a random number of change-points and outliers. *Ann. Inst. Statist. Math.* **49**(4): 761-775.
- Caussinus H, Mestre O. 2004. Detection and correction of artificial shifts in climate series. *J. Roy. Stat. Soc. C* **53**: 405-425.
- Craddock JM. 1979. Methods of comparing annual rainfall records for climatic purposes. *Weather*, **34**: 332-346.
- Della-Marta PM, Wanner H. 2006. A method of homogenizing the extremes and mean of daily temperature measurements. *J. Clim.* **19**: 4179–4197.
- Domonkos P. 2008: Testing of homogenisation methods: purposes, tools and problems of implementation. *Proceedings of the 5<sup>th</sup> Seminar for Homogenisation and Quality Control in Climatological Databases*. (Ed. Lakatos, M., Szentimrey, T., Bihari, Z. and Szalai, S.), WCDMP-No. 71, WMO/TD-NO. 1493, 126-145.
- Domonkos P. 2011a. Efficiency evaluation for detecting inhomogeneities by objective homogenisation methods. *Theor. Appl. Climatol.* **105**: 455-467.
- Domonkos P. 2011b. Adapted Caussinus-Mestre Algorithm for Networks of Temperature series (ACMANT). *Int. J. Geosci.* **2**: 293-309, doi: 10.4236/ijg.2011.23032.
- Domonkos P. 2013a. Efficiencies of inhomogeneity-detection algorithms: comparison of different detection methods and efficiency measures. *Journal of Climatology* pp15, doi:10.1155/2013/390945.
- Domonkos, P. 2013b: Measuring performances of homogenization methods. *Időjárás* 117, 91-112.
- Domonkos P, Efthymiadis D. 2013. After HOME: Progress in the practical application of homogenization methods. 11<sup>th</sup> ECAM conference, 09-13 Sept 2013, Reading, England, EMS2013-99.
- Domonkos P, Štěpánek P. 2009. Statistical characteristics of detectable inhomogeneities in observed meteorological time series. *Studia Geoph. Geod.* **53**: 239-260, doi: 10.007/ s11200-009-0015-9.
- Domonkos P, Venema V, Mestre O. 2011. Efficiencies of homogenisation methods: our present knowledge and its limitation. In: *Seventh Seminar for Homogenisation and Quality Control in Climatological Databases* (Eds. Lakatos M, Szentimrey T, Vincze E.), WCDMP-78, WMO, Geneva, 11-24.
- Drogue G, Mestre O, Hoffmann L, Iffly J-F, Pfister L. 2005. Recent warming in a small region with semi-oceanic climate, 1949-1998: what is the ground truth? *Theor. Appl. Climatol.*, **81**: 1-10.
- Hawkins DM. 1972. On the choice of segments in piecewise approximation. *J. Inst. Math. Appl.* **9**: 250–256, doi:10.1093/imamat/9.2.250
- Menne M, Williams Jr CN. 2009. Homogenization of temperature series via pairwise comparisons. *J. Clim.* **22**: 1700-1717.
- Menne MJ, Williams Jr CN, Vose RS. 2009. The U.S. Historical Climatology Network monthly temperature data, version 2. *Bull. Amer. Meteor. Soc.*, **90**: 993-1007.



- Mestre O. 2004: Correcting climate series using ANOVA technique. In: Fourth Seminar for Homogenization of Surface Climatological Data (Eds. Szalai S, Szentimrey T) WCDMP 56, WMO-TD 1236, WMO, Geneva, 93-96.
- Mestre O, Domonkos P, Picard F, Auer I, Robin S, Lebarbier E, Böhm R, Aguilar E, Guijarro J, Vertacnik G, Klancar M, Dubuisson B, Štěpánek P. 2013. HOMER: homogenization software in R – methods and applications, *Időjárás* **117**: 47–67.
- Peterson TC, Easterling DR. 1994. Creation of homogeneous composite climatological reference series. *Int. J. Climatol.* **14**: 671-679.
- Peterson TC, Easterling DR, Karl TR, Groisman P, Nicholls N, Plummer N, Torok S, Auer I, Böhm R, Gullett D, Vincent L, Heino R, Tuomenvirta H, Mestre O, Szentimrey T, Salinger J, Førland EJ, Hanssen-Bauer I, Alexandersson H, Jones P, Parker D. 1998. Homogeneity adjustments of in situ atmospheric climate data: a review. *Int. J. Climatol.* **18**: 1493–1517.
- Picard F, Lebarbier E, Hoebeke M, Rigail G, Thiam B, Robin S. 2011. Joint segmentation, calling and normalization of multiple CGH profiles. *Biostatistics*, **12**: 413-428. doi:10.1093/biostatistics/kxq076
- Rienzner M, Gandolfi C. 2011. A composite statistical method for the detection of multiple undocumented abrupt changes in the mean value within a time series. *Int. J. Climatol.* **31**: 742-755.
- Štěpánek P, Zahradníček P, Skalák P. 2009. Data quality control and homogenization of the air temperature and precipitation series in the Czech Republic in the period 1961–2007. *Adv. Sci. Res.*, **3**: 23–26.
- Szentimrey T. 1999: Multiple Analysis of Series for Homogenization (MASH). In: Second Seminar for Homogenization of Surface Climatological Data (Eds. Szalai S, Szentimrey T, Szinell Cs.) WCDMP 41, WMO-TD 962, WMO, Geneva, 27-46.
- Venema V, Mestre O, Aguilar E, Auer I, Guijarro JA, Domonkos P, Vertacnik G, Szentimrey T, Štěpánek P, Zahradníček P, Viarre J, Müller-Westermeier G, Lakatos M, Williams CN, Menne M, Lindau R, Rasol D, Rustemeier E, Kolokythas K, Marinova T, Andresen L, Acquafredda F, Fratianni S, Cheval S, Klancar M, Brunetti M, Gruber C, Duran MP, Likso T, Esteban P, Brandsma T. 2012. Benchmarking monthly homogenization algorithms. *Climate of the Past* **8**: 89-115, doi:10.5194/cp-8-89-2012.
- Wang XLL. 2008. Accounting for autocorrelation in detecting mean shifts in climate data series using the penalized maximal t or F test. *J. Appl. Meteor. Climatol.* **47**: 2423–2444.
- [www.c3.urv.cat/data.html](http://www.c3.urv.cat/data.html): Software package and Manual of ACMANT2.

## **APPENDIX I. CHANGES FROM ACMANT1 IN HOMOGENIZING MONTHLY MEANS OF DAILY TMEAN OR TMAX.**

### **I-1. Replacing missing data and outliers with interpolated values**

In ACMANT2 this routine is applied three times, while in ACMANT1 it was run only twice. There is no change of the parameterization for the first two running of the routine. In the third running, the parameterization is the same as for the second routine.

### **I-2. Construction of relative time series**

Let the candidate series and reference series be denoted by **A** and **F**, respectively, then relative time series (**Q**) are defined as their difference series:  $Q = A - F$ .

#### **I-2.1. Parameterization**

Minimum length of **A**: 10 years

Minimum length of **F**: 10 years (Note: it implies that each **F** must have at least 10 year common section with **A**).

Minimum number of reference composites: 4

Maximum number of reference composites: not limited

Minimum number of months for which observed values are available both in **A** and in a reference composite: 50

Minimum threshold of spatial correlation ( $r_0$ ): it depends on the number of reference composites ( $J$ ):

if  $J = 4$  then  $r_0 = 0.6$

if  $J = 5$  then  $r_0 = 0.48$

if  $J > 5$  then  $r_0 = 0.4$

Further restrictions: If the sum of the accumulated weight ( $w$ ) of reference composite  $j$  reaches 4.0 without reference composites with length ( $n$ ) of shorter than 20 years ( $n_j < 20$ ), then the composites of  $n_j < 20$  are excluded. Similarly, if  $w \geq 7.0$ , then composites of  $n_j < 30$  are excluded.

#### **I-2.2. Constructing different relative time series for different sections of the candidate series**

(a) First the longest section of **A** that can be homogenized is paired with at least one **F** series in a way that the result **Q** series are the possible longest.

(b) In the second step, for each section of **A**, the **Q** series with the maximal accumulated weight of reference composites is created.

(c) If after (a) and (b) the number of **Q** series for a specific **A** would exceed 80, then phase (b) is repeated in a modified way, namely sections with the highest overlap are not distinguished, thus the number of **Q** series is reduced. It is managed in a way that the time lapse between similar sections is monitored and a parameter indicating the threshold degree of time lapse between sections is increasing gradually from 0 as long as the number of created **Q** series remains under 80.

(d) If the accumulated weight of the reference composites for two **Q** series differs with at least 5%, then both series are retained, irrespectively to the degree of time lapse between sections. This rule overwrite rule (c), but the total number of retained **Q** series cannot be higher than 80.

### **I-2.3. Unified relative time series**

In ACMANT2 it is not used.

### **I-2.4. Selection of relative time series**

There are often more than one relative time series are available for the examination of a specific section of the candidate series. As a principal rule, the **Q** with the highest accumulated weight ( $w$ ) of the reference composites is used only, but  $w$  is modified ( $w^*$ ) according to the length of **Q** (17), since the use of relatively long **Q** series is preferred.

$$w^* = w \log(6n_Q) \quad (17)$$

Only in the Secondary Detection, all the **Q** series are used in the check of the maximal accumulated anomalies. The rules of the harmonization of section-examinations in the detection on low time resolution are unchanged (D2011).

## **I-3. Detection on low time resolution**

### **I-3.1. Pre-homogenization**

(i) Ranking of time series according to the inhomogeneous character is not applied in ACMANT2.

(ii) Pre-detection is done according to the rules of the detection on low time resolution (sect. 3.2). In the penalty-term of the Caussinus – Lyazrhi criterion an empirical coefficient ( $p$ ) is applied (eq. 10). In the Pre-detection, this coefficient depends on the accumulated weight of the reference components (18).

$$p = p_1 - \left( \frac{p_2 \log(10w - 7)}{\log(10w)} \right)^{p_3} \quad (18)$$

$$p_1 = 2.0; \quad p_2 = 1.05; \quad p_3 = 17.0$$

(iii) Parameter  $c_0$  (vs. Eq. 7):

In ACMANT2  $c_0 = 5^{-0.5}$  both in the Pre-homogenization and in the Main Detection.

### I-3.2. Monthly precision of breaks

The width of the window, in which break is searched in the monthly precision (step III/4 in the algorithm of D2011) is 29 months (it was 25 months in ACMANT1).

## I-4. Inhomogeneity detection on short time-scale

### I-4.1. Filtering of outlier periods

Outlier periods could also be referred to as short-term inhomogeneities, since their model is a short-term, platform-like bias from the correct values. In this model the bias is constant for the outlier period. Both the detection and the adjustment of outlier periods are more similar to outlier filtering than to the detection and adjustment of long-term biases.

Filtering of outlier periods is applied for 2 - 27 month long periods, always after the routine of common outlier filtering. In switched off outlier filtering mode, the minimum duration of outlier periods is 5 months.

In searching outlier-periods, relative time series (**Q**) are used on monthly scale and the values are transformed to standard anomalies (**B**). Further denotations:  $l$  – length of outlier period,  $h_1$  and  $h_2$  – starting and ending months (respectively) of the outlier-period in the first estimation,  $l_A$  and  $l_B$  are lengths of outer sections of  $(h_1, h_2)$  before that and after that, respectively,  $\text{int}$  and  $\text{sgn}$  – integer part and sign of arithmetic expression, respectively,  $\text{mod}$  – function of modulo,  $\lambda$  statistic of significance.

The detection of outlier periods is a step-by-step procedure, since only one outlier-period is identified in a particular step, i.e. the one with the highest  $\lambda$  (19). The mean value of a potential outlier-period is compared with the mean value of the adjacent outer sections in both sides of the potential outlier-period (20). Once an outlier-period has been selected, its values are adjusted to make it possible searching the next most significant outlier-period. Here, temporal adjustments are applied, which are valid during the operations of this routine only. The temporal adjustments eliminate the difference between the means of the outlier-period and its outer sections, and thus the next most significant outlier period can be selected in the next round. The procedure stops when  $\lambda < 30$ .

The identification of an outlier-period comprises two phases. In the first phase (i), the most significant outlier-period of the time series is selected and a first estimate is made for its position. In the second phase (ii) the starting and ending months of the outlier-period are determined.

Phase (i): the outlier-period with the maximal  $\lambda$  is searched for each  $h_1, h_2$  pairs ( $2 \leq h_2 - h_1 < 27$ ) of standardised relative time series.

$$\lambda = l'^{0.75} d^2 \quad (19)$$

where  $d$  (magnitude-characteristic) and  $l'$  (duration-characteristic) are determined by Eqs. (20) and (21):

$$d = \overline{b_{h_1-l_A}, b_{h_1-1}} + \overline{b_{h_2+1}, b_{h_2+l_B}} - \overline{b_{h_1}, b_{h_2}} \quad (20)$$

$$l' = \text{int}(\max \left\{ l - \frac{0.75}{3.5} \sum_{[h_1, h_2]} c_m, 1 \right\}) \quad (21)$$

Further conditions are that

$$\text{sgn}(\overline{b_{h_1}, b_{h_2}} - \overline{b_{h_1-l_A}, b_{h_1-1}}) = \text{sgn}(\overline{b_{h_1}, b_{h_2}} - \overline{b_{h_2+1}, b_{h_2+l_B}}) \quad (22)$$

$$\text{mod}(l_A, 12) = 0 \quad \text{mod}(l_B, 12) = 0 \quad (23)$$

The usual length of the outer periods is 24 months in both sides of the potential outlier-period. However, if an outlier-period is close to an endpoint of  $\mathbf{B}$ ,  $l_B$  or  $l_A$  can be 12 or even 0. The two outer periods together must contain at least 36 months for providing statistical sample of adequate size for the calculations. For the fulfilment of this condition, if  $l_B = 0$  then  $l_A = 36$  and if  $l_A = 0$  then  $l_B = 36$ . In (21), the sum of  $c_m$  within the outlier period is included in order to take into account the seasonal imbalance of the period. It is necessary, because biases due to breaks seasonally vary, and thus a long-standing bias with enhanced seasonal cycle could be detected as short-term outlier-period when a seasonal peak of the long-term bias and random noise accidentally add up. The sum of  $c_m$  is an indicator of the seasonal imbalance and it is normalised with the absolute value of sum of  $c_m$  over a half year (i.e. 3.5, see the denominator of the coefficient). The 0.75 in the counter is an empirical constant.

Phase (ii): The first and last months of the outlier-period are re-estimated with fitting optimal step-function in window  $[b_{h_1-l_B}, b_{h_2+l_A}]$ . For longer than 9 month sections harmonic functions

are fitted instead of constant values and from this point of view the procedure is the same as the break detection part of Secondary Detection (D2011). Differing from Secondary Detection, solutions with exactly two breaks are accepted only, and the first and second breaks are expected in the periods  $[h_1 - 14, h_1 - 1]$  and  $[h_2, h_2 + 13]$ , respectively. So that, the final duration of an outlier-period is equal or greater than the pre-estimated duration. If  $h_1$  or  $h_2$  coincides with one endpoint of **B**, then one only break is searched, since the other endpoint of the outlier-period is defined by the endpoint of **B**.

#### **I-4.2. Secondary detection**

The parameterized penalty term of the Caussinus – Lyazrhi criterion is applied (eqs. 9-10), and here  $p = 1.8$ . There is no other change relative to ACMANT1.

#### **I-5. Exclusion of detected breaks**

##### **I-5.1 Causes of possible exclusions**

Detected breaks may be deleted for the following reasons:

i) If all the time series have break at the same time, the equation system of ANOVA is non-determined.

ii) A large number of simultaneous breaks reduce the reliability of the results, since the basic theory of the statistical homogenization is that the break is individual and the other station series are free from break at the same time. If the number of simultaneous breaks approaches to or exceeds the half of the number of time series within network, then there is a high risk that the correct series will be adjusted instead of the biased series.

iii) The inclusion of very close breaks (in time) or breaks with insignificant shift-size reduces the accuracy, since the random error of  $v'$  in eq. (13) increases with the shortening of the homogeneous period.

iv) When a break is detected with bivariate detection, it is possible that the shift-size is significant only in one of the variables examined.

v) The calculated shift sizes by ANOVA may indicate that a detected break is not significant statistically, in spite of it seemed to be significant during the detection phase.

In accordance with i), ii) and iii), the number of simultaneous breaks is limited in ACMANT2 even when all breaks seem to be significant, and, on the other hand, breaks with non-significant shift sizes are always deleted. The technical solution of the exclusion of breaks is as follows.

##### **I-5.2. Significance of breaks**

For the possible exclusion of one or more breaks, a significance order must be determined. For breaks which are detected in different phases (by different routines) of ACMANT2, the phase determines the order of significance, i.e. breaks detected in later phases of the procedure are considered to be more significant than those that detected in earlier phases, independently from other characteristics of the breaks. As a consequence, breaks of Secondary Detection are more significant than the breaks of Main Detection, the results of Main Detection overwrite the results of Pre-detection and the breaks of Main Detection are more significant than the breaks detected by Filtering of outlier periods.

Sometimes the rank order of significance must be determined for breaks detected by the same routine. As the assessment of significance of break  $k$  is limited to the examination of the period between  $j_{k-1}$  and  $j_{k+1}$  (which includes 1 detected break), the single break model can be applied for these assessments and thus the use of  $t$ -test and its modified versions are appropriate here.  $t$ -test is applied with the simplification that the standard deviation ( $\sigma$ ) is considered to be constant for a given series, it is because the signal-to-noise ratio is generally too low to estimate specific  $\sigma$  values for individual homogeneous sections with sufficient confidence. With this simplification, the calculation of  $t$ -statistic ( $\tau$ ) is given by (24).

$$\tau = \frac{|d|\sqrt{l_1 l_2 (l-2)}}{l\sigma} \quad (24)$$

In (24),  $l_1 = j_k - j_{k-1}$ ,  $l_2 = j_{k+1} - j_k$ ,  $l = l_1 + l_2$  and  $d$  denotes the shift size. For determining the order of significances only, instead of absolute significances,  $\tau^*$  (25) can be used instead of  $\tau$ .

$$\tau^* = \frac{l_1 l_2 d^2}{l} \quad (25)$$

### I-5.3. Reduction of the number of synchronous breaks

In ACMANT2, the number of synchronous breaks is not allowed to reach the 50% of the number of time series which are homogenized together at the section including the synchronous break. This rule is valid both for Pre-Homogenization and Main Homogenization. In Pre-homogenization the number of breaks is checked separately for variables  $E$  and  $Z$ , while in Main Homogenization the break-list is always identical for  $E$  and  $Z$  until the final filtering of insignificant breaks, for technical reasons. For the limitation of synchronous breaks the least significant breaks are deleted from the break-list when it is necessary. The rank order of significance is determined by the origin of the break or with the calculation of  $\tau^*$ . If the timings of two breaks have 1 month difference only, then they are considered to be synchronous. In the Main Homogenization, a combination of the shift-sizes of  $E$  and  $Z$  is considered in  $d^2$  (26).

$$d^2 = d(E)^2 + d(Z)^2 \quad (26)$$

#### **I-5.4. Exclusion of breaks due to too short section of homogeneous period**

The minimum distance between two adjacent breaks of the break-list is 5 months. The origin of the breaks is considered for determining and excluding the less significant break, if it is necessary for complying with the rule. Note that Filtering of outlier periods and Secondary Detection can detect shorter IHs than 5 months, but the biases due to such IHs are treated with interpolation and the breaks bordering such IHs are never included in break lists.

#### **I-5.5. Exclusion of insignificant breaks**

The significance of breaks is checked by  $t$ -test, separately for  $E$  and  $Z$ , both in Pre-homogenization and Main Homogenization. In Main Homogenization, breaks with  $\tau$  indicating insignificant shift-size at the 0.05 level are considered to be insignificant. If two adjacent breaks are insignificant, then only one break (with the smaller  $\tau$ ) is excluded in one particular step, thereafter the check of significance is repeated applying the reduced break-list. If more than two sequent breaks are insignificant, then the first and last insignificant breaks are excluded, thereafter the check of significance is repeated.

In Pre-homogenization, only one break can be excluded in one specific step, i.e. the one with the lowest statistical significance. Here, modified  $t$ -statistics are calculated, i.e.  $d\tau$  is examined instead of  $\tau$  and the threshold statistic is lower (higher) for  $E$  ( $Z$ ) than the relevant threshold of  $\tau$  by the multiplier 0.842 (1.786). These modifications are based on test experiments.

#### **I-6. Adjustments before the homogenized section**

Relative homogenization often cannot be applied for early sections of time series when the density of observing network is inadequate. However, breaks of the homogenized section might be responsible for biases both within the homogenized section and before that. Thus adjustments can be applied for early sections of time series, even when break detection was not performed for them. Such adjustments will improve the data accuracy when the impact of detected breaks is not overwritten by some undetected breaks in the early sections.

##### **I-6.1. Concepts of treated section and homogenized section**

For treated sections the ratio of missing data and the length of data gaps are limited. Observed data out of the treated section do not take part in any calculation and they remain unchanged during the homogenization procedure. The homogenized section can be shorter than the treated section when the number of comparable time series and their spatial correlations are inadequate to create reliable reference series for some periods of the treated section. Each input time series has 0 or 1 treated section including 0 or 1 homogenized section (see their deduction in the Manual). Periods of the treated section out of the homogenized section are not subjected to break detection and outlier filtering, but the data of such periods may take part in gap filling, and they might be adjusted as well.

##### **I-6.2. Deduction of the adjustment terms for data before the homogenized section**



In ACMANT2 the mean estimated bias of the first 30 years of the homogenized section is considered to be “persistent bias” and the relevant adjustment is applied for the treated section before the homogenized section. If the length of the homogenized section is shorter than 30 years, then the persistent bias is zero by definition. The change of the adjustment-terms close to the beginning of the homogenized section is gradual to avoid creating seeming breaks due to the rapid alteration of adjustment terms. For this reason, adjustment terms change linearly in the 3 years before the homogenized section (from the adjustment terms of the first year of the homogenized section to those of the persistent bias). From the fourth years before the homogenized period the adjustment terms of the persistent bias will be applied backwards until the beginning of the treated section.

### **I-7. Deduction of daily adjustment terms**

Monthly adjustment terms equal with the daily adjustment term in a middle day of months, more precisely: on 15 January, 14 February, and on 16 (15) of other months in non-leap years (leap years). These days are named middle days. For any other day of the year the adjustment term is calculated with linear interpolation between the adjustment terms of the two closest middle days.

## **APPENDIX II. DIFFERENCES IN HOMOGENIZING T<sub>min</sub> RELATIVE TO HOMOGENIZING T<sub>max</sub> OR T<sub>mean</sub>**

### **II-1. Detection on low time resolution**

In homogenizing T<sub>min</sub> (or homogenizing temperatures of tropical or monsoonal regions) always the univariate detection (formulas 3, 4 and 5) is applied.

In the Main Detection, the coefficient in the modified Caussinus – Lyazrhi criterion (eq. 10)  $p = 1.4$ .

In the Pre-homogenization eq. (18) is applied for determining  $p$ , but with modified parameterization:

$$p_1 = 2.6; \quad p_2 = 1.3; \quad p_3 = 20.0$$

### **II-2. Monthly precision**

Two-phase step function is fitted instead of harmonic functions. There is no change here in the parameterization.

### **II-3. Filtering of outlier periods**

Eqs. (19) and (21) are simplified here to (27), as the modification of  $l$  due to seasonal accumulation of biases is not applicable here.

$$\lambda = l^{0.75} d^2 \quad (27)$$

Another change is that in phase ii of this routine always step functions are fitted (and never harmonic functions).

There is no change here in the parameterization.

### **II-4. Secondary Detection**

In searching the most likely break positions around the maximum of accumulated anomalies, always step functions are fitted.

Changes in the parameterization:

Threshold for 5-month accumulated anomalies: 2.15

Threshold for 10-month accumulated anomalies: 1.5

$p$ -coefficient of the modified Caussinus – Lyazrhi criterion: 2.0.

### **II-5 Calculation of adjustment terms**

The biases of  $E$  only are calculated. In the model of  $T_{min}$  homogenization the annual cycle of biases is zero, therefore the adjustment term is constant within homogeneous subperiods.

# **HOMOGENIZATION OF MONTHLY TEMPERATURE SERIES IN ISRAEL - AN INTEGRATED APPROACH FOR OPTIMAL BREAK-POINTS DETECTION**

**Yizhak Yosef, Isabella Osetinsky-Tzidaki, Avner Furshpan**

Israel Meteorological Service, P.O.B. 25 Bet-Dagan, 5025001, Israel  
(yosefy@ims.gov.il)

## **Abstract**

In 2013 the Israel Meteorological Service (IMS) began using the homogenization methods systematically. After an examination of several common homogenization methods recommended by WMO and ACTION COST-ES0601, a procedure for optimal break-points detection has been developed for the monthly maximum and minimum temperature time series. The present work describes this procedure along with a few results obtained for the period of 1950-2012.

In our first experiments, it was found out that the absolute homogeneity tests applied to the temperature series as recorded in the Israeli meteorological stations gave insufficient results. Therefore, the relative methods which refer to the reference stations have been chosen.

Our approach for optimal break-points detection integrates a number of advanced homogenization methods: ACMANT, HOMER, RHtestsV3 and AnClim. The reference series were based on more than 30 stations. A cluster analysis was applied to find the most suitable reference stations for each base station. In making the final decisions on the break-points' locations, we were relying on the exclusive reliable metadata found in the IMS archive. Sometimes, however, the finally established location of a break-point was not among the events documented in a station's recorded history. After establishing the optimal (most approved) break-points' locations, the adjustment step of the homogenization procedure was carried out.

## **1. INTRODUCTION**

Israel is located in the subtropical region next to the southeastern corner of the Mediterranean Sea, and its climate is varying from the Mediterranean climate in the northern and central parts of the country through the semiarid to arid climate in the southern and southeastern parts. Israel's climate is affected as well by the complex topography over a very small area: A coastal plain in the west through a mountain range that goes from north to south, in the central part of the country, and a deep depression - the Jordan – Dead Sea Valley, in the east. All these factors produce a wide climatic variety all over the country. It should be mentioned that the dynamics of urban development and intensive industrialization along with the expansion of agricultural activity and afforestation in the last sixty years, have produced dramatic changes in the country's landscape. All these complex factors make the analysis and the construction of a homogeneous series in Israel quite a challenge.

In general, inhomogeneity in the time series can be caused by several factors. In Israel, the typical factors causing inhomogeneity in the temperature data series are:

- Relocation – almost all of our stations have changed their location during the stations' history, sometimes more than once.
- Instrumentation – there were many thermometer replacements and calibrations recorded in the stations' history. A common replacement of manual instruments with electronic sensors that began during 1990s caused a break-point in most cases. Upgrading the electronic sensors in the following years was sometimes a source for additional break-points as well.
- Change in screen design – a replacement of the original Stevenson screen with another type may cause a break-point in the temperature series.
- In addition to these key factors, there were maintenance problems and changes in the station's vicinity including gradual changes like urbanization.

In light of the aforementioned problems, a systematic use of the homogenization methods was adopted at the IMS in 2013. After examining several common homogenization methods recommended by WMO and ACTION COST-ES0601, a procedure for optimal break-points detection has been developed for the monthly maximum and minimum temperature time series.

Three main problems have been found while carrying out the homogenization procedure in the relative mode: (a) scarcity of neighboring stations from the same climate region as that of the base station, (b) lack of stations with long temperature records, especially during the 1950s and backward and (c) discontinuity of the data. In some cases there were just fragments of records, stations and periods. This last issue made it difficult and sometimes even impossible to construct a reference series because there was no common period for all the neighboring stations' time series (hereinafter NSTS).

The aim of this work was to develop a technique enabling the best break-point location through an integration of the most suitable features of several homogeneity methods. The integrated homogenization model accompanied with a few examples of its application is described in Sections 2 and 3.

## **2. METHODOLOGY**

### **2.1. Quality control**

Most of the temperature data (from both the base and neighboring stations) were undergoing a systematic quality control procedure. In addition, the HOMER fast quality control tool (Mestre et al., 2013) based on CLIMATOL (Guijarro, 2011) was applied to analyze the outliers, histograms and boxplots. The outliers detected by ACMANT (Domonkos, 2011) were analyzed as well.

## 2.2. Homogeneity methods and software

The integrated approach proposed here is based on a combination of four main methods:

- AnClim (Štěpánek, 2008) – This software contains several common homogeneity tests such as the SNHT (Alexandersson, 1986), Easterling-Peterson test (Easterling and Peterson, 1995) and Vincent test (Vincent, 1998). In our study, this software mainly serves for building the reference series and performing some basic absolute and relative homogeneity tests. The final decision on a break-point location with AnClim is being made after at least two different tests (using different methods) located at the same specific break-point.
- RHtestsV3 (Wang and Feng, 2010) – This method is based on the penalized maximal  $t$  or  $F$  test (Wang et al., 2007; Wang, 2008). These tests are applied in both the absolute (PMF) and relative (PMT) mode. The RHtestsV3 preliminary results of the statistically identified dates of the break-points are verified versus the documented dates and fixed where needed. Then the significance of small shifts is reassessed.
- ACMANT (The Adapted Caussinus-Mestre Algorithm for Networks of Temperature series) – One of the most recommended methods by WMO and COST ACTION and was found among those achieving very good results in the 2012 benchmark (Venema et al., 2012). This method is based on a bivariate detection of changes that includes a penalty term (Caussinus and Mestre, 2004). The ACMANT works fully automatically and its results are based on the stations combining the network.
- HOMER (HOMogenization software in R) – One of the latest advanced methods that includes the finest features from several leading methods like PRODIG (Caussinus and Mestre, 2004), ACMANT and joint segmentation method (Picard et al., 2011). HOMER is an interactive method, which takes advantage of metadata. There are some subjective decision parts where an expert intervention is required.

## 2.3. The homogenization model at the IMS

The integrated model includes four main methods as described above: AnClim, ACMANT, RHtestsV3 and HOMER. The first step is an application of the absolute tests using AnClim and RHtestsV3. The second step is using the relative methods which are the core of this procedure. The whole model is presented in Figure 1:

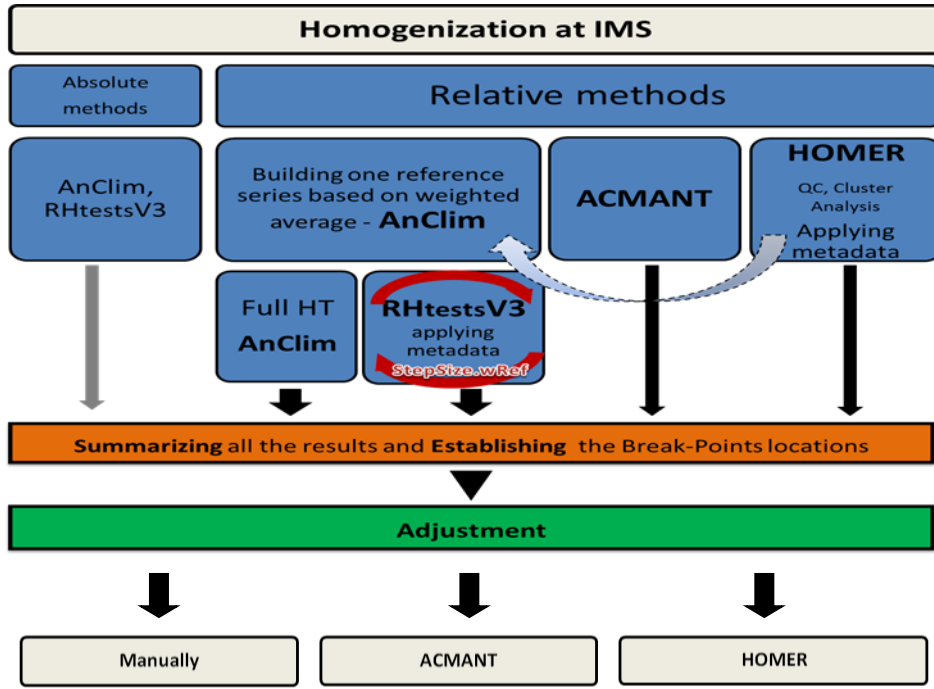


Fig. 1. The IMS homogenization model. "HT" stands for "homogeneity tests" and "QC" for "quality control".

### 2.3.1. Building the reference series

First, an initial set of the reference stations, located in the same climate region as that of the base station, is built. Secondly, the boxplot and cluster analysis are applied to this initial set in order to eliminate the less related stations, using the '*fast climatol check*' of HOMER. Then a weighted average based on the squared correlation coefficient ( $r^2$ ) is calculated between the base and each of the NSTS:

$$Q_i = y_i - \left\{ \frac{\sum_{j=1}^k r_j^2 [x_{ji} - \bar{x}_j + \bar{y}]}{\sum_{j=1}^k r_j^2} \right\}$$

$y_i$  – base series value at  $i$ -th time.  $i=1, \dots, n$ .

$\bar{y}$  – average temperature of the base series.

$Q_i$  – temperature value at  $i$ -th time.

$x_j$  – the  $j$ -th neighboring station time series.  $j=1, \dots, k$

$k$  – total number of neighboring stations.

$\bar{x}_j$  – average temperature of the  $j$ -th neighboring station.

$r_j$  – correlation coefficient between the base and  $j$ -th neighboring station.

After building the reference series (using AnClim), it is transformed into the reference temperature anomaly series.

### **2.3.2. Homogeneity tests using AnClim, RHtestsV3 and ACMANT**

At this step, the relative tests are applied to the base and reference temperature anomaly series. The significant outputs of AnClim are summarized and then RHtestsV3 is applied. The RHtestsV3 allows a user to modify manually the break-points' locations according to metadata. After the objective break-points detection with RHtestsV3, a modification of the dates is done on the base of the reliable metadata. The '*StepSize.wRef*' function is used to reassess the significance of the updated break-points' dates.

In parallel with the AnClim and RHtestsV3, ACMANT is also applied. The ACMANT performance is fully automatic and it uses a composite reference series for spatial comparisons. Application of ACMANT is done twice: first, with an automatic outliers' filtering and then, without it, after the removal of the manually approved ones. The outputs of the second run are taken into account as the final results of this test.

### **2.3.3. Homogeneity tests using HOMER**

The use of HOMER obliges an expert to make some subjective decisions. The consequences of wrong decisions may lead to a false break-point detection and an impaired adjustment. After gaining experience with the three methods described in 2.3.2, we have a good knowledge about the break-points' locations, so it can be assumed that we are capable to make better decisions at the subjective parts of HOMER.

HOMER is used mainly for verification of our results through (a) comparison with other methods, (b) analysis of our NSTS using a pairwise detection, and (c) comparison of the calculated correction factors at the adjustment step.

### **2.3.4. Summarizing the outputs and establishing the final break-points' locations**

At this step, we summarize all the described outcomes and cross-check them with our metadata. It should be noted that the metadata comes into consideration only after the detection phase in order to validate and support the results. At this step, the final establishment of the optimal break-points' locations is made. The IMS archive contains exclusively reliable metadata. However, this archive is not complete, therefore several final break-points have no metadata support.

### **2.3.5. Adjustment**

After the final establishment of the break-points' locations, we proceed to the adjustment step. It can be performed either manually or automatically with ACMANT or HOMER. The manual adjustment is based on the mean differences between the base and reference series. The same principle is applied in RHtestsV3. The distinction between the RHtestsV3 and the manual technique is that in the latter, each month is associated with its specific correction

factor while the RHtestsV3 uses the same mean correction factor (an annual average) for all the months. Both the ACMANT and HOMER make an automatic adjustment. Normally, we prefer the manual adjustment where (a) there is a lack of stations with a common period to perform a full ACMANT/HOMER run and/or (b) the final break-points' dates are resulted from a combination of different methods. However, in few cases the ACMANT adjustments are being used, especially where there are short time fragments of the NSTS making the building of the long united reference series and the derivation of the correction factors (for the entire period) almost impossible.

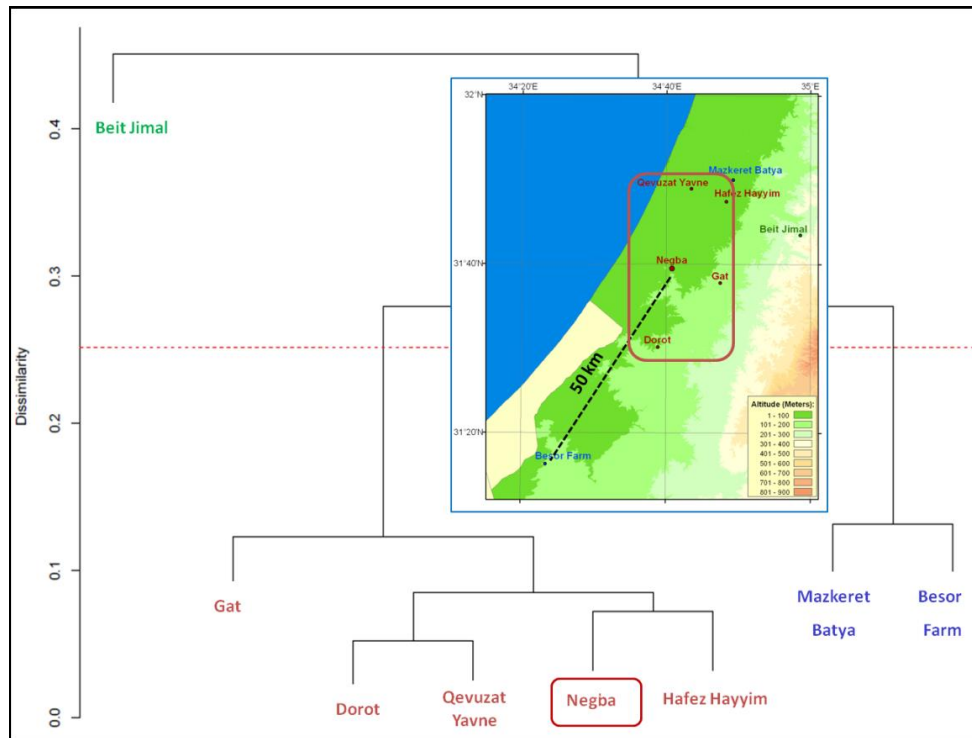
### **3. RESULTS**

In this section, applications of several specific blocks of the homogenization model (Figure 1) are presented: use of the cluster analysis, RHtestsV3, HOMER, establishing the break-points' locations and adjustment. Also shown are the final results for the Negba maximum and minimum temperatures.

#### **3.1. Cluster analysis**

After choosing the most correlated neighboring stations located in the same climate region of the base station, a cluster analysis was applied to improve the reference series. In Figure 2, such application to the Negba minimum temperatures is presented. The annual mean correlation coefficients between the base station and each of the neighboring stations of Beit Jimal, Besor Farm, Mazkeret Batya are quite high: 0.78, 0.84, 0.91, respectively. Despite the fact that we were intuitively tempted to use these stations, due to their proximity to Negba, a cluster analysis brought into consideration other stations as the preferred ones. It was found out that it is quite a frequent case where for a minimum temperature series there is no intuitive "hint". The spatial distribution of minimum temperatures typically has a local character, while that of the maximum temperatures usually represent a much wider area. In some cases there was no alternative, and due to a scarcity of the neighboring stations with common time periods, the less climatologically suitable stations were used. In such cases, the data adopted from a less suitable station (but still quite well correlated) were taken for the shortest time period as possible, only to complete the calculation.

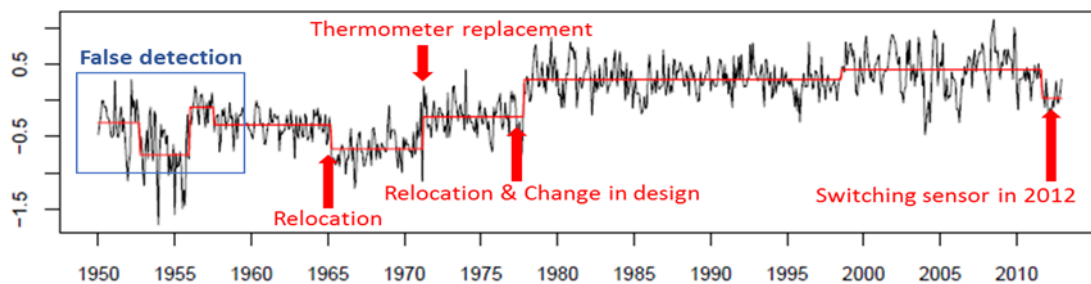




**Fig. 2. Cluster analysis for the Negba minimum temperatures. The red rectangle in the map embraces the stations selected by the cluster analysis.**

### 3.2. Analysis with RHtestsV3

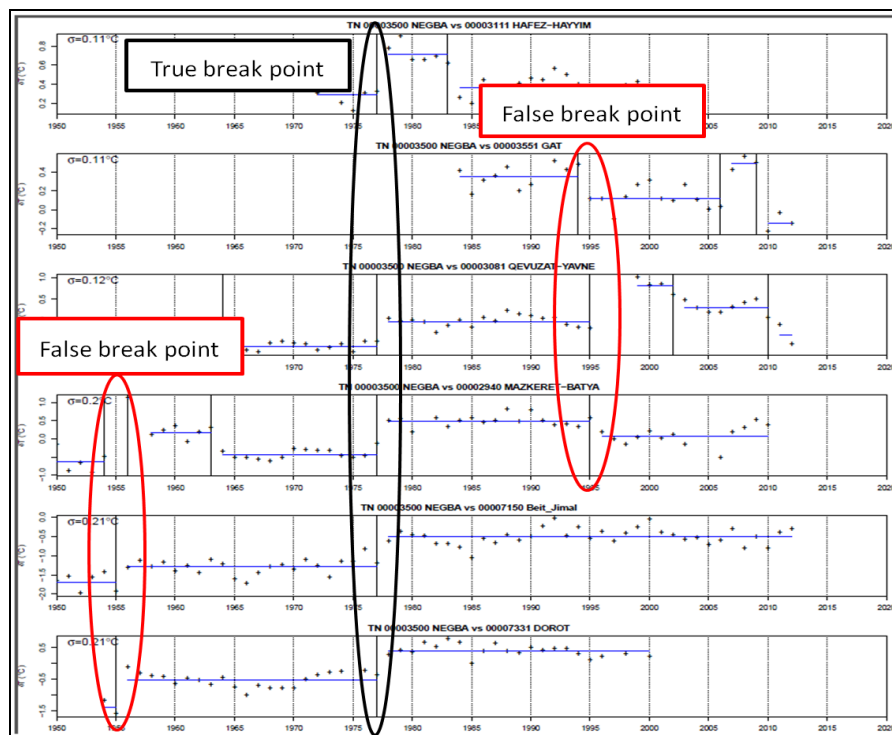
At this step, a base series meets with its reference series. The reference series consists of several NSTS. It is very important to know the number and quality of the neighboring stations included in each part of the reference series. The temperatures differences between the base and reference series obtained with the PMT method (Wang et al, 2007; Wang, 2008) are shown in Figure 3. In this case, it was decided to cut off the time series for the homogenization procedure by 2011, due to a break-point detected in 2012 which was caused by an electronic sensor replacement (the new sensor was found to be more sensitive). When the tested period was truncated in 2011, the break-point in 1998 became insignificant. Moreover, a metadata support was found for almost all the detected break-points. Several break-points detected at the beginning (up to the first 8 years) of the tested period were finally defined as false due to the inhomogeneity found in one, two, or three NSTS that comprised the reference series for that period (see 3.3). This is an example of how a small number of neighboring stations and their inhomogeneity can have a negative impact on the reference series that eventually may lead to a false outcome.



**Fig. 3. Negba: Monthly differences between the base and reference minimum temperature series (black), break-points and metadata (red).**

### 3.3. Pairwise detection using HOMER

Figure 4 introduces the pairwise detection (univariate detection) performed by HOMER. Each panel shows the annual differences between the base station and one of its neighbors. According to this figure, it seems like there are break-points in 1955, 1977, and 1995. The break-point of 1977 was found to be true, while those of 1955 and 1995 eventually appeared to be false. It was found out that for these two latter break-points, the source for inhomogeneity was in the neighboring time series whereas the base series was detected as homogeneous for those periods. This was concluded through analyzing the results of the pairwise detection for the corresponding neighboring stations (not shown). These examples show how the NSTS may influence the reference series and lead to false detections in the base series.



**Fig. 4. Pairwise detection using HOMER for the Negba minimum temperatures. Each panel shows the difference in the minimum temperatures between the Negba and each of its neighboring stations. Vertical black lines represent the break-points.**

### 3.4. Establishing the optimal break-point location

This is the final step of the break-points detection. Two examples are given in Tables 1 and 2. According to these tables, it is possible to obtain the optimal locations of the break-points. Table 1 summarizes the break-points, relevant metadata and methods for the Negba annual minimum temperature (Tn). As mentioned above, the metadata was used to validate the results, only after the detection phase. The break-points' locations were not forced by the metadata, to avoid any influence on the objective detection of the significant changes. In this case, the finally established break-points' locations were according to ACMANT, because there were (a) metadata support, (b) similarity to the results obtained with other methods and

(c) agreement with the definition of the 1955 and 1957 break-points as false (caused by inhomogeneity in the reference series, Figure 3).

*Table 1.* The final break-points' locations for the Negba Tn detected by different methods, and relevant metadata. "E.P" stands for Easterling and Peterson (1995), which was the only method in AnClim that spotted the 1964 and 1971 break-points. Bold 'V' represents the finally chosen (optimal) break-points.

Break points	AnClim	RHtestsV3	ACMANT	Metadata
1955	V	V		
1957		V		
1964	E.P	V	<b>V</b>	Relocation
1971	E.P	V	<b>V</b>	Thermometer replacement
1977	V	V	<b>V</b>	Relocation & Change in screen design

In Table 2, showing the results for the Zefat minimum temperature series, we established the optimal break-points' locations according to RHtestsV3 for three main reasons:

- 1) There was reliable metadata support for almost all the break-points.
- 2) The significant results were in common with other methods.
- 3) For the period 1990-2012, a reliable reference series has been obtained, comprising of 3 to 5 NSTS and resembling quite well the climate signal in that period.

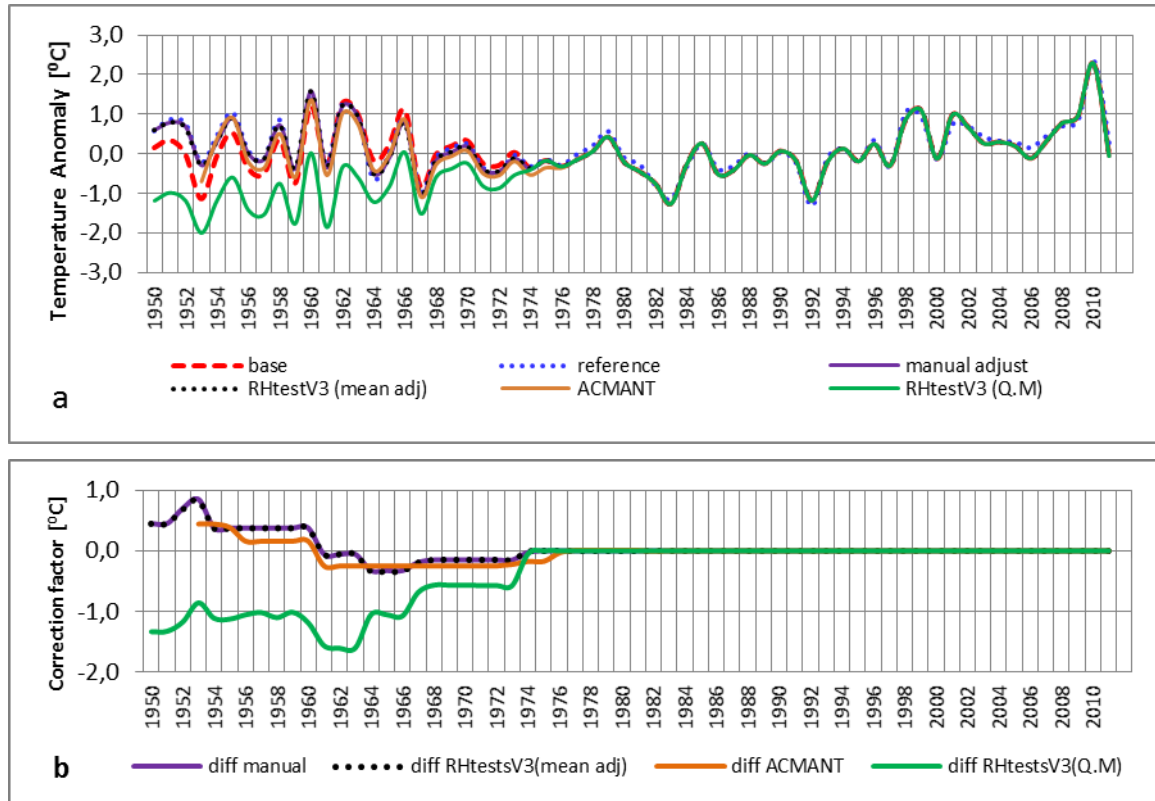
*Table 2.* The final break-points' locations for the Zefat Tn series, detected by different methods and relevant metadata. AWS – automatic weather station. Bold 'V' represents the finally chosen optimal break-points.

Break points	AnClim	RHtestsV3	ACMANT	Metadata
1990	V	<b>V</b>	V	No metadata
1992		<b>V</b>	V	Thermometer replacement
1995		<b>V</b>		Thermometer replacement
2000	V	<b>V</b>	V	Calibration
2004	V	<b>V</b>		Starting the use of the AWS data
2008		<b>V</b>		Electronic sensor replacement

### 3.5. Comparison of different adjustment methods

The results obtained with different adjustment methods for the Negba maximum temperature anomalies are displayed in Figure 5. It should be noticed that the signs and magnitudes of the correction factors are quite similar for all the methods (Figure 5b), except the quantile matching (RHtestsV3) which was found to be inadequate for our temperature series. In addition, it should be mentioned that the ACMANT considered our data only from 1953, because there were less than four NSTS for the period 1950 to 1952, which is not enough for

ACMANT's requirements. The manual and RHtestsV3 annual correction factors were identical (dissimilarities between these methods exist for the seasonal and the monthly adjustments, see 2.3.5).



**Fig. 5. The Negba annual maximum temperature anomalies' adjustment with different methods. Panel (a) shows four adjustment methods: manual, RHtestsV3 (both are based on mean adjustment), ACMANT (ANOVA) and RHtestsV3 (quantile matching). Panel (b) shows the annual correction factors [°C].**

### 3.6. Final results for the Negba maximum and minimum temperature series

The final results for the Negba temperature anomaly series are presented in Figure 6. The graphs show the base vs. the adjusted series for the maximum (Figure 6a) and minimum (Figure 6b) temperatures. The maximum temperature series has six break-points (green vertical lines), while the minimum temperature series has three break-points. The annual correction factors are summarized in Table 3.

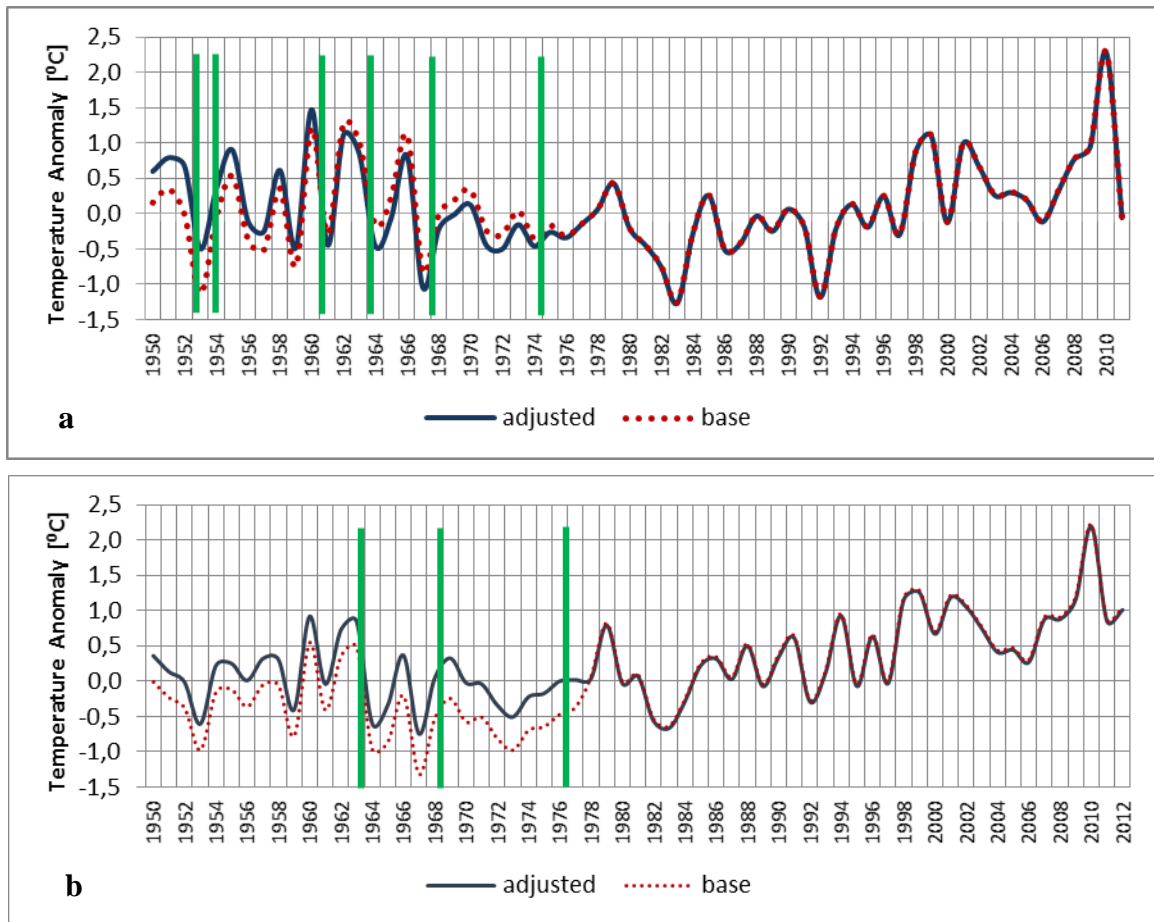


Fig. 6. The Negba temperature anomaly series, base vs. adjusted, (a) for maximum temperature and (b) for minimum temperature. The break-points' locations are marked with green vertical lines.

Table 3. The annual correction factors [°C] for the Negba maximum and minimum temperatures.

Parameter	Break-point	Correction factor [°C]
Tx	1952	0.40
	1953	0.97
	1960	0.40
	1963	-0.10
	1967	-0.33
	1974	-0.15
Tn	1965	0.39
	1970	0.63
	1977	0.47

## 4. CONCLUSIONS AND SUMMARY

This work presents the homogenization model developed at the IMS. This homogenization model is based on an integration of several advanced homogeneity methods. Such an approach enables raising, to the best of our knowledge, the reliability of break-points' locations. The absolute homogeneity tests were found to be insufficient for the Israeli long temperature series since they detected real climate signals as break-points. The relative homogeneity methods produced good results, especially when a cluster analysis was applied. An integrated approach allows merging the results obtained with different methods, getting the optimal break-points' locations, and minimizing the risk of a false break-point detection.

The adjustment may be performed either manually, subject to the possibility of building a long and reliable reference series, or with ACMANT or HOMER if the time series of the neighboring reference stations have too short common periods. In addition, these two latter methods helped us to improve the estimates for the correction factors.

The location of Israel in the subtropical region and the complexity of its climatic regime with several climate regions over such a small and narrow country make the homogenization procedure to be quite a challenge. This forces us to use different methods, according to availability and reliability of our data. With the integrated approach described in this paper, we can analyze and fix the long temperature series for different regions to find the optimal break-points' locations and to apply proper adjustments. That will enable us to construct a reliable long-term base series aiming to best understand the climate change in our region.

## References

- Alexandersson, H., 1986: A homogeneity test applied to precipitation data. *International Journal of Climatology*, 6, 661–675.
- Caussinus, H. and Mestre, O., 2004: Detection and correction of artificial shifts in climate series. *Applied Statistics*, 53, 405–425.
- Domonkos, P., 2011: Adapted Caussinus-Mestre Algorithm for Networks of Temperature series (ACMANT). *International Journal of Geosciences*, 2, 293–309, DOI: 10.4236/ijg.2011.23032.
- Easterling, D.R. and Peterson, T.C., 1995: A new method for detecting undocumented discontinuities in climatological time series. *International Journal of Climatology*, 15, 369–377.
- Guijarro, J. A., 2011: User's guide to climatol. An R contributed package for homogenization of climatological series. State Meteorological Agency, Balearic Islands Office, Spain. <http://www.climatol.eu/climatol-guide.pdf>
- Mestre, O., Domonkos, P., Picard, F., Auer, I., Robin, S., Lebarbier, E., Böhm, R., Aguilar, E., Guijarro, J., Vertachnik, G., Klancar, M., Dubuisson, B. and Štěpánek, P., 2013: HOMER: A Homogenization Software - Methods and Applications. *Idojaras, Quarterly journal of the Hungarian Meteorological Service*, Vol. 117, No. 1, 47–67.
- Picard, F., Lebarbier, E., Hoebeker, M., Rigai, G., Thiam, B., and Robin, S., 2011: Joint segmentation, calling, and normalization of multiple CGH profiles. *Biostatistics*, 12, 413–428.

- Štěpánek, P., 2008: AnClim - software for time series analysis. Dept. of Geography, Fac. of Natural Sciences, MU, Brno. 1.47 MB. <http://www.climahom.eu/AnClim.html>
- Venema, V., Mestre, O., Aguilar, E., Auer, I., Guijarro, J.A., Domonkos, P., Vertacnik, G., Szentimrey, T., Štěpánek, P., Zahradnicek, P., Viarre, J., Müller-Westermeier, G., Lakatos, M., Williams, C.N., Menne, M., Lindau, R., Rasol, D., Rustemeier, E., Kolokythas, K., Marinova, T., Andresen, L., Acquaotta, F., Fratianni, S., Cheval, S., Klancar, M., Brunetti, M., Gruber, C., Duran, M.P., Likso, T., Esteban, P. and Brandsma, T., 2012: Benchmarking monthly homogenization algorithms. *Climate of the Past*, 8, 89–115.
- Vincent, L. A., 1998: A technique for the identification of inhomogeneities in Canadian temperature series. *Journal of Climate*, 11, 1094–1104.
- Wang, X. L., Wen, Q.H., Wu, Y., 2007: Penalized maximal t test for detecting undocumented mean change in climate data series. *Journal of Applied Meteorology and Climatology*, 46 (No. 6), 916–931. DOI:10.1175/JAM2504.1
- Wang, X. L., 2008: Accounting for autocorrelation in detecting mean shifts in climate data series using the penalized maximal t or F test. *Journal of Applied Meteorology and Climatology*, 47, 2423–2444.
- Wang, X. L. and Y. Feng, published online January 2010: RHtestsV3 User Manual. Climate Research Division, Atmospheric Science and Technology Directorate, Science and Technology Branch, Environment Canada.

# THE WMO/MEDARE INITIATIVE: BRINGING AND DEVELOPING HIGH-QUALITY HISTORICAL MEDITERRANEAN CLIMATE DATASETS INTO THE 21<sup>ST</sup> CENTURY

**Khalid Elfadli<sup>1</sup> & Manola Brunet<sup>2</sup>**

<sup>1</sup>Climate & Climate Change adviser at Libyan National Meteorological Center + Astronomical & Meteorological Dep. At Cairo university.

<sup>2</sup>MEDARE co-chair and WMO/CCI OPACE2 co-chair + Centre for Climate Change (C3), Dep. of Geography, University RoviraiVirgili, Tarragona, Spain + Climatic Research Unit, School of Environmental Sciences, University of East Anglia, Norwich, UK

## **Abstract**

The Greater Mediterranean Region (GMR) has a very long and rich history in monitoring of the atmosphere, going back in time several centuries in some countries and at least to the mid of 19<sup>th</sup> century across much of the GMR.

However, despite the efforts undertaken by National Meteorological and Hydrological Services (NMHS), research centres, universities and motivated individuals in Data Rescue (DARE) activities, available and accessible digital climate data are still mostly restricted to the second half of the 20<sup>th</sup> century for a few countries and since 1970's for most of the GMR. This reality is preventing the region from developing more robust, accurate and reliable assessments of climate variability and change and its adverse impacts on the socio-ecosystems of the Mediterranean Basin, at the same time it is impeding the development of optimum strategies to mitigate and/or adapt the countries to the current and future impacts of climate change.

In addition, the fragmentation and scarcity of long-term and high-quality surface climate records is hampering our ability for better detecting, predicting and adapting the countries to present and future impacts of climate variability and change as well. This is particularly over this climate change 'hot-spot' region.

The WMO/Mediterranean (climate)DATA REscue (MEDARE) Initiative was set up to address developing, accessible and traceable comprehensive long and high-quality instrumental surface climate datasets for the GMR.

The MEDARE community exercises and implements its functions and actions throughout (4) working groups (WG1-WG4) under leading of rotational steering group (SG).

This structure has allowed the MEDARE community to undertake many other organizational, implementation and dissemination activities in order to raise awareness on the importance of bringing historical climate datasets into the 21<sup>st</sup> century, which is paving the way to get achieved the MEDARE's end-goal and objectives.

Among very important objectives of MEDARE initiative are represented in the following lines:

- To develop comprehensive, long and high-quality surface climate datasets for the GMR with a focus on the relevant essential climate variables (i.e Temperature, precipitation, air and sea pressure, .. etc.) of the Global Climate Observing System (GCOS) at different scales of time, which are currently required to support the work of the UNFCCC, the IPCC and the WMO/World Climate Program (WCP);



- To seek and mobilise resources and efforts at the national, regional and international scales in support of Data Rescue and Homogenisation (DARE&H) of long and key climate records over the GMR.

MEDARE web-site for linking the MEDAREcommunity already implemented, updated and maintained, while the on-line MEDARE portal metadata base infrastructure for the longest and key Mediterranean climate records: about 700 sites documented for mainly Tx(max temp)/Tn(min temp), RR(rainfall) and SLP(sea level pressure) at daily (sub-daily) scales, populated to be used by scientists, stakeholders, policy-makers and the general public within the region.

Other efforts for recovering, digitising, quality controlling and homogenising total of 38 daily Tx and Tn time-series for various locations in the southern and eastern parts of the Mediterranean Basin, where their recent part extends into the first decade of the 21<sup>st</sup> century while for some of them data are available since the late part of the 19<sup>th</sup> century, are being completed under the EU-funded European Reanalysis and Observations for Monitoring (EURO4M) project, linked to the World Meteorological Organization (WMO) Mediterranean (climate) Data REscue (MEDARE) Initiative.

Finally, build up the Mediterranean climate databases for GMR is being the end goal of the MEDARE Initiative.

## 1. INTRODUCTION

The Mediterranean (climate) Data REscue (MEDARE) Initiative is:-

- A joint-WMO effort (established on November, 2007) whose common goals being the enhancement of bringing historical climate datasets into the 21<sup>st</sup> century, which is paving the way to get achieved the MEDARE's end-goal of building up the Mediterranean climate databases of Greater Mediterranean Region (GMR);
- Following the MEDARE recipe: bringing together climatologists and scientists from Mediterranean NMHS & Academia to exchange their experiences (both theoretical and operational) on DARE;
- promoting a new culture of data and knowledge sharing within GMR;
- Non-regularly-funded WMO project and run on a volunteer basis.

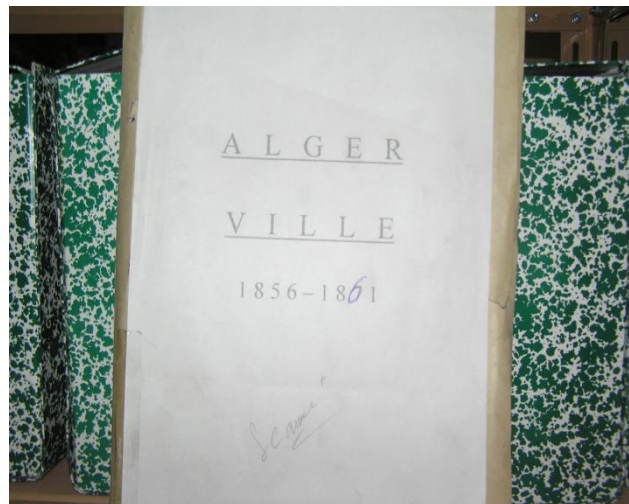
MEDARE also integrated by most of the Mediterranean NMHS (Albania & Portugal not included yet) and endorsed by WMO EC-60 (June, 2008) and quoted by GFCS (2013) as one of DARE initiatives to be supported.

## 2. THE RATIONALE & NEED FOR MEDARE

Mediterranean region has a long and rich history in monitoring of the atmosphere, going back in time several centuries in some countries and at least to the mid of 19<sup>th</sup> century across much

of the GMR. However, limited availability and accessibility of long and high-quality climate series represents the biggest challenge in the region.

This is hampering progress on our capability to detect, predict & adapt the countries to the impacts of climate variability & change and it is limiting the timely delivery of



**Oldest climate data source in the Algerian NMHS archive (1856-1861). Courtesy of Mehdi Kerrouche**

climate products and services. Following factors could be considered the key challenges:

- short period climate records (e.g. from 1970's onwards) availability and accessibility;
- poor spatial coverage (limited observing stations), especially over southern and south-eastern Mediterranean countries ;
- lack of quality climatic time series .

### **3. MEDARE COMPOSITION**

MEDARE community composed of 37 organizations, including 25 Mediterranean NMHS and 11 research centres with about 100 individual members and 4 working groups (WG) as follows:

- WG1. interests with inventorying/assessing/approaching of old material sources and holders;

- WG2. interests with DARE techniques and procedures (including digitization);
- WG3. interests with approaches on best practices for quality controlling and homogenizing specific climate variables;
- WG4. interests with promotional activities, bringing MEDARE to the wider scientific and other communities.

Steering Group (SG) leads all MEDARE community activities; the 2<sup>nd</sup> SG is composed of:

- *Manola Brunet & SerhatSensoy (Co-chairs)*
- *Victor Venema (University of Bonn)*
- *Athanasios Sarantopoulos (Greece NMHS)*
- *Fatima Elguelai (Morocco NMHS)*
- *Khalid Elfadli (Libya NMHS)*
- *Yolanda Luna (Spain NMHS)*
- *JanjaMilkovic (Croatia NMHS)*
- *DjamelBoucherf (Algeria NMHS)*
- *MesutDemircan (Turkey NMHS)*
- *Marius Theophilou (Cyprus NMHS)*



(MEDARE 2<sup>nd</sup> SG (Istanbul, Turkey, 27-28/Sep/2012))

#### 4. MEDARE OBJECTIVES

MEDARE has a wide spectrum of goals and objectives on regional and national scales which could be briefed as following:

- Fostering DARE projects at national, sub-regional and regional scales;
- Mobilising resources (human and financial) to undertake DARE projects over the GMR;
- Innovating on DARE techniques (from efficient data transfer into digital format to time-series QC and homogenisation);

- Capacity building through training activities (regional workshops and tailored training programs);
- Increasing awareness on the need for DARE among stake-holders and decision-makers (several dissemination material elaborated and distributed);
- Linking MEDARE members to other DARE initiatives to better coordination and avoid duplication with special links to ECA&D and ICA&D and WMO DARE-I.

## 5. MAIN PROGRESS & ACTIVITIES MADE SINCE THE LAST (3) YEARS

- i. Updating and maintaining the MEDARE web-site for linking the MEDARE Community
- ii. Defining, implementing and populating the MEDARE portal Metadata Base (<http://app.omm.urv.cat/urv>) with country and research projects metadata on long and key Mediterranean climate records: about 700 sites documented for mainly min. & max. temperature (Tn/Tx) and rainfall (RR) as well as sea level pressure (SLP) variables at daily (sub-daily) scales
- iii. Paving the way for MEDARE becoming a WMO/WIS Data Collection and Production Centre (DCPC)
- iv. Undertaking DARE activities over North Africa and Middle East countries under the opportunity that brought us by EU-EURO4M project (*EUropean Reanalysis and Observations for Monitoring project*) and in cooperation with NMHS in these areas

## 6. MEDARE

### 6.1.MEDARE web-page

MEDARE main online web-page established since 2008's for linking the MEDARE Community & joint users.



MEDARE main web-page

### 6.2.Medare portal metadata base

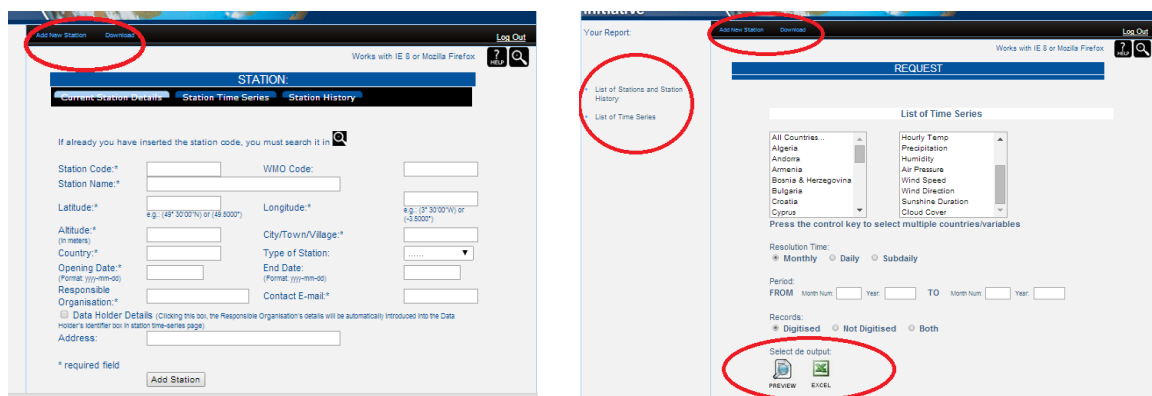
It is managed by C3/URV (*Centre for Climate Change (C3), University Rovira i Virgili, Tarragona, Spain*) and contributed mainly by Med. NMHS, but with a significant input from the DARE component of the EURO4M project: 35 metadata providers & 261 users.

Country name	No. of observing sites in MEDARE metadata base
Algeria	190
Andorra	7
Bulgaria	10
Croatia	13
Egypt	62
France	14
Greece	44
Israel	15
Italy	54
Jordan	12
Lebanon	3
Libya	28
Morocco	30
FYR of Macedonia	56
Slovenia	20
Spain	72
Tunisia	18
Turkey	8
<b>Total</b>	<b>656</b>

The portal is the on-line accessibility (psw protected), but accessible through Toulouse GISC, with remarkable improvement in coverage over southern and south-eastern areas. It is being also useful for identifying the “TARGET” records to be developed (digitised and homogenised) but this only contains METADATA, NO DATA.



### The MEDARE Metadata Base: contributions & access

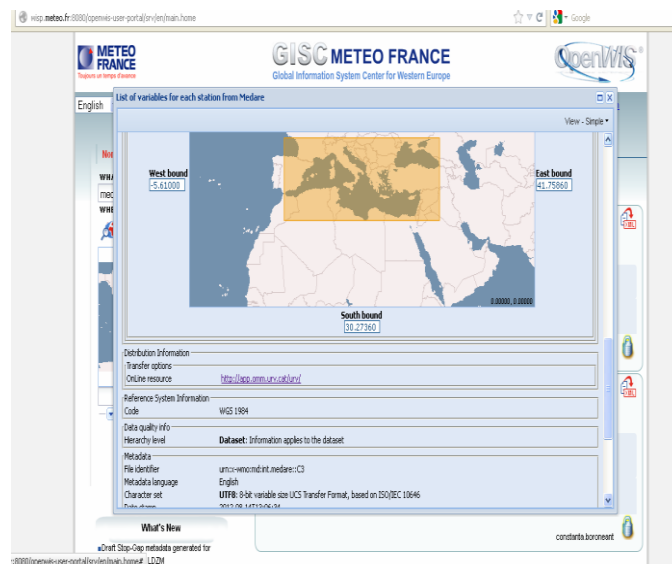


### The MEDARE portal and Metadata Base: easy access and using

MEDARE initiative												
Country:		Libya										
Climate Variables:		All Climate Variables										
Resolution Time:		Monthly										
Period:		All										
Digitised Record:		BOTH										
Country	City/Town/Village	Station Code	WMO Code	Station Name	Climate Variable	Start Date	End Date	Data Source	Latitude	Longitude	Altitude (m)	Digitised Record
Libya	Agedabia	96910009	62055	Agedabia	Air Pressure	1961.01.01	2009.12.31	Meteorological Department	30.7167	20.1667	07	Yes
Libya	Agedabia	96910009	62055	Agedabia	Cloud Cover	1946.05.01	2009.12.31	Meteorological Department	30.7167	20.1667	07	Yes
Libya	Agedabia	96910009	62055	Agedabia	Hourly Temp	1951.01.01	2010.12.31	Meteorological Department	30.7167	20.1667	07	Yes
Libya	Agedabia	96910009	62055	Agedabia	Humidity	1946.04.01	2009.12.31	Meteorological Department	30.7167	20.1667	07	Yes
Libya	Agedabia	96910009	62055	Agedabia	Maximum Temp	1946.04.01	2010.12.31	Meteorological Department	30.7167	20.1667	07	Yes
Libya	Agedabia	96910009	62055	Agedabia	Mean Temp	1946.04.01	2010.12.31	Meteorological Department	30.7167	20.1667	07	Yes
Libya	Agedabia	96910009	62055	Agedabia	Minimum Temp	1946.04.01	2010.12.31	Meteorological Department	30.7167	20.1667	07	Yes
Libya	Agedabia	96910009	62055	Agedabia	Precipitation	1946.05.01	2010.12.31	Meteorological Department	30.7167	20.1667	07	Yes
Libya	Agedabia	96910009	62055	Agedabia	Sunshine Duration	1965.06.01	2009.12.31	Meteorological Department	30.7167	20.1667	07	Yes
Libya	Agedabia	96910009	62055	Agedabia	Wind Direction	1956.01.01	2009.12.31	Meteorological Department	30.7167	20.1667	07	Yes
Libya	Agedabia	96910009	62055	Agedabia	Wind Speed	1949.01.01	2009.12.31	Meteorological Department	30.7167	20.1667	07	Yes
Libya	Bengazi	96910008	62053	Benina	Air Pressure	1961.01.01	2009.12.31	Meteorological Department	32.0833	20.2667	132	Yes
Libya	Bengazi	96910008	62053	Benina	Cloud Cover	1945.01.01	2009.12.31	Meteorological Department	32.0833	20.2667	132	Yes
Libya	Bengazi	96910008	62053	Benina	Hourly Temp	1951.01.01	2010.12.31	Meteorological Department	32.0833	20.2667	132	Yes
Libya	Bengazi	96910008	62053	Benina	Humidity	1945.03.01	2009.12.31	Meteorological Department	32.0833	20.2667	132	Yes
Libya	Bengazi	96910008	62053	Benina	Maximum Temp	1945.03.01	2010.12.31	Meteorological Department	32.0833	20.2667	132	Yes
Libya	Bengazi	96910008	62053	Benina	Mean Temp	1945.03.01	2010.12.31	Meteorological Department	32.0833	20.2667	132	Yes
Libya	Bengazi	96910008	62053	Benina	Minimum Temp	1945.03.01	2010.12.31	Meteorological Department	32.0833	20.2667	132	Yes
Libya	Bengazi	96910008	62053	Benina	Precipitation	1945.03.01	2010.12.31	Meteorological Department	32.0833	20.2667	132	Yes
Libya	Bengazi	96910008	62053	Benina	Sunshine Duration	1963.09.01	2009.12.31	Meteorological Department	32.0833	20.2667	132	Yes
Libya	Bengazi	96910008	62053	Benina	Wind Direction	1989.01.01	2009.12.31	Meteorological Department	32.0833	20.2667	132	Yes
Libya	Bengazi	96910008	62053	Benina	Wind Speed	1949.01.01	2009.12.31	Meteorological Department	32.0833	20.2667	132	Yes

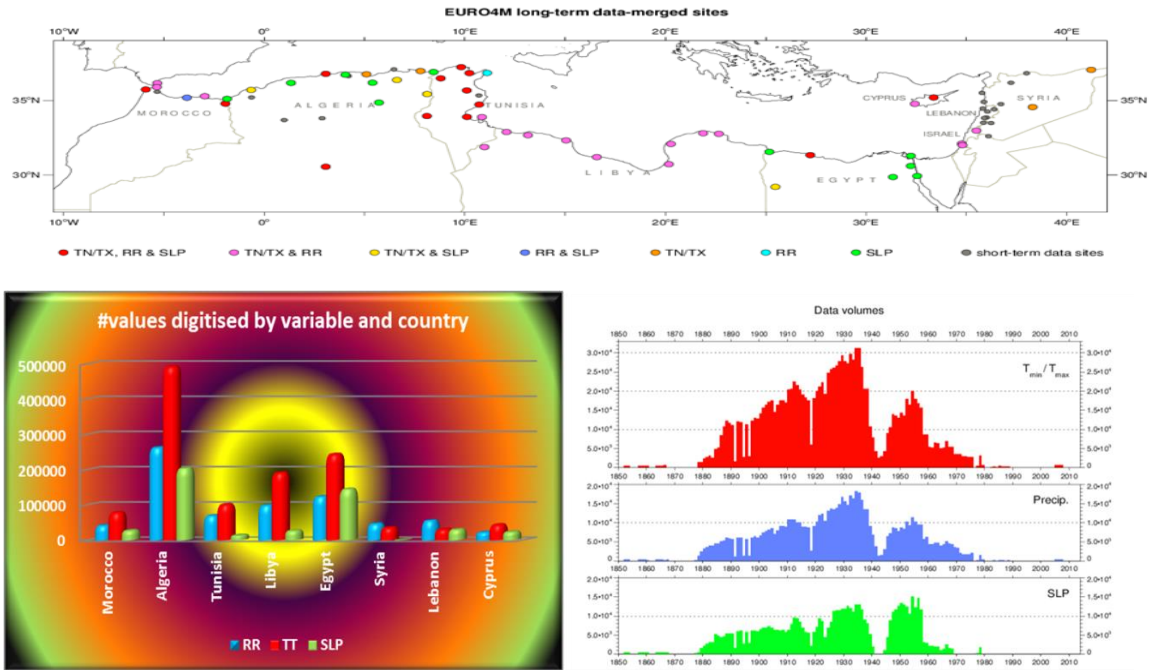
### 6.3.MEDARE as WIS/DATA collection and production centre (DCPC): steps taken &status

MEDARE endorsed by the *Spanish PR* (16<sup>th</sup> March 2011) and starting the process and fulfilling in the CBS Expert Team on GISC-DCPC Demonstration Process (ET-GDDP) questionnaire. A test account & the ET-GDDP audit of MEDARE metadata (in compliance with WMO/WIS standards) completed on July 2012, on August 2012, MEDARE metadata publicly available on the Toulouse GISC site: (<http://wisp.meteo.fr:8080/openwis-user-portal/srv/en/main.home>),and on the MEDARE portal: (<http://app.omm.urv.cat/urv/>).



### 6.4.The MEDARE datasets under development

Based on MEARE formula (bringing together NMHS & Academia) and under the EU-FP7 EURO4M project, the first ancient series are being recovered and populated at the MEDARE database.

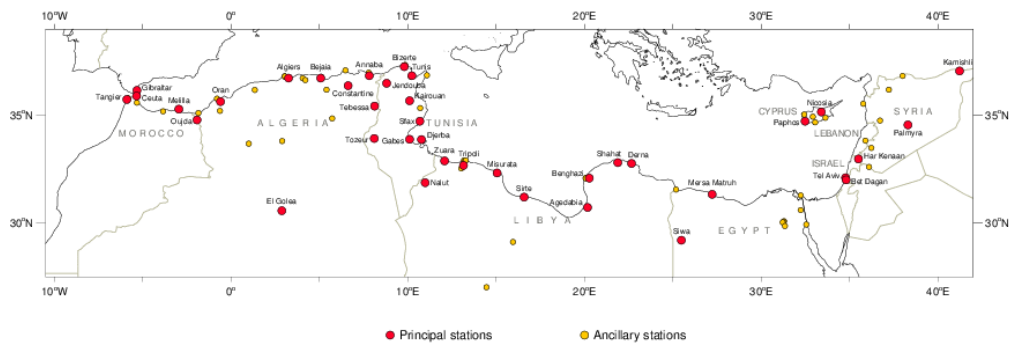


**MEDARE DARE activities over southern and south-eastern Mediterranean countries under EURO4M (recovering, digitising and quality controlling processes)**

The focus put on southern (North Africa) and eastern Mediterranean countries, involving the recovery & development of ancient climate daily (Tx/Tn, RR (about 65 series) and hourly SLP (38 series) from various sources (from digitisation to QC & homogenisation). Combining process of the ancient parts with recently observations fractions by data exchange agreements with several NMHS (e.g. Algeria, Cyprus, Libya, Jordan...) for developing long and high-quality climate time-series already started.

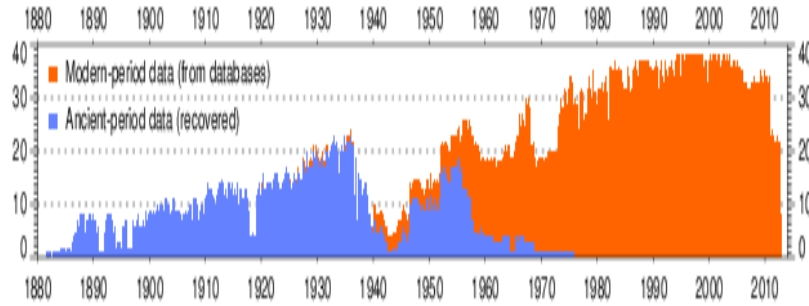
**6.4.1. MEDARE dare activities over southern and south-eastern Mediterranean countries under EURO4M (recovering, digitising and quality controlling processes)**

A total of 38 daily Tx and Tn time-series for various locations in the southern and eastern parts of the Mediterranean Basin have been selected; their recent part extends into the first decade of the 21<sup>st</sup> century, while for some of them data are available since the late part of the 19<sup>th</sup> century.



**Location of sites with Tx/Tn data series that have been merged and homogenised**





**Data availability for the sites selected**

### 6.4.2. Homogenisation methods & results (Cited by: D. Efthymiadis et al. 2013)

The daily time series selected have been converted into monthly means and then subjected to homogenisation following two approaches:

- (1) the ACMANT method;

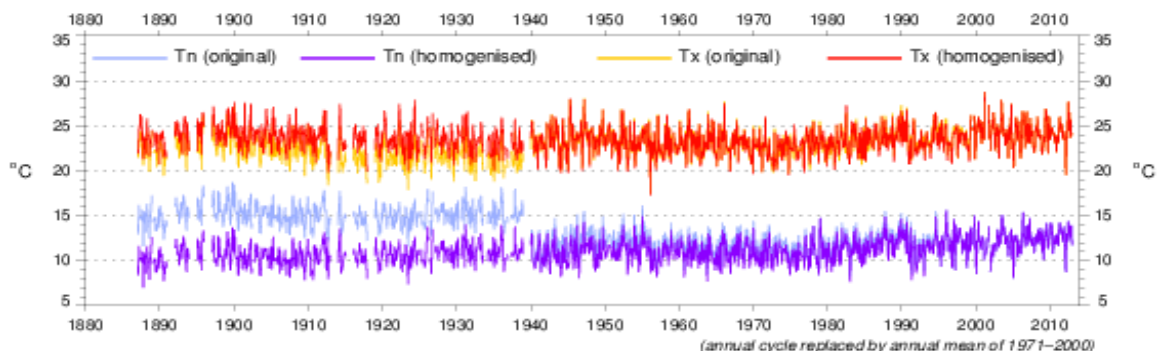
(<http://www.c3.urv.cat/data.html>)

- (2) the HomeR method.

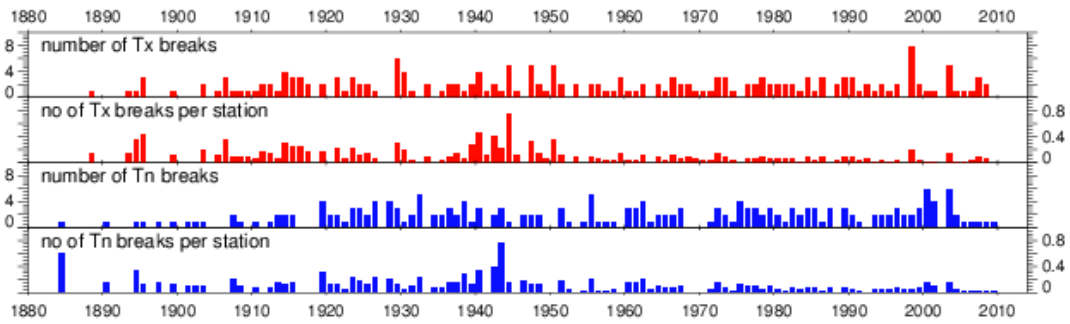
([http://www.homogenisation.org/v\\_02\\_15/index.php?option=com\\_content&view=article&id=93:homer&catid=1:general&Itemid=1](http://www.homogenisation.org/v_02_15/index.php?option=com_content&view=article&id=93:homer&catid=1:general&Itemid=1))

The homogenisation methods have identified a series of breaks and estimated adjustment factors which are necessary for making the original time series homogeneous over their overall span.

The density of breaks detected, i.e. the number of breaks per year, is similar in both the modern and ancient periods of data. However, since the stations network declines back in time the number of breaks per station available is higher in the early data periods and especially before the mid of 20<sup>th</sup> century.

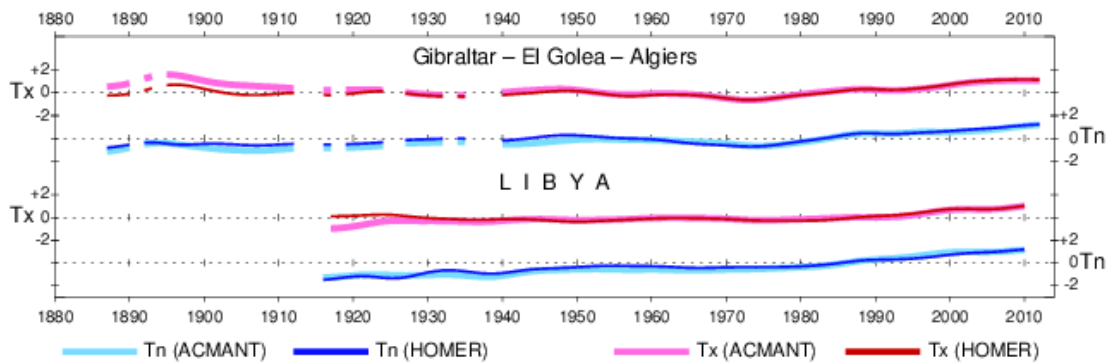


**Original and ACMANT-homogenised times series for Algiers (Algeria)**



**Number of ACMANT-detected breaks per year and per station for Tx and Tn.**

The stations meta-data over the ancient period are poorly documented in the data sources used for the data recovery, making difficult the attribution of breaks. Nevertheless, certain breaks coincide with stations relocation or stations-data merging points within the merged-data series been homogenized. While the two homogenization methods produce comparable results, discrepancies are also observed, especially in the data-sparse decades of the late of 19<sup>th</sup> and early of 20<sup>th</sup> centuries.



**Low-pass filtered temperature anomalies (°C) of homogenised series for the average of selected south-western Mediterranean stations (upper panel) and Libyan stations (lower panel)**

Another factor affecting the homogenization effectiveness is the various data gaps and the intermittent character of the data time series recovered. It is expected that as more stations data may additionally be recovered in this geographical area, the two methods' homogenized products will further converge resulting in time series of increasing reliability and thus suitable for long-term studies.

## 7. CONCLUSIONS

- MEDARE wants to contribute by enhancing GMR climate data availability and accessibility;
- Long-term & high-quality climate series are the basic input that underpin climate products and services;

- The elaboration of some changes in mean and extreme states of the climate or decadal climate prediction, demand the longest and more reliable climate information;
- Historical climate series should be contemplated as global cultural heritage to be preserved, developed and made freely accessible;
- The application of the ACMANT and HomeR methods on the long-term data-merged series leads to similar homogenized and Tx products;
- The data recovered together with existing data bases and other ongoing data-rescue efforts will provide an insight in the historical climatic variations over the southern and eastern parts of the Mediterranean Basin and will shed more light on the origins and the potential response of the overall Mediterranean climate to natural and anthropogenic forcing.

# HOMOGENIZATION OF SPANISH MEAN WIND SPEED MONTHLY SERIES

**José A. Guijarro**

Meteorological Agency (AEMET), Balearic Islands office, Spain  
(jguijarrop@aemet.es)

## **Abstract**

Monthly mean wind speed data were gathered from all Spanish series with a minimum of 10 years of data in the period 1951-2013, resulting in the selection of 233 series. Monthly wind speed averages were initially drawn from daily wind runs, but since they had too many missing data, mean wind speed recorded at 07, 13 and 18 UTC were obtained as well. These datasets were homogenized by means of the R package *Climatol* twice: 1) using a ratio normalization of the data; 2) applying a cubic root transformation to the data and standardizing them. Around two thirds of the series were found inhomogeneous through both normalization methods, which gave also similar results in terms of mean RMSE when estimating the series from the neighboring stations and mean SNHT of the final homogenized series. But the overall correlations of the wind series were not good enough, and showed a poor spatial coherence. Wind speed series were then extracted from the NCEP reanalysis to explore their potential value as reference series, but more than 80% of them were found inhomogeneous, probably because of their less noisy nature. Therefore, wind speed seems an element very prone to inhomogeneities, since it is very sensitive to obstacles and surface roughness changes in the surroundings of the observatories, and at the same time difficult to homogenize, because local air circulations as thermal winds may contribute to a significant part of the wind speed values, worsening the correlations between neighboring stations in complex regions. Anyway, wind speed trends were computed from these preliminary homogenization exercises, yielding negative figures mostly ranging between -1 and -2 m/s/century in the colder months of the year.

## **1. INTRODUCTION**

Wind is an important climatic element for many economic areas: agriculture (modulating evapotranspiration), water resources (controlling evaporation from dams and natural surfaces), leisure (outdoor activities, sailing, etc), and renewable energy production. For this reason, many work has been devoted to study its spatial and temporal variability (McVicar *et al.*, 2012, refer 148 papers on wind speed trends).

Wind speed has been traditionally measured in meteorological observatories with cup anemometers, although FUESS type used differential dynamic air pressure, and in recent times sonic anemometers are deployed as well. Changes of instrumentation or calibration drifts (e. g., increase of friction in the rotating axis) are a source of inhomogeneities in the series, as are instrument relocation or changes in the surroundings (new buildings, growing trees, etc), since wind is very sensitive to obstacles, orography, and surface roughness.

Yet many variability studies do not try to homogenize these series, but just to select those having a long period of observation which appear of a reasonable quality according to meta-data, visual inspection or basic comparison with a suitable reference (Dadaser-Celik and Cengiz, 2014).

Wan *et al.* (2010) did a thorough adjustment and homogenization of 117 Canadian wind stations with a minimum of 45 years of observation using the package RHtestV2 (Wang and Feng, 2007), and a recent paper by Azorín-Molina *et al.* (2014) also applied an homogenization package (AnClim, by Stepanek, 2004, using MM5 output as reference series) in their study of 67 wind speed series from Portugal and Spain selected for completeness in the period 1961-2011.

In this work, a more extensive homogenization is applied to most Spanish wind speed series, testing different approaches whose results are discussed, to end with a preliminary evaluation of the trends of the homogenized series.

## **2. METHODOLOGY**

### **2.1. Data**

Monthly mean wind speed data were gathered from all Spanish series with a minimum of 10 years of data in the period 1951-2013, resulting in the selection of 233 series. Monthly wind speed averages were initially drawn from daily wind runs, but since they had too many missing data, mean wind speed recorded at 07, 13 and 18 UTC were obtained as well. These latter values were an 8 % higher in average than those computed from daily wind runs. Figure 1 shows the number of data from both origins, the sharp increase in 1961 being due because data digitization from that year on were prioritized.

To complement observational series, wind speed monthly averages from NCEP reanalysis (Kalnay *et al.*, 1996) were also downloaded from NOAA servers.

### **2.2. Homogenization method**

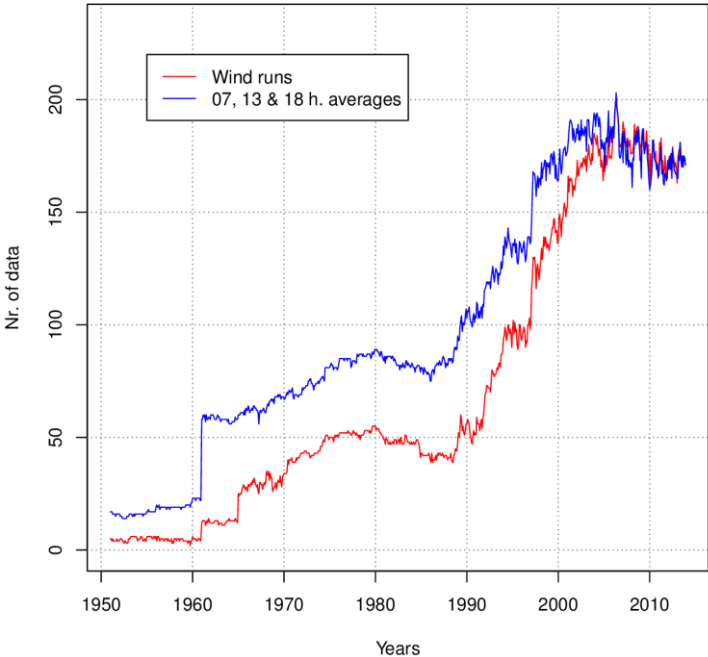
These series were homogenized with the 'Climatol' R package (Guijarro, 2014), that provides automatic quality control (outlier correction), homogenization (shift correction) and missing data attribution. The package begins by normalizing all data and computing a reference series for each observed series by averaging up to 10 data (if available) at every time step. As reference data are chosen by proximity, nearest data can be used even without any common period of observation with the problem series, taking advantage of short observational series that otherwise would be disregarded.

Series of anomalies are then computed by subtracting the reference series from the original series, allowing a simple detection of outliers (which are rejected if lying beyond a prescribed threshold) and breaks (shifts in the mean). Shift detection is performed by the well known SNHT test (Alexandersson, 1986), applied in stepped windows first to cope with multiple breaks, and then on the whole series to get all the power of the test.

These reference series are not assumed to be homogeneous, but only significantly less inhomogeneous than the original series. Therefore, an iterative application of the detection algorithm from big to small inhomogeneities in successive passes is performed, splitting the series at each noticeable break. Finally, newly computed reference series are straightforwardly used to fill any missing data in the series, including the reconstruction of the split series generated in the break detection process.

This methodology is able to yield results of a quality comparable with other good methods (as shown in <http://www.climatol.eu/DARE/testhomog.html>), and was applied to monthly wind speed series from wind runs (WRun), wind speed measured three times per day (WSm3) and NCEP reanalysis (WSRe), with two kind of normalizations: ratio to the series means, and full standardization. Ratio normalizations are normally used with variables with a zero lower limit and an L-shape probability distribution, while full standardization (removing the mean and dividing by the standard deviation) is applied to variables with a (near) normal distribution. Therefore, wind speed data were cubic root transformed (when greater than 1.0) in order to normalize their probability distribution.

Finally, trends of the homogenized series were obtained by regression with time, with the help of a post-processing function of the same computer package.

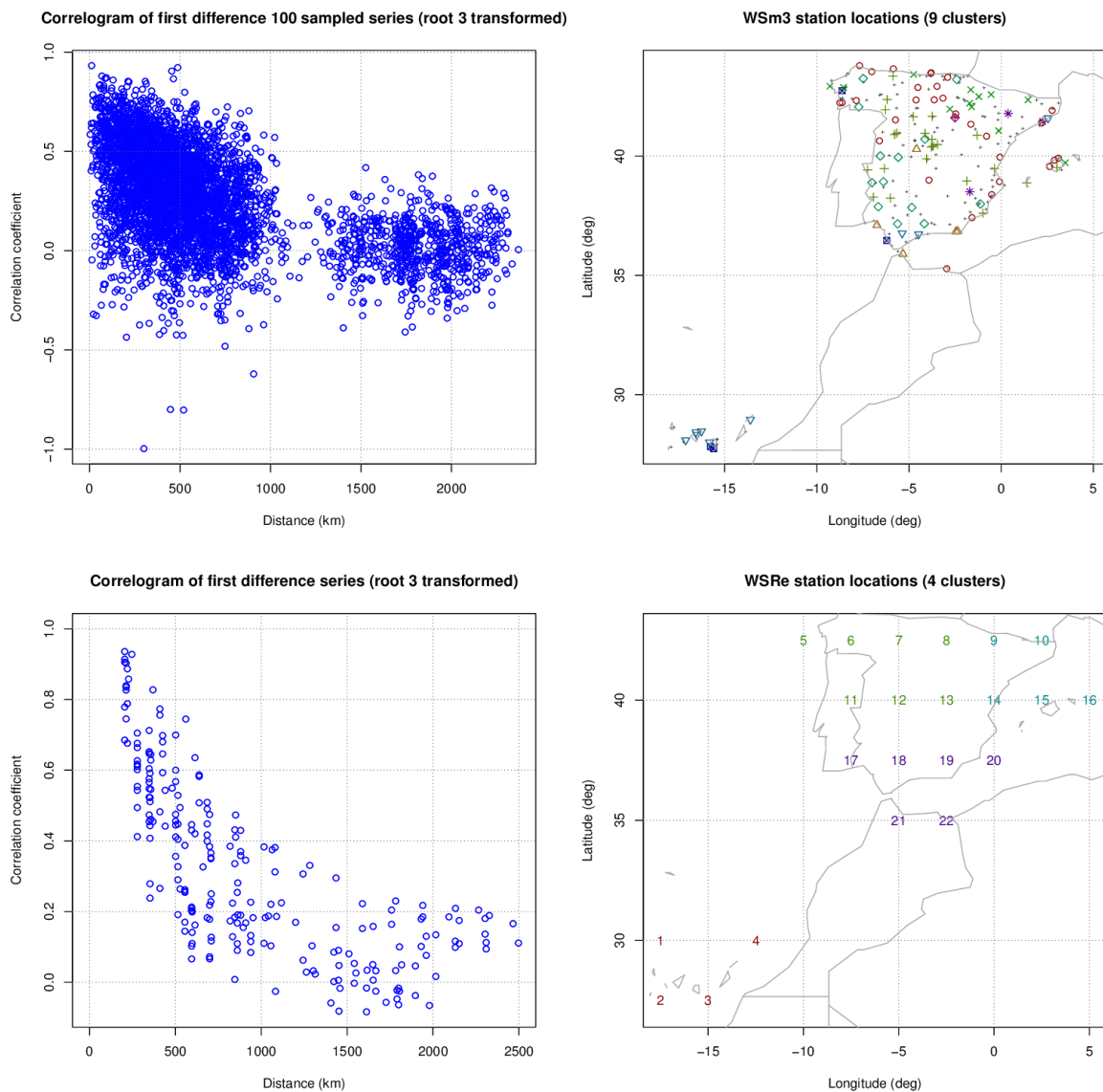


**Fig. 1. Number of average monthly wind speed data available from daily wind runs and from observations at 07, 13 and 18 hours UTC.**

### 3. RESULTS AND DISCUSSION

#### 3.1. Homogenization results

The first exploratory graphics yielded by Climatol show correlograms quickly decaying with distance, resulting in a spatially incoherent distribution of stations clustered according to their inter-correlations (Figure 2, upper row). This points to a high influence of topography and other features of the surroundings of the observatories on their wind measurements, precluding the use of nearby series as the better references. Wang (2008) already noticed that a reference wind speed series built by averaging neighboring stations gave worse results than another of geostrophic wind calculated from homogenized series of pressure. Yet the use of pressure gradients does not account for local thermal winds (see or valley breezes) that may contribute to a high portion of the average wind speed in regions with complex orography and coastal configuration.



**Fig. 2. Correlograms (left) and spatial distribution of clustered stations (right) of observed (up) and reanalysis (down) wind speed series. Cluster analysis was limited to a maximum of 100 stations, the other 133 being represented by dots in the upper right map.**

To account for this local wind circulations, outputs from mesoscale model simulations would be a better reference, as those from MM5 model used by Azorín-Molina *et al.* (2014), although its 10 km resolution is insufficient to capture most small scale thermal winds. Resolutions of 1 km or less would be needed to achieve a full picture of air circulation near the ground, but these simulations are very costly in computer requirements, hindering their use as references for the homogenization of long wind series. Therefore, reanalysis products are a more affordable source of reference series for wind homogenization studies, and the NCEP series gathered here display a better spatial consistency than the observational series (Figure 2, lower row). However, their density is much lower than that of the observational series, and then a direct application of Climatol to a joint (observed plus analyzed) data-set would be using as references more nearby measured series than those more distant of the reanalysis. For these reason, homogenization has been applied separately to each dataset as a first approach in this work.

Results of the different homogenizations performed are summarized in Table 1, containing the number of corrected outliers and breaks, the percentage of inhomogeneous series, the mean RMSE of the data when computed from nearby stations and the mean SNHT of the homogenized series. Both ratio and standardization normalization types gave similar results in the wind run series, with slightly better (lower) values of RMSE and SNHT averages with the ratio normalization, more breaks and less outliers, making this the preferred normalization strategy for this variable. But this is not so clear in the series computed from three hourly observations (WSm3): Mean RMSE is also slightly lower with the ratio normalization (R), but the mean SNHT of the homogenized series is lower with the full standardization of cubic root transformed data (S3r). This is probably due to the higher number of breaks corrected, that could be explained by a lower noise in the series of cubic root transformed data. The number of outliers is also noticeable, more than doubling that of the ratio normalization. Around two thirds of the observational series appear as inhomogeneous, with one or more breaks corrected, while only about one third of the Spanish series analyzed by Azorín-Molina *et al.* (2014) were found inhomogeneous during 1961-2011.

No outlier was detected in the reanalysis series with any of the normalization types, but they are not free from shifts in the mean, with 31 and 36 breaks detected and corrected in the two homogenization processes. Moreover, as there are only 22 series coming from reanalysis, the percentage of inhomogeneous series is far higher than expected: 81.8% with the ratio normalization and 90.9% with the full standardization. Most of the breaks are detected in the second stage of the process, when SNHT is applied to the whole series of anomalies, since only 2 and 3 breaks are detected in the first stage respectively, with the stepped windows SNHT. A possible explanation, to be further investigated, is that the presumed lower noise of the reanalysis series allows the test to achieve significant values that would not be reached in more irregular observational series. As to RMSE and SNHT figures, the full standardization of cubic root values strategy yield better results in this case.



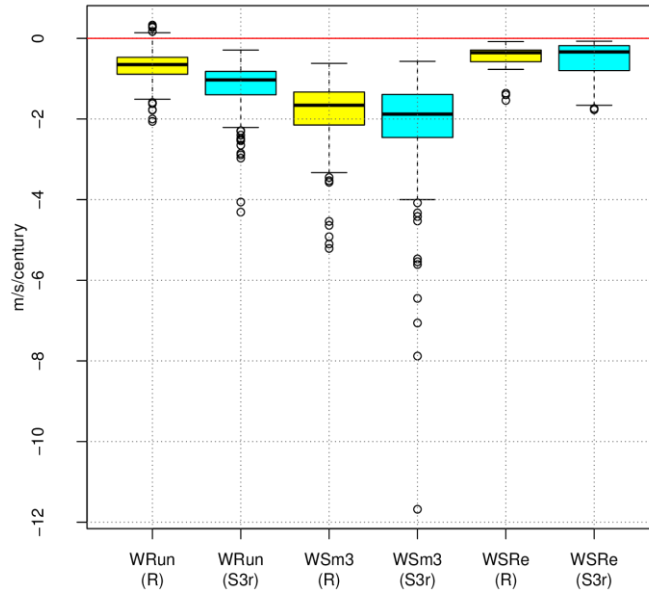
*Table 1.* Outliers and breaks corrected, percentage of inhomogeneous series, mean RMSE of the data when computed from nearby stations and mean SNHT of the resulting homogenized series, for the three data-sets WRun (wind speeds computed from wind daily runs), WSm3 (average wind speed measured three times per day) and WSRe (wind speed from reanalysis). Homogenizations were applied with two different settings: ratio normalization of original data (R) and full standardization of cubic root transformed data (S3r).

	<b>Outliers</b>	<b>Breaks</b>	<b>% Inhom.</b>	<b>Mean RMSE</b>	<b>Mean SNHT</b>
WRun (R)	71	268	64.4	0.38	8.30
WRun (S3r)	75	240	60.1	0.41	9.24
WSm3 (R)	38	360	66.5	0.46	10.64
WSm3 (S3r)	97	409	68.2	0.48	9.50
WSRe (R)	0	31	81.8	0.42	10.2
WSRe (S3r)	0	36	90.9	0.40	8.28

### 3.2. Trends of the homogenized wind speed series

Annual trends computed from the three homogenized monthly wind speed datasets and both methods of normalization are shown in Figure 3, displaying a majority of decreasing values between -0.02 and -2.50 m/s/century. Wind runs present less negative trends than wind observed three times per day, and the ratio normalization also yield less negative trends than the full standardization, which generates some very negative outliers. This fact makes the ratio normalization to be preferred to the standardization of cubic root transformed data, although both gave similar results in terms of RMSE and SNHT of the homogenized series.

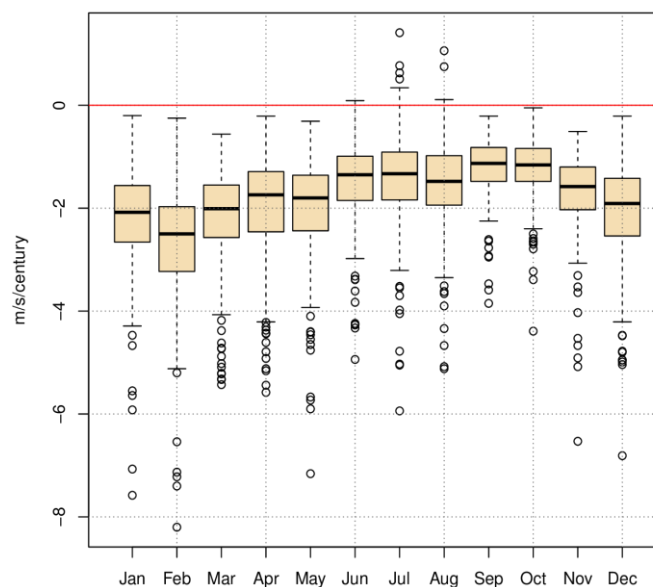
On the other hand, reanalysis series have less negative trends than the observational datasets, backing the hypothesis of the influence of increasing surface roughness on the negative trends of wind speed series observed in many regions (Vautard *et al.*, 2010; Wever, 2012).



**Fig. 3. Annual trends computed from the three homogenized monthly wind speed datasets and both methods of normalization.**

Monthly trends of the three observations per day wind speed monthly averages are presented in Figure 4, showing the higher wind decreases of around -2 m/s/century from November-December until May, while in the warmer months, from June to October, trend values are near -1.5 m/s/century.

This seasonal distribution of trends is in accordance with Azorín-Molina *et al.* (2004) results, although their values were weaker and even positive in summer. But they used a lower number of stations (less than 50 from Spain), did not include the Canary islands, and the period of study was shorter (1961-2011).



**Fig. 4. Monthly trends of the homogenized (with ration normalization method) wind speed averages of three observations per day.**

## 4. CONCLUSIONS AND FUTURE WORK

The Climatol package has allowed an easy homogenization of two datasets of 233 wind speed Spanish series with two different normalization methods.

Wind series appear to be very sensitive to changes and local influences, and are difficult to homogenize, especially in regions with complex orography and coastal configuration, because nearby stations may be poorly correlated.

Most wind speed trends are negative, especially in winter, with typical values between -1 and -2 m/s/century. Trends of reanalysis series are less negative than the observational series, pointing at a possible influence of an increasing roughness in the surroundings of the observatories.

Future work will be devoted to further investigating the benefits of using reanalysis products as a source of reference series to improve the homogenization of the wind speed climatological series, and also to study the geographical distribution of wind speed trends on land and sea, to ascertain the influence of roughness changes on the observed trends.

### Acknowledgements

NCEP Reanalysis data provided by the NOAA/OAR/ESRL PSD, Boulder, Colorado, USA, from their Web site at <http://www.esrl.noaa.gov/psd/>

### References

- Alexandersson H (1986): A homogeneity test applied to precipitation data. *Jour. of Climatol.*, 6, 661-675.
- Azorín-Molina C, Vicente-Serrano SM, McVicar TR, Jerez S, Sánchez-Lorenzo A, López-Moreno JI, Revuelto J, Trigo RM, López-Bustins JA, Espirito-Santo F (2014): Homogenization and Assessment of Observed Near-Surface Wind Speed Trends over Spain and Portugal, 1961-2011. *J. of Climate*, 27:3692-3712.
- Dadaser-Celik F, Cengiz E (2014): Wind speed trends over Turkey from 1975 to 2006. *Int. J. Climatol.*, 34:1913-1927.
- Guijarro JA (2014): User's Guide to Climatol. 40 pp., <http://www.climatol.eu/climatol-guide.pdf>
- Kalnay E, Kanamitsu M, Kistler R, Collins W, Deaven D, Gandin L, Iredell M, Saha S, White G, Woollen J, Zhu W, Chelliah M, Ebisuzaki W, Higgins W, Janowiak J, Mo KC, Ropelewski C, Wang J, Leetmaa A, Reynolds J, Jenne R, Joseph D (1996): The NCEP/NCAR 40-year reanalysis project. *Bull. Amer. Meteor. Soc.*, 77, 437-470.
- Li Z, Yan Z, Tu K, Liu W, Wang Y (2011): Changes in wind speed and extremes in Beijing during 1960-2008 based on homogenized observations. *Advances in Atmospheric Sciences*, 28:408-420.

- McVicar TR, Roderick ML, Donohue RJ, Li LT, Niel TGV, Thomas A, Grieser J, Jhajharia D, Himri Y, Mahowald NM, Mescherskaya AV, Kruger AC, Rehman S, Dinpashoh Y (2012): Global review and synthesis of trends in observed terrestrial near-surface wind speeds: Implications for evaporation. *J. Hydrol.*, 416-417: 182-205.
- Stepanek P (2004): *AnClim: Software for time series analysis and homogenization*. Department of Geography, Faculty of Natural Sciences, Masaryk University.
- Vautard R, Cattiaux Yiou P, Thépaut JN, Ciais P (2010): Northern Hemisphere atmospheric stilling partly attributed to an increase in surface roughness. *Nat. Geosci.*, 3, 756–761, doi:10.1038/ngeo979.
- Wever N (2012): Quantifying trends in surface roughness and the effect on surface wind speed observations. *J. Geophys. Res.*, 117, D11104, doi:10.1029/2011JD017118.
- Wan H, Wang XL, Swail VR (2010): Homogenization and Trend Analysis of Canadian Near-Surface Wind Speeds. *J. of Climate*, 23:1209-1225.
- Wang XL, Feng Y (2007): RHtestV2 user manual. Climate Research Division, Science and Technology Branch, Environment Canada, 19 pp.
- Wang XL (2008): Accounting for Autocorrelation in Detecting Mean Shifts in Climate Data Series Using the Penalized Maximal t or F Test. *Jour. Appl. Meteor. and Climatol.*, 47:2423-2444.

# MATHEMATICAL QUESTIONS OF SPATIAL INTERPOLATION OF CLIMATE VARIABLES

**Tamás Szentimrey, Zita Bihari, Mónika Lakatos**

Hungarian Meteorological Service  
(szentimrey.t@met.hu)

## **Abstract**

We focus on the basic mathematical and theoretical questions of spatial interpolation of meteorological elements. Nowadays in meteorology the most often applied procedures for spatial interpolation are the geostatistical interpolation methods built also in GIS software. The mathematical basis of these methods is the geostatistics that is an exact but special part of the mathematical statistics. However special meteorological spatial interpolation methods for climate variables also can be developed on the basis of the mathematical statistical theory. The main difference between the geostatistical and meteorological interpolation methods can be found in the amount of information used for modeling the necessary statistical parameters. In geostatistics the usable information or the sample for modeling is only the actual predictors, which are a single realization in time. While in meteorology we have spatiotemporal data, namely the long data series which form a sample in time and space as well. The long data series is such a speciality of the meteorology that makes possible to model efficiently the statistical parameters in question. The planned topics to be discussed are as follows.

- Interpolation formulas depending on the spatial probability distribution of climate variables.
- Estimation and modeling of statistical parameters (e.g.: spatial trend, covariance or variogram) for interpolation formulas using spatiotemporal sample and supplementary model variables (topography).
- Use of background information (e.g.: dynamical forecast results, satellite, radar data) for spatial interpolation, data assimilation.
- Creation of gridded climatological databases.

## **1. INTRODUCTION**

First let us consider the abstract schema of the meteorological examinations. The initial stage is the meteorology that means the qualitative formulation of the given problem. The next stage is the mathematics in order to formulate the problem quantitatively. The third stage is to develop software on the basis of the mathematics. Finally the last stage is again the meteorology that is the application of the developed software and evaluation of the obtained results. In the practice however the mathematics is sometimes neglected. Instead of adequate mathematical formulation of the meteorological problem ready-made software are applied to solve the problem. Of course in this case the results are not authentic either.

Concerning our topic we have the following question. What kind of mathematics of spatial interpolation is adequate for meteorology? Nowadays the geostatistical interpolation methods built in GIS software are applied in meteorology. The mathematical basis of these methods is the geostatistics that is an exact but special part of the mathematical statistics. The speciality is connected with the assumption that the data are purely spatial. Consequently, as we see it, the geostatistical methods cannot efficiently use the meteorological data series while the data series make possible to obtain the necessary climate information for the interpolation in meteorology.

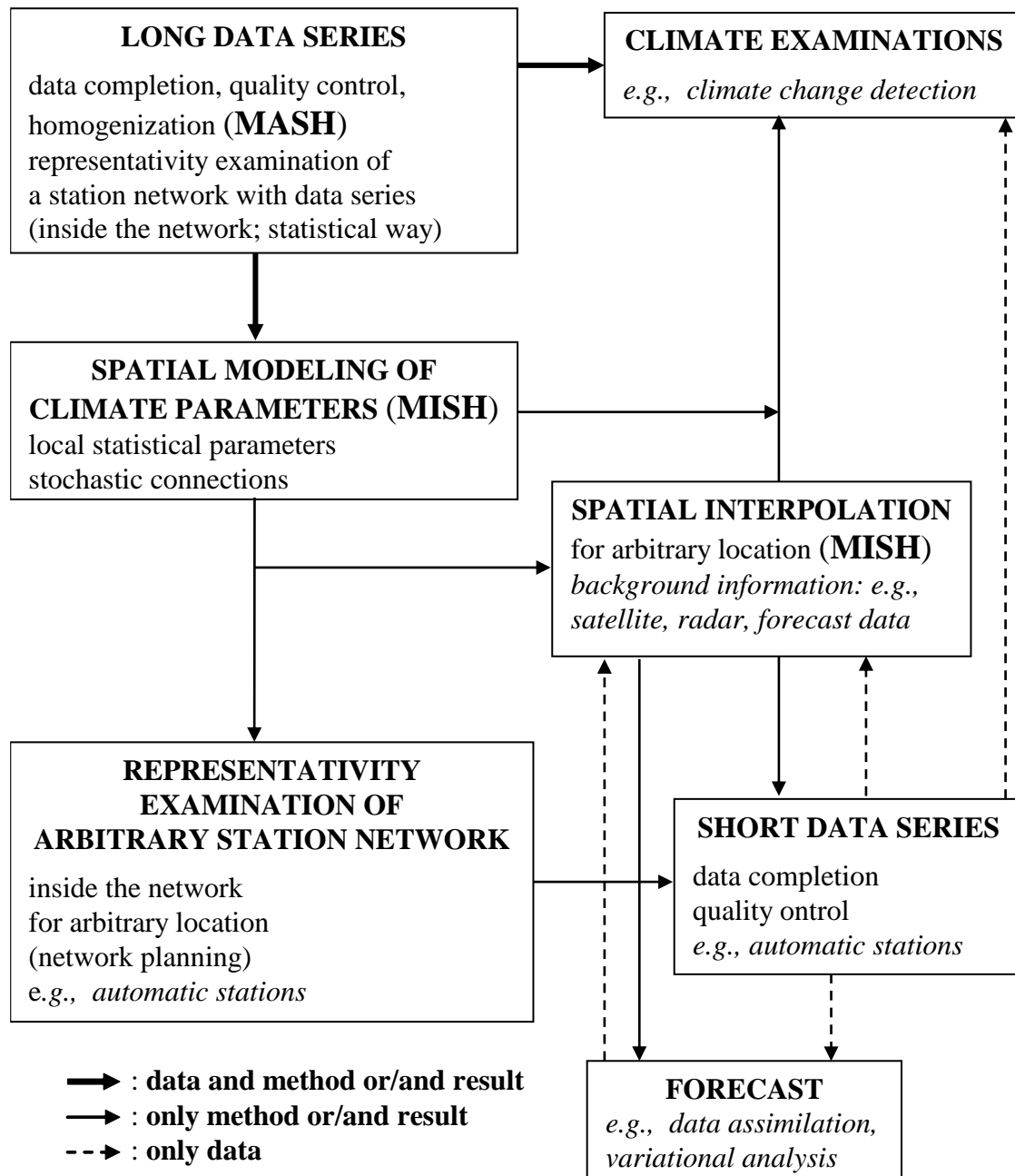


Fig. 1. Block diagram for the possible connection between various basic meteorological topics and systems

Modeling of the climate statistical parameters is a key issue to the interpolation of meteorological elements and that modeling can be based on the long data series. However the data series are usually affected by inhomogeneities (artificial shifts), due to changes in the measurement conditions (relocations, instrumentation) therefore the direct analysis of the raw data series can lead to wrong conclusions. In order to deal with this crucial problem many statistical homogenization procedures have been developed for detection and correction of these inhomogeneities. Similarly to the connection of interpolation and homogenization in our conception the meteorological questions cannot be treated separately. Therefore we present a block diagram (*Fig. 1*) to illustrate the possible connection between various important meteorological topics. The software MASH (Multiple Analysis of Series for Homogenization; *Szentimrey*, 1999, 2014) and MISH (Meteorological Interpolation based on Surface Homogenized Data Basis; *Szentimrey* and *Bihari*, 2007, 2014) were developed by us. These software were applied also in CARPATCLIM project.

## 2. MATHEMATICAL OVERVIEW OF SPATIAL INTERPOLATION PROBLEM IN METEOROLOGY

According to the interpolation problem the unknown predictand  $Z(\mathbf{s}_0, t)$  is estimated by use of the known predictors  $Z(\mathbf{s}_i, t)$  ( $i = 1, \dots, M$ ) where the location vectors  $\mathbf{s}$  are the elements of the given space domain  $D$  and  $t$  is the time.

### 2.1 Additive model of spatial interpolation

The type of the adequate interpolation formula depends on the probability distribution of the meteorological variable. Assuming normal distribution (e.g. temperature) the additive (linear) formula is adequate.

#### 2.1.1 Statistical parameters

In general the interpolation formulas have some unknown interpolation parameters which are known functions of certain statistical parameters. At the additive interpolation formulas the basic statistical parameters can be divided into two groups such as the local and the stochastic parameters. The local parameters are the expected values  $E(Z(\mathbf{s}_i, t))$  ( $i = 0, \dots, M$ ). The stochastic parameters are the covariance or the variogram values belonging to the predictand and predictors such as,

- $\mathbf{c}$ : predictand-predictors covariance vector,
- $\mathbf{C}$ : predictors-predictors covariance matrix,
- $\gamma$ : predictand-predictors variogram vector,
- $\Gamma$ : predictors-predictors variogram matrix.

The covariance is preferred in mathematical statistics and meteorology while the variogram is preferred in geostatistics.

### 2.1.2 Linear meteorological model for expected values

At the statistical modeling of the meteorological elements we have to assume that the expected values of the variables are changing in space and in time alike. The spatial change means that the climate is different in the regions. The temporal change is the result of the possible global climate change. Consequently in case of linear modeling of expected values we assume that

$$E(Z(\mathbf{s}_i, t)) = \mu(t) + E(\mathbf{s}_i) \quad (i = 0, \dots, M) \quad (1)$$

where  $\mu(t)$  is the temporal trend or the climate change signal and  $E(\mathbf{s})$  is the spatial trend.

### 2.1.3 Additive (Linear) Interpolation Formula

Assuming the linear model (1) the appropriate additive meteorological interpolation formula is as follows,

$$\hat{Z}(\mathbf{s}_0, t) = \lambda_0 + \sum_{i=1}^M \lambda_i \cdot Z(\mathbf{s}_i, t)$$

where  $\sum_{i=1}^M \lambda_i = 1$  because of unknown  $\mu(t)$ .

The optimal interpolation parameters  $\lambda_0, \lambda_i$  ( $i = 1, \dots, M$ ) minimize the root-mean-square error,  $RMSE = \sqrt{E\left(\left(Z(\mathbf{s}_0, t) - \hat{Z}(\mathbf{s}_0, t)\right)^2\right)}$ .

These optimal parameters are known functions of statistical parameters!

The optimal constant term is:  $\lambda_0 = \sum_{i=1}^M \lambda_i (E(\mathbf{s}_0) - E(\mathbf{s}_i))$

The vector of weighting factors  $\boldsymbol{\lambda}^T = [\lambda_1, \dots, \lambda_M]$  can be written in covariance form

$$\boldsymbol{\lambda}^T = \left( \mathbf{c}^T + \mathbf{1}^T \frac{(\mathbf{1} - \mathbf{1}^T \mathbf{C}^{-1} \mathbf{c})}{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1}} \right) \mathbf{C}^{-1},$$

or equivalently in variogram form

$$\boldsymbol{\lambda}^T = \left( \boldsymbol{\gamma}^T + \mathbf{1}^T \frac{(\mathbf{1} - \mathbf{1}^T \boldsymbol{\Gamma}^{-1} \boldsymbol{\gamma})}{\mathbf{1}^T \boldsymbol{\Gamma}^{-1} \mathbf{1}} \right) \boldsymbol{\Gamma}^{-1}.$$

Consequently the unknown statistical parameters are the spatial trend differences  $E(\mathbf{s}_0) - E(\mathbf{s}_i)$  ( $i = 1, \dots, M$ ) and the covariances  $\mathbf{c}, \mathbf{C}$ . In essence these parameters are climate parameters which in fact means we could interpolate optimally if we knew the climate.



### Remark

Unfortunately inadequate formulas are often applied in the practice:

- Inverse Distance Weighting (IDW):  
 $\lambda_0 = 0$  that is excluding spatial trend, and  $\lambda_i$  ( $i = 1, \dots, M$ ) based on distances are not optimal weighting factors.
- Ordinary kriging:  $\lambda_0 = 0$  excludes the spatial trend.

Adequate formulas are in meteorology:

- Universal kriging formula,
- Regression (residual, detrended) kriging formula.

But in geostatistics modeling of statistical parameters is based on only the actual predictors.

### 2.1.4 Possibility for modeling of unknown statistical parameters in Meteorology

The special possibility in meteorology is to use the long meteorological data series for modeling of the climate statistical parameters in question. The data series make possible to know the climate in accordance with the fundament of statistical climatology!

The main difference between geostatistics and meteorology can be found in the amount of information being usable for modeling the statistical parameters. In geostatistics the usable information or the sample for modeling is only the actual predictors  $Z(\mathbf{s}_i, t)$  ( $i = 1, \dots, M$ ) which belong to a fixed instant of time, that is a single realization in time. While in meteorology we have spatiotemporal data, namely the long data series which form a sample in time and space as well and make possible to model the climate statistical parameters in question. If the meteorological stations  $\mathbf{S}_k$  ( $k = 1, \dots, K$ ) ( $\mathbf{S} \in D$ ) have long data series then spatial trend differences  $E(\mathbf{S}_k) - E(\mathbf{S}_l)$  ( $k, l = 1, \dots, K$ ) as well as the covariances  $\text{cov}(Z(\mathbf{S}_k), Z(\mathbf{S}_l))$  ( $k, l = 1, \dots, K$ ) can be estimated statistically. Consequently these parameters are essentially known and provide much more information for modeling than the predictors  $Z(\mathbf{s}_i, t)$  ( $i = 1, \dots, M$ ) only. However nowadays unfortunately the geostatistical interpolation methods built in GIS software are applied in meteorology mostly.

### 2.2 Multiplicative model of spatial interpolation

In this paper only the linear or additive model was described in detail which is appropriate in case of normal probability distribution. However perhaps it is worthwhile to remark that for case of a quasi lognormal distribution (e.g. precipitation sum) we deduced a mixed additive multiplicative formula which is used also in our MISH system and it can be written in the following form,

$$\hat{Z}(\mathbf{s}_0, t) = \mathcal{G} \cdot \left( \prod_{q_i \cdot Z(\mathbf{s}_i, t) \geq \mathcal{G}} \left( \frac{q_i \cdot Z(\mathbf{s}_i, t)}{\mathcal{G}} \right)^{\lambda_i} \right) \cdot \left( \sum_{q_i \cdot Z(\mathbf{s}_i, t) \geq \mathcal{G}} \lambda_i + \sum_{q_i \cdot Z(\mathbf{s}_i, t) < \mathcal{G}} \lambda_i \cdot \left( \frac{q_i \cdot Z(\mathbf{s}_i, t)}{\mathcal{G}} \right) \right)$$

where the interpolation parameters are  $\mathcal{G} > 0$ ,  $q_i > 0$ ,  $\lambda_i \geq 0$  ( $i = 1, \dots, M$ ) and  $\sum_{i=1}^M \lambda_i = 1$ .

### 3. INTERPOLATION WITH BACKGROUND INFORMATION

The background information e.g. forecast, satellite, radar data can be efficiently used to decrease the interpolation error. In this paper only the interpolation based on additive model or normal distribution is presented.

According to the section 2.1.3 let us assume that,

$Z(\mathbf{s}_0, t)$ : predictand,

$$\hat{Z}(\mathbf{s}_0, t) = \lambda_0 + \sum_{i=1}^M \lambda_i Z(\mathbf{s}_i, t): \text{interpolated predictand,}$$

moreover there is given,

$\mathbf{G} = \{G(\mathbf{s}, t) | \mathbf{s} \in D\}$  : background information on a dense grid.

#### 3.1 The principle of interpolation with background information

The interpolated predictand given  $\mathbf{G}$  can be expressed as,

$$\hat{Z}_{\mathbf{G}}(\mathbf{s}_0, t) = \hat{Z}(\mathbf{s}_0, t) + \mathbb{E}\left(Z(\mathbf{s}_0, t) - \hat{Z}(\mathbf{s}_0, t) | \mathbf{G}\right)$$

where  $\mathbb{E}\left(Z(\mathbf{s}_0, t) - \hat{Z}(\mathbf{s}_0, t) | \mathbf{G}\right)$  is the conditional expectation of  $Z(\mathbf{s}_0, t) - \hat{Z}(\mathbf{s}_0, t)$ , given  $\mathbf{G}$ .

#### 3.2 Reanalysis data, Data Assimilation

The reanalysis data are based on the data assimilation which procedure is in strong relationship with the methodology of interpolation with background information. The Bayes estimation theory is the mathematical background of the data assimilation and the following variational cost function has to be minimized in order to estimate the analysis field,

$$J(\mathbf{z}) = (\mathbf{z} - \mathbf{g})^T \mathbf{Q}^{-1} (\mathbf{z} - \mathbf{g}) + (\mathbf{y}_0 - \mathbf{Fz})^T \mathbf{P}^{-1} (\mathbf{y}_0 - \mathbf{Fz}) \quad (2)$$

$\mathbf{z}$  : analysis field, predictand (grid),

$\mathbf{g}$  : given background field (forecast),

$\mathbf{y}_0$  : given observations, predictors;  $\mathbf{Fz} = \mathbb{E}(\mathbf{y}_0 | \mathbf{z})$ ,

$\mathbf{Q}$  : background error covariance matrix,

$\mathbf{P}$  : observation error covariance matrix.

It can be proved that this procedure is essentially an interpolation with background information including a quality control part for the predictors.

However there are several problems with the reanalysis data in the practice:

- i, Inhomogeneous predictor station data series are used.
- ii, Few stations are used with little spatial representativity.
- iii, There are also some problems with the data assimilation formula (2):
  - Lack of good climate statistical parameters in matrix  $\mathbf{Q}$ .
  - Formula (2) includes an implicit assumption of  $E(\mathbf{z}|\mathbf{g}) = \mathbf{g}$ .

## 4. GRIDDED DATABASES

We emphasize the importance of gridded databases based on observations with good quality! For example there is the CARPATCLIM project implemented during the last years. The main aim of this project was to produce gridded climatological database for the Carpathian Region using unified methods. The grids cover the area between latitudes 44°N and 50°N, and longitudes 17°E and 27°E. Daily values of more than ten meteorological variables were calculated on a 0.1° spatial resolution grid for the period 1961-2010. Climate statistics (monthly and annual values) and different climate indices were also determined from the daily grids. For ensuring the usage of largest possible station density the necessary work phases were implemented on national level but by the same methods and software. The commonly used methods and software were the method MASH (Multiple Analysis of Series for Homogenization) for homogenization, quality control, completion of the observed daily data series; and the method MISH (Meteorological Interpolation based on Surface Homogenized Data Basis) for gridding of homogenized daily data series. Besides the common software, the harmonization of the results across country borders was promoted also by near border data exchange.

CARPATCLIM homepage: <http://www.carpatclim-eu.org/pages/home/>

### 4.1. Software MISH

Our method MISH (Meteorological Interpolation based on Surface Homogenized Data Basis) for the spatial interpolation of surface meteorological elements was developed (*Szentimrey and Bihari, 2007, 2014*) according to the mathematical background that is outlined in Sections 2, 3. This is a meteorological system not only in respect of the aim but in respect of the tools as well. It means that using all the valuable meteorological information – e.g. climate and possible background information – is required.

The last software version MISHv1.03 consists of two units that are the modeling and the interpolation systems. The interpolation system can be operated on the results of the modeling system. We summarize briefly the most important facts about these two units of the developed software.

#### **Modeling subsystem for climate statistical (local and stochastic) parameters:**

- Based on long homogenized data series and supplementary deterministic model variables. The model variables may be such as height, topography, distance from the sea etc.. Neighbourhood modeling, correlation model for each grid point.

- Benchmark study, cross-validation test for interpolation error or representativity.
- Modeling procedure must be executed only once before the interpolation applications!

**Interpolation subsystem:**

- Additive (e.g. temperature) or multiplicative (e.g. precipitation) model and interpolation formula can be used depending on the climate elements.
- Daily, monthly values and many years' means can be interpolated.
- Few predictors are also sufficient for the interpolation and no problem if the greater part of daily precipitation predictors is equal to 0.
- The interpolation error or representativity is modeled too.
- Capability for application of supplementary background information (stochastic variables) e.g. satellite, radar, forecast data.
- Data series completion that is missing value interpolation, completion for monthly or daily station data series.
- Interpolation, gridding of monthly or daily station data series for given predictand locations. In case of gridding the predictand locations are the nodes of a relatively dense grid.

Our MISH-MASH software can be downloaded from:

[http://www.met.hu/en/omsz/rendezvenyek/homogenizationand\\_interpolation/software/](http://www.met.hu/en/omsz/rendezvenyek/homogenizationand_interpolation/software/)

**References**

Szentimrey, T., 1999: Multiple Analysis of Series for Homogenization (MASH), Proceedings of the Second Seminar for Homogenization of Surface Climatological Data, Budapest, Hungary; WMO, WCDMP-No. 41, pp. 27-46.

Szentimrey, T., Bihari, Z., 2007: Mathematical background of the spatial interpolation methods and the software MISH (Meteorological Interpolation based on Surface Homogenized Data Basis), Proceedings of the Conference on Spatial Interpolation in Climatology and Meteorology, Budapest, Hungary, 2004, COST Action 719, COST Office, 2007, pp. 17-27

Szentimrey, T, Bihari, Z., Lakatos, M., Szalai,S., 2011: Mathematical, methodological questions concerning the spatial interpolation of climate elements. Proceedings of the Second Conference on Spatial Interpolation in Climatology and Meteorology, Budapest, Hungary, 2009, Időjárás 115, 1-2, 1-11

Szentimrey, T. 2013: Theoretical questions of daily data homogenization, Időjárás Vol. 117. No. 1, January-March 2013. pp. 113-122.

Szentimrey,T., Bihari, Z., 2014: Manual of interpolation software MISHv1.03, Hungarian Meteorological Service, p. 60.

Szentimrey,T.,2014: Manual of homogenization software MASHv3.03, Hungarian Meteorological Service, p. 69.

# **PRACTICAL ASPECTS OF RAW, HOMOGENIZED AND GRIDDED DAILY PRECIPITATION DATASETS**

**Predrag Petrović, Gordana Simić, Ivana Kordić**

Republic Hydrometeorological Service of Serbia, Kneza Višeslava 66, 11030 Belgrade, Serbia, tel. +381/11/3050-803, +381/11/3050-802  
(predrag.petrovic@hidmet.gov.rs, gordana.simic@hidmet.gov.rs, ivana.kordic@hidmet.gov.rs)

## **THE USE OF DAILY DATA IN CLIMATE ASSESSMENTS**

In order to take advantage of the best possible temporal resolution of climate data, there is an increased need for including daily series into climate assessments. However, raw (observed) data are not always convenient due to data gaps and inhomogeneities.

Construction of quality controlled and homogenized datasets with gaps filled in was the first step in solving this problem. Still, uneven spatial distribution of measurement sites could not represent spatial distribution of values of a climate element adequately. Thus, spatial interpolation techniques had to be applied as necessary step in obtaining the best possible climate maps.

### **Types of available daily data**

Hence, there are three types of climate data:

- Raw (observed) data feature gaps and inhomogeneities in series;
- Homogenized data, where the gaps are filled in and inhomogeneities are mostly eliminated, but series are with "spotty" data, not representative for any wider area;
- Gridded (spatially interpolated) data, where any wider area is uniformly covered with data.

Both homogenization and spatial interpolation procedures are designed to obtain the best possible output from the climate data, without gaps, artificial changes in measurements and various data density. Still, every procedure involve techniques that modify the data.

Although every type of data has its own advantages, results of such climate assessments depend upon the choice of type of data.

## **PRACTICAL USE OF DAILY PRECIPITATION DATA IN DIFFERENT STUDIES**

In order to examine to what extent this choice might affect the results, a comparison of data processing products from daily precipitation series has been made between the raw, homogenized and gridded datasets.

Raw (observed) data are the primary source of data for climate assessments. These data are in daily or sub-daily temporal resolution, depending on the observation schedule of a weather station. Climate assessments require calculations of various values of lower temporal resolution (i.e. daily from sub-daily, monthly, annual resolution). These are rather simple calculations of average, maximum, minimum values, counts and standard deviations. Calculated values are further used in an assessment.

Naturally, homogenized and gridded data might also be used the same way as raw data, deriving the same products of calculations. However, homogenization techniques might modify daily data on a basis of calculated values of lower temporal resolution and/or other statistical tools.

Since precipitation data have the greatest spatial and temporal variability, modifications of their daily values might produce the most obvious modifications in products that come from calculated values.

Two types of data processing products from daily precipitation data were examined: climate indices and extreme daily precipitation.

## Data and methods

In its various stages, project CARPATCLIM has dealt with the three types of data. Therefore, raw and homogenized daily precipitation data from 73 stations in Serbia are used. Data referring to the grid point nearest to the measurement sites are used for gridded data series. All series involve the 50-year period, 1961-2010.

Homogenization and filling in the data gaps is performed using MASH method and software (version v3.03), where 11 raw data series from neighboring countries and 38 raw data series from central Serbia, all within the 50 km distance from the territory of Serbia north of 44 N, are included as bordering territory. Spatial interpolation is performed using MISH method and software.



**Fig. 1. Network of precipitation stations**

## CALCULATION OF CLIMATE INDICES

Climate indices are simply calculated values used for climate change assessment. They are calculated on daily data basis and, as a result, annual index values are returned. There are 27 core indices, derived from daily maximum and minimum temperature and daily precipitation data. The definitions for a core set of 27 descriptive indices of extremes are defined by the Joint CCI/CLIVAR/JCOMM Expert Team on Climate Change Detection and Indices (ETCCDMI) and they are recommended by WMO.

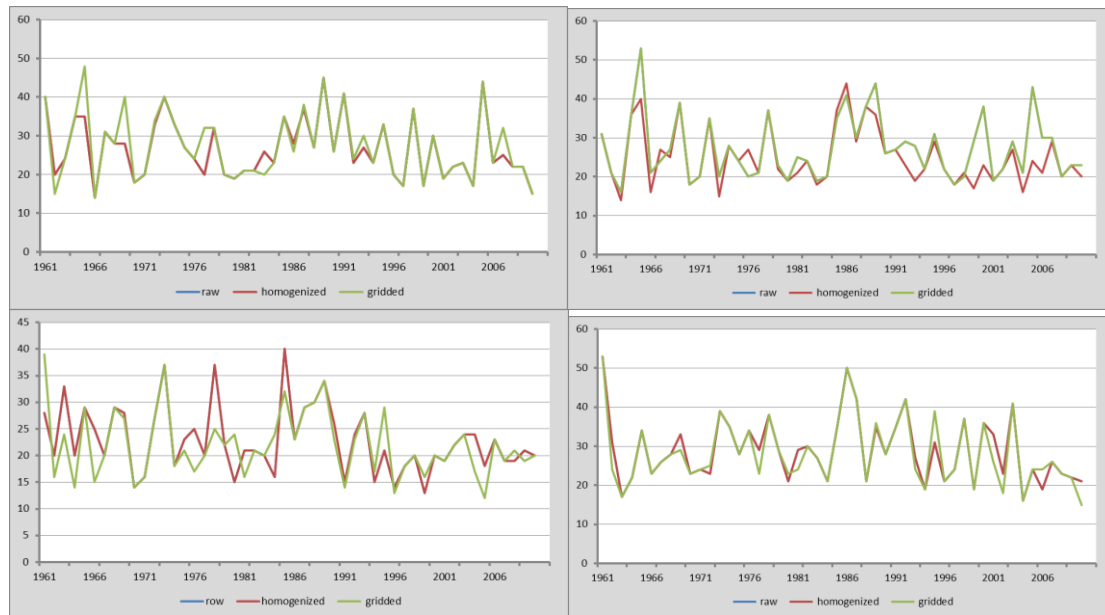
A subset of 11 indices processed in this survey are derived from daily precipitation data only (Table 1).

*Table 1.* A subset of CLIVAR/ETCCDMI indices derived from daily precipitation data only

<b>Index</b>	<b>Index full name</b>	<b>Definition</b>	<b>Unit</b>
CDD	Consecutive dry days	Maximum number of consecutive dry days with RR<1mm	days
CWD	Consecutive wet days	Maximum number of consecutive wet days with RR>1mm	days
PRCPTOT	Annual total wet-day precipitation	Annual total PRCP in wet days (RR>=1mm)	mm
R10	Number of heavy precipitation days	Annual count of days when PRCP>10mm	mm
R20	Number of very heavy precipitation days	Annual count of days when PRCP>20mm	mm
R25	Number of days above 25mm	Annual count of days when PRCP>25mm	mm
R95p	Very wet days	Annual total PRCP when RR>95th percentile	mm
R99p	Extremely wet days	Annual total PRCP when RR>99th percentile	mm
Rx1day	Max 1-day precipitation amount	Monthly max 1-day precipitation	mm
Rx5day	Max 5-day precipitation amount	Monthly max 5-day consecutive precipitation	mm
SDII	Simple daily intensity index	Annual total precipitation divided by the number of wet days (defined by PRCP >=1mm) in the year	mm/day

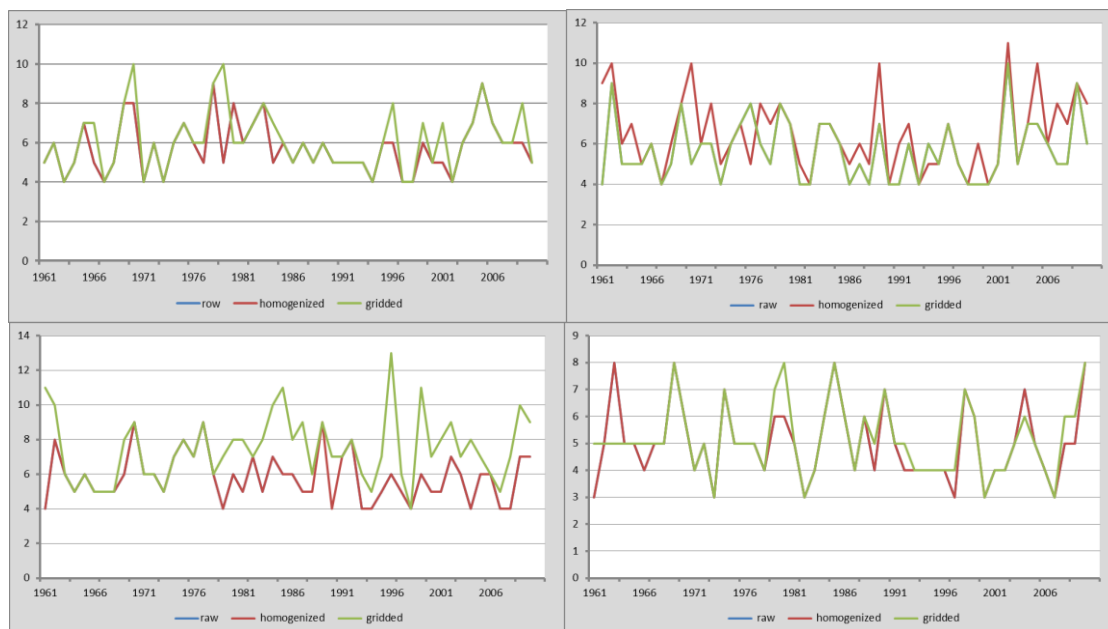
Indices are calculated for the three types of data from all series. The results are compared between the types of data. In order to present spatial pattern of the results, a comparison of series is given for four stations from different parts of northern Serbia: Beograd (central part), Kragujevac (southern part), Valjevo (western part) and Palić (northern part).

**Consecutive dry days (CDD)** practically match between raw and homogenized series, while gridded series are slightly different. This difference might come from small changes in values around the threshold of 1 mm within gridding procedure. However, trends in values are mostly unchanged.



**Fig. 2. Consecutive dry days (CDD) calculated from raw, homogenized and gridded datasets for Beograd (upper left), Kragujevac (upper right), Valjevo (lower left) and Palić (lower right), 1961-2010**

**Consecutive wet days (CWD)** also have a good match between raw and homogenized series, while gridded series still differ slightly from the two other series. Since the same threshold is used as the criterion for determining values of this index and CDD, the problem of small changes of values around 1 mm might also be attributed to CWD.

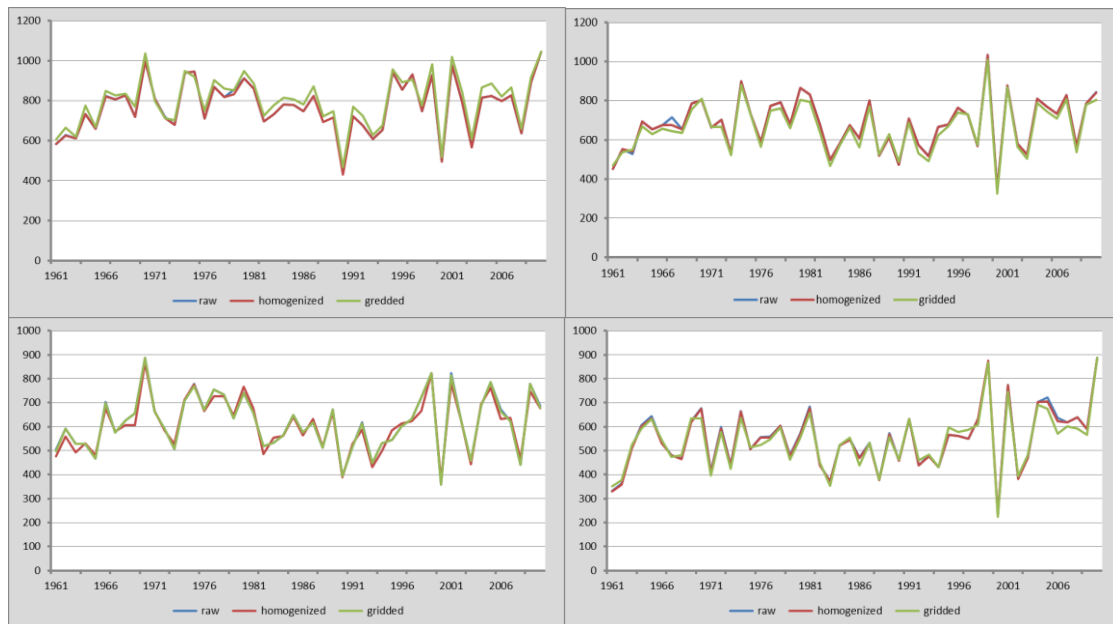


**Fig. 3. Consecutive wet days (CWD) calculated from raw, homogenized and gridded datasets for Beograd (upper left), Kragujevac (upper right), Valjevo (lower left) and Palić (lower right), 1961-2010**

**Annual total wet day precipitation (PRCPTOT)** has almost the same values in raw and homogenized series, while gridded series are slightly different. The match between raw and homogenized series practically show the true effect of homogenization procedure, since the

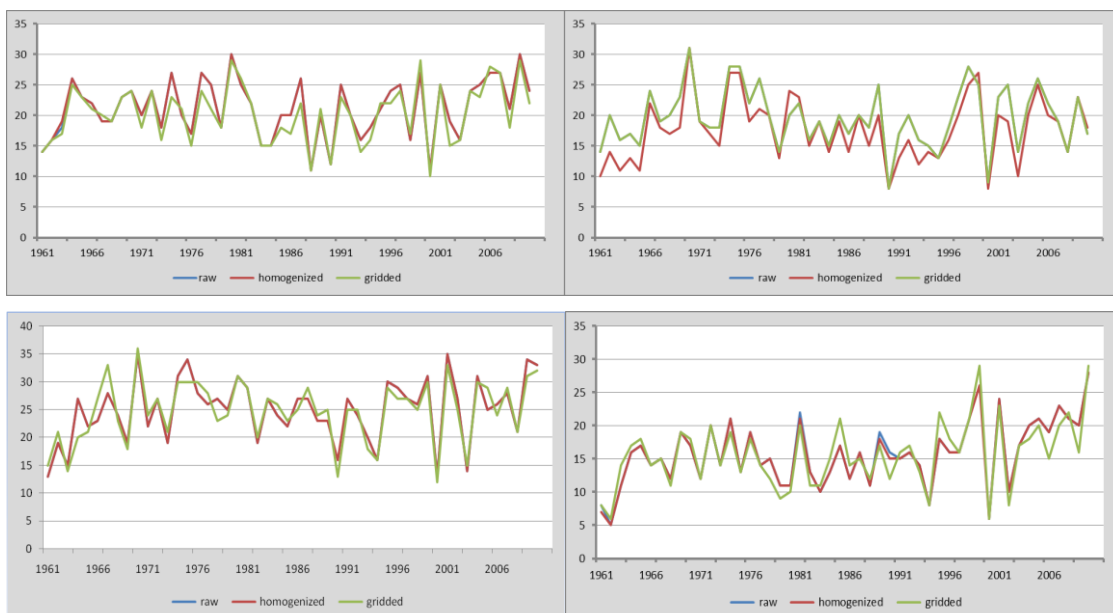


MASH method is based on calculations on monthly temporal resolution. Here, homogenization has performed well, correcting suspicious values only.



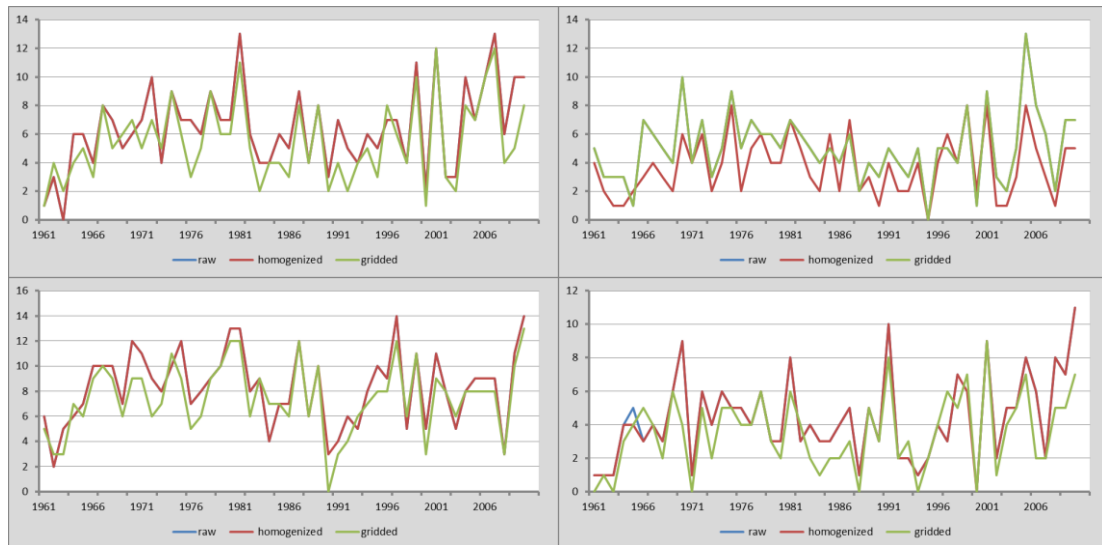
**Fig. 4. Annual total wet day precipitation (PRCPTOT) calculated from raw, homogenized and gridded datasets for Beograd (upper left), Kragujevac (upper right), Valjevo (lower left) and Palić (lower right), 1961-2010**

**Number of heavy precipitation days (R10)** matches quite well between all series, except in Kragujevac. Gridding must have been under influence of sparse station network in the southern part of examined part of Serbia, producing differences in values around 10 mm or more.



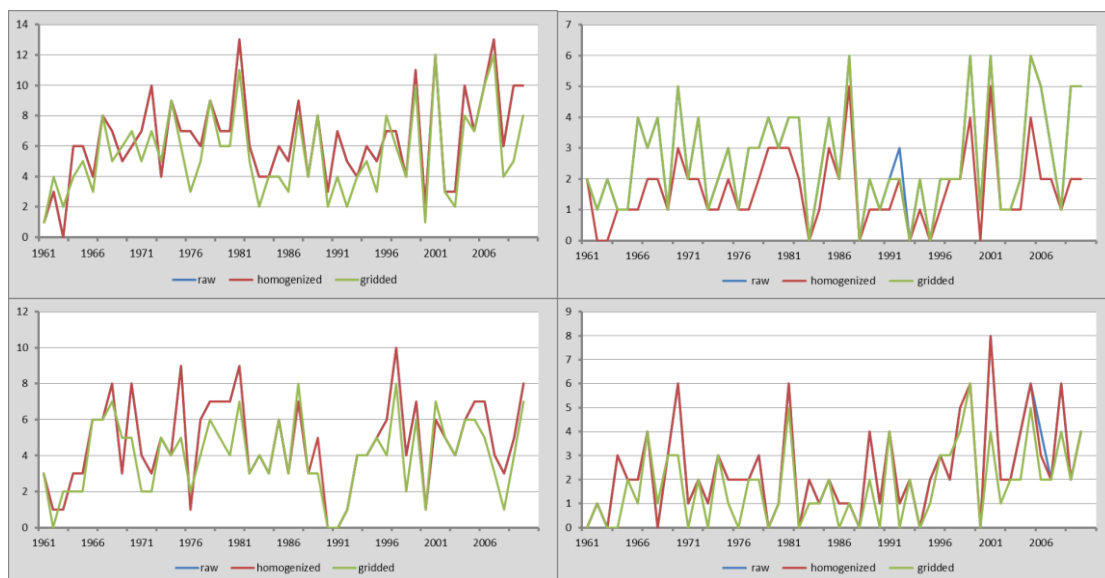
**Fig. 5. Number of heavy precipitation days (R10) calculated from raw, homogenized and gridded datasets for Beograd (upper left), Kragujevac (upper right), Valjevo (lower left) and Palić (lower right), 1961-2010**

**Number of very heavy precipitation days (R20)** experiences practically the same problem as with R10. However, due to smaller values of R20, these differences between gridded and other series are relatively more significant.



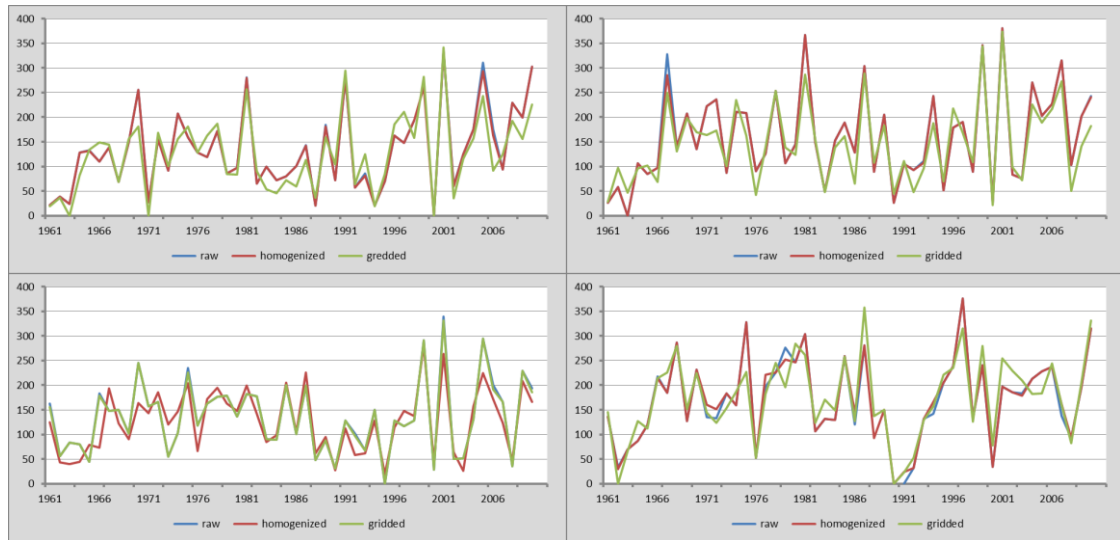
**Fig. 6. Number of very heavy precipitation days (R20) calculated from raw, homogenized and gridded datasets for Beograd (upper left), Kragujevac (upper right), Valjevo (lower left) and Palić (lower right), 1961-2010**

**Number of days above 25mm (R25)** behaves the same way as R20, but with a little more emphasis on differences between gridded and other series, relatively significant due to low index values.



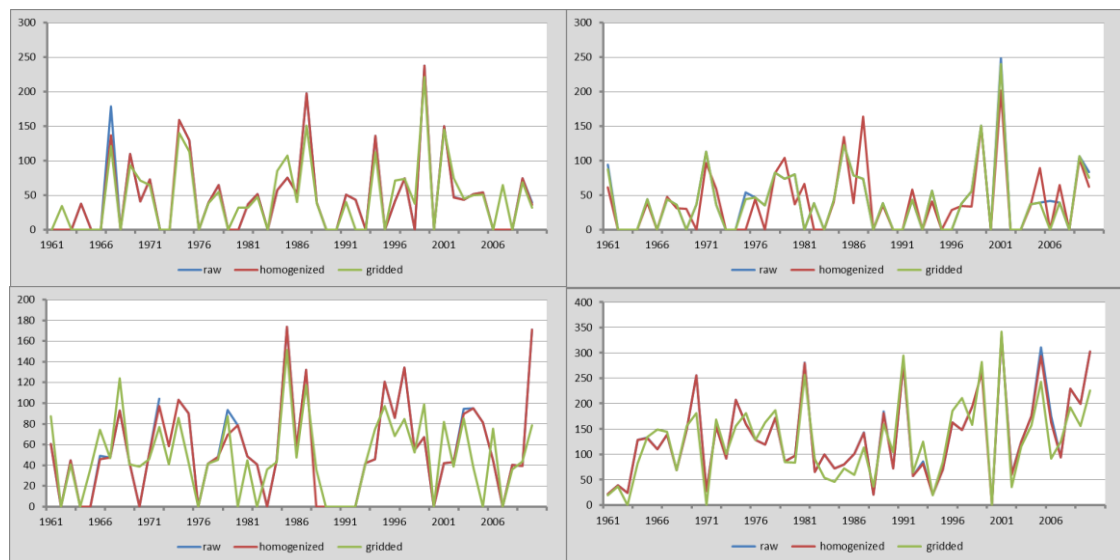
**Fig. 7. Number of days above 25mm (R25), calculated from raw, homogenized and gridded datasets for Beograd (upper left), Kragujevac (upper right), Valjevo (lower left) and Palić (lower right), 1961-2010**

**Very wet days (R95p)** show seldom differences between raw and homogenized series. Most probably, these differences are due to elimination of outliers and correction of suspiciously high values in homogenization procedure. Gridded data also differ from two other series, but in more or less the same magnitude as in PRCPTOT.



**Fig. 8. Very wet days (R95p), calculated from raw, homogenized and gridded datasets for Beograd (upper left), Kragujevac (upper right), Valjevo (lower left) and Palić (lower right), 1961-2010**

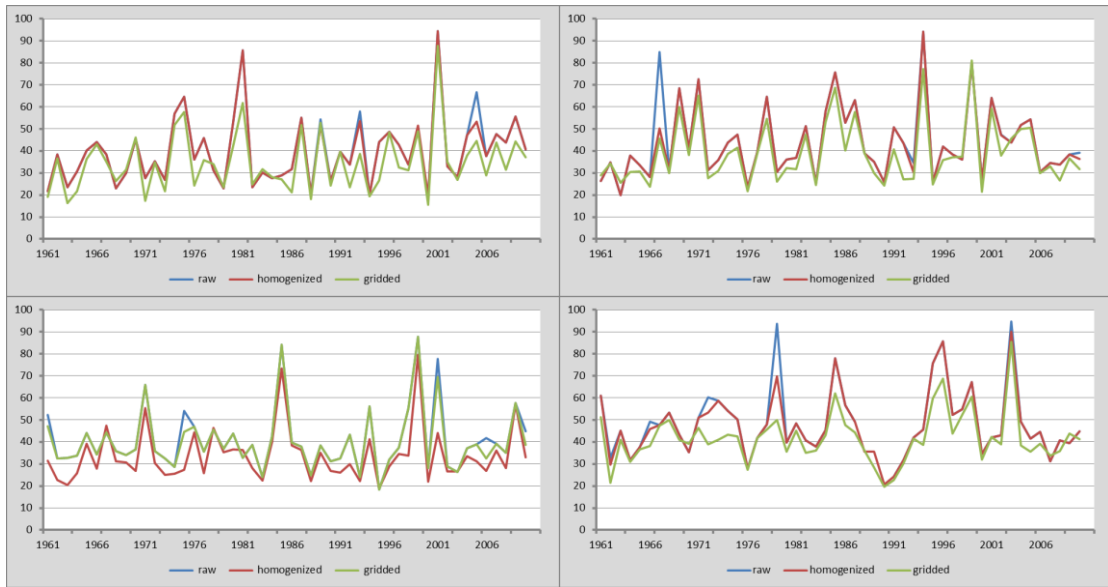
**Extremely wet days (R99p)** are similar to R95p, showing the same differences of a similar magnitude. Outliers from raw series are still visible through mismatches with homogenized series.



**Fig. 9. Extremely wet days (R99p), calculated from raw, homogenized and gridded datasets for Beograd (upper left), Kragujevac (upper right), Valjevo (lower left) and Palić (lower right), 1961-2010**

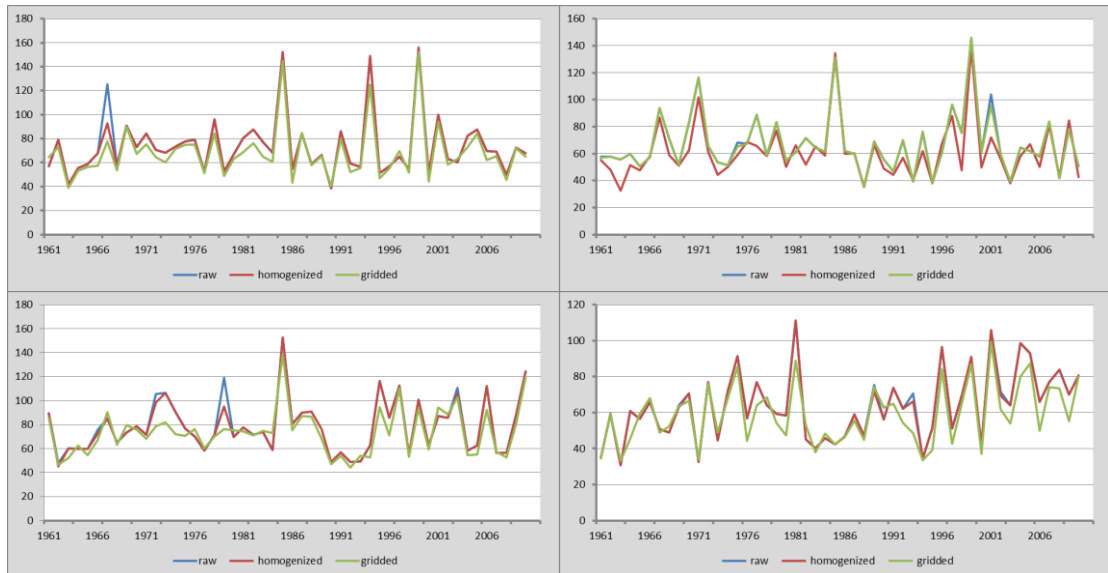
**Max 1-day precipitation amount (Rx1day)** has clear differences between raw and homogenized series at outlier data. However, these differences do not change index trends

significantly, for very few points being replaced by lower values. Gridded series produce similar or lower data values.



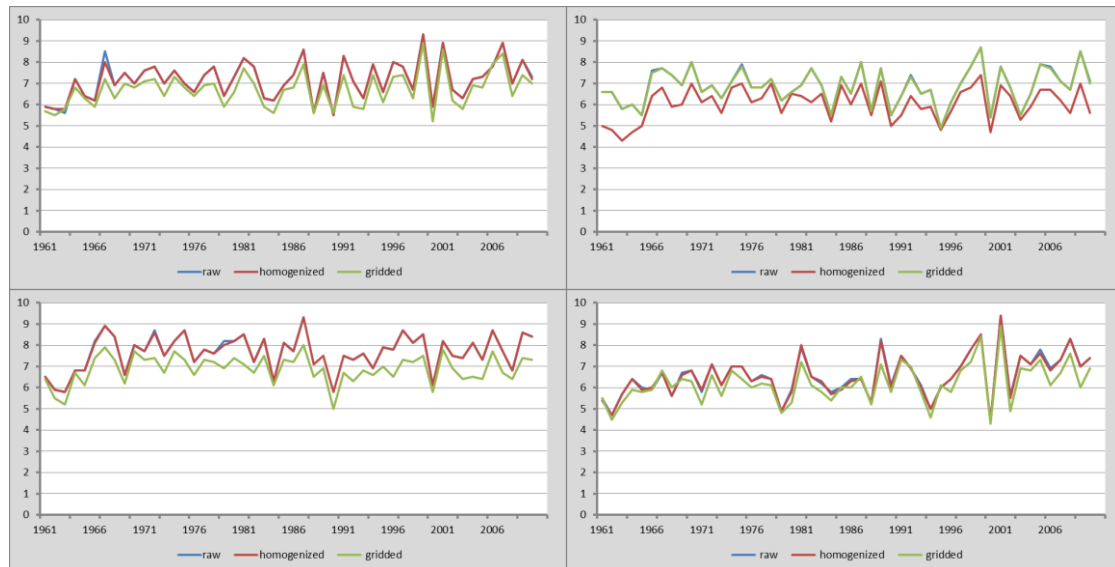
**Fig. 10. Max 1-day precipitation amount (Rx1day), calculated from raw, homogenized and gridded datasets for Beograd (upper left), Kragujevac (upper right), Valjevo (lower left) and Palić (lower right), 1961-2010**

**Max 5-day precipitation amount (Rx5day)** has significantly lower differences between the series than Rx1day. This comes from the fact that Rx1day is only a single selected value, while Rx5day index values include other, lower values than maximum. Still, descriptions of changes are the same as Rx1day, since both indices deal with extreme precipitation values.



**Fig. 11. Max 5-day precipitation amount (Rx5day), calculated from raw, homogenized and gridded datasets for Beograd (upper left), Kragujevac (upper right), Valjevo (lower left) and Palić (lower right), 1961-2010**

**Simple daily intensity index (SDII)** is the most complex index of the eleven selected indices, taking both values and number of days into account. Therefore, differences between the series are not necessarily in favor of any series. While Kragujevac has values from gridded series higher than from raw or homogenized series, Valjevo presents the opposite case. Nevertheless, raw and homogenized series show very good match.



**Fig. 12. Simple daily intensity index (SDII), calculated from raw, homogenized and gridded datasets for Beograd (upper left), Kragujevac (upper right), Valjevo (lower left) and Palić (lower right), 1961-2010**

As could be seen from the presented results, all three types of data can be used for calculation of climate indices. Raw and homogenized series match in most cases, while series from gridded data are different from the original series. Larger differences occur with indices that are derived from very high and extreme values, which is due to smoothing effect of spatial interpolation of data.

Daily precipitation data are amongst the most variable climate parameters both in time and in space. Very high and extreme values might cover a wide range of values that rarely and sparse occur almost independently. Thus, it is a tricky task to confirm and verify such values. Interpolated data always intercept between two or more source values, and thus never bring an extreme value as a result. This is the main cause of the smoothing effect of spatial interpolation that cut down most of very high and especially extreme values.

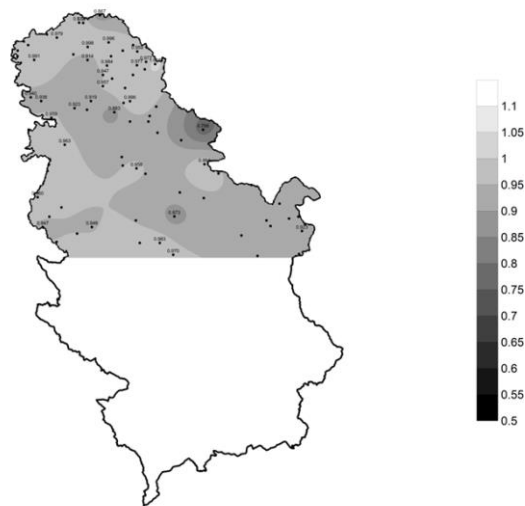
Bearing all these findings in mind, it is recommended to use raw and quality controlled or homogenized data rather than spatially interpolated data for calculation of climate extremes. Spatial interpolation does not support preservation of any spatial pattern that come from raw or homogenized values.

## **CALCULATION OF EXTREME VALUES FOR A RETURN PERIOD**

Extreme precipitation values are important not only in climatology. These values are amongst the main inputs in civil engineering for designing buildings and other objects. They also

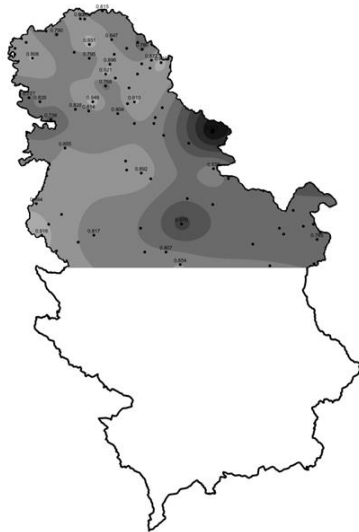
present one of the main input data for hydrological forecasts. Since extreme precipitation values are of a crucial interest in many fields, a special care has to be taken to perform their calculation. One of the most common approaches is to calculate maximum expected precipitation value for a return period.

Calculation of extreme daily precipitation is performed using Gumbel method via extRemes software (R-platform). The return period was set to 100 years. Since both software and method require the complete series of 50 values (one value for each year of the series), series with gaps of at least one year are discarded from this survey. Thus, a comparison of the results could be performed for 32 remaining stations with no gaps in raw series.

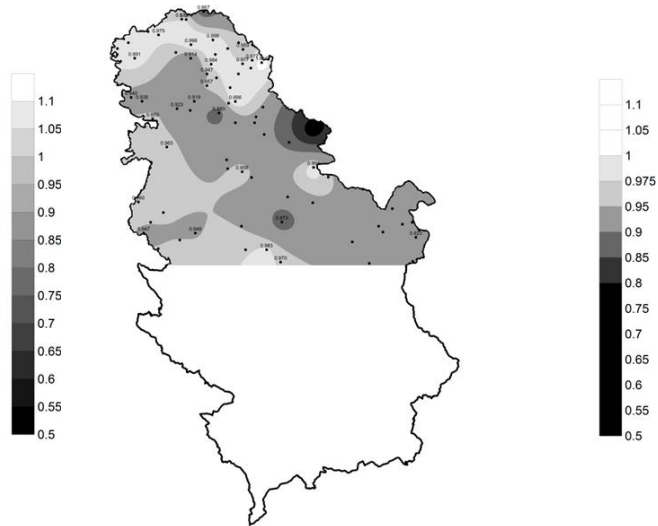


**Fig. 13. Extreme precipitation calculated for 100-year return period, homogenized vs. raw series (mm/mm)**

Comparison of the calculated extreme precipitation from the three types of data is given as a proportion of one value vs. another, using a dimensionless value (mm/mm). Similarly to calculation of climate indices, extreme daily precipitation was calculated from selected raw and homogenized series as well as series from the grid point nearest to the measurement location. In order to get an insight of any possible spatial pattern of values, the proportion values are given as a map. Such map does not represent a true spatial distribution of values for every pixel of the map; it only emphasizes differences between the results derived from the series.



**Fig. 14. Extreme precipitation calculated for 100-year return period, gridded vs. homogenized series (mm/mm)**



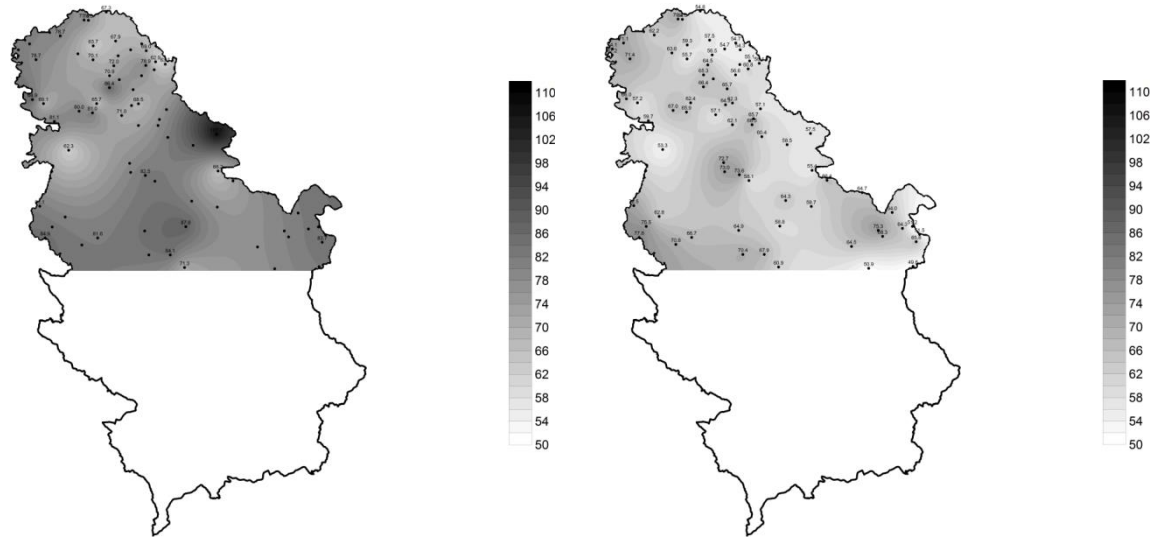
**Fig. 15. Extreme precipitation calculated for 100-year return period, gridded vs. homogenized series (mm/mm)**

Effects of the homogenization procedure is presented through **comparison of homogenized vs. raw series** (Figure 13). Homogenized series reduced maximum precipitation to 85% of values from raw series. Spatial pattern of changes seem to depend upon network density. The highest ratios are calculated for northern area with dense network, while lower values are returned from areas with sparse network (eastern and central part of the examined area). Another important influence is noted for Vršac, close to the eastern border, where an excessive maximum value is present in war series. Homogenization procedure reduced its maximum value to 77% of raw series value. This reduction is the most important effect of homogenization reflected upon extreme values.

Effects of spatial interpolation is presented through **comparison of gridded vs. homogenized series** (Figure 14). Similarly to comparison of homogenized vs. raw series, the highest ratios are in high network density area, while low ratios maintain in the area with sparse stations. Also, the lowest ratio is calculated for Vršac, due to very high maximum precipitation value in homogenized series. The value from gridded series reduced the value from homogenized series by 60%. Smoothing effect of spatial interpolation techniques is a main cause of differences between gridded and homogenized series.

After the two techniques of data processing, homogenization and gridding, a **comparison between gridded vs. raw series** shows their joint effect (Figure 15). Since the two effects are similar, it is obvious that their joint effect produces similar pattern of change.

Comparison of maps of extreme precipitation values for 100-year period from raw and gridded series show two main features. First, all values are reduced significantly, making values from gridded series only 55% to 95% of values from raw series. This is very serious problem, since neither climate assessment nor other activities (civil engineering, hydrology) would claim such reduction of extreme values to be valid. Second, spatial pattern of these values is not preserved (Figure 16). In fact, it is greatly distorted, since high values in raw series are not necessary the same in gridded series. Practically, any possible map of calculated extreme precipitation values makes no sense.



**Fig. 16. Extreme precipitation calculated for 100-year return period from raw series (left) and gridded series (right)**

## **DISCUSSION OF THE RESULTS AND CONCLUSIONS**

The comparison analysis of these indices showed that there are certain differences between raw and homogenized data. Grid point data that showed more significant differences from other two types of data. These differences include lower extreme values after homogenization and change of spatial distribution pattern.

The differences come from various factors, which include:

- elimination of outlying values (not necessarily errors, but real values) during homogenization,
- smoothing effects in spatial interpolation,
- low network density feature the highest magnitude of changes, while dense networks suffer minor losses.

The best results are derived from raw series. If raw series are not complete, filling in the gaps as a part of a homogenization procedure is quite good solution. However, completely homogenized and especially spatially interpolated series return less satisfactory results for shown practical use of daily precipitation data.

On the other hand, both homogenization and spatial interpolation procedure are optimal on a monthly basis, since performed calculations are based on monthly temporal resolution. This basis is satisfactory for the most features of climate assessments.

This survey shows the good and poor sides of homogenization and gridding daily precipitation datasets. The given results show that choice of dataset should depend upon the purpose of future surveys that engage daily precipitation data. Climatological assessments should use daily data only as a source for calculation of data in more robust temporal distribution. Therefore, homogenized and spatially interpolated values might perform better



than raw series. On the other hand, climate indices and calculation of extreme precipitation values for a return period should prefer using raw daily data.

## References

- Gilleland, E., R. Katz, and G. Young, 2004: The Extremes Toolkit (extrRemes): Weather and climate applications of extreme value statistics. useR! 2004 - The R User Conference, Vienna, Austria, 20-22 May
- Klein Tank, A.M.G, Zwiers F.W. Guidelines on Analysis of extremes in a changing climate in support of informed decisions for adaptation, WCMDP-No. 72, WMO-TD No. 1500, WMO, 2009
- Szalai, S., Auer, I., Hiebl, J., Milkovich, J., Radim, T. Stepanek, P., Zahradnicek, P., Bihari, Z., Lakatos, M., Szentimrey, T., Limanowka, D., Kilar, P., Cheval, S., Deak, Gy., Mihic, D., Antolovic, I., Mihajlovic, V., Nejedlik, P., Stastny, P., Mikulova, K., Nabyvanets, I., Skyryk, O., Krakovskaya, S., Vogt, J., Antofie, T., Spinoni, J.: Climate of the Greater Carpathian Region. Final Technical Report. [www.carpatclim-eu.org](http://www.carpatclim-eu.org).
- Szentimrey, T. Multiple Analysis of Series for Homogenization (MASH v3.03), Hungarian Meteorological Service, 2014
- Szentimrey, T. Bihari Z. Meteorological Interpolation based on Surface Homogenized Data Basis (MISH v1.03), Hungarian Meteorological Service, 2014
- Zhang, X, Yang, F. RClimDex (1.0) User Manual, Climate Research Branch, Environment Canada, 2004

# HOMOGENIZATION OF MONTHLY AIR TEMPERATURE AND MONTHLY PRECIPITATION SUM DATA SETS COLLECTED IN UKRAINE

Skrynyk O.<sup>1</sup>, Savchenko V.<sup>2</sup>, Radchenko R.<sup>2</sup> and Skrynyk O.<sup>3</sup>

<sup>1</sup> National University of Life and Environmental Sciences of Ukraine, Kyiv, Ukraine

<sup>2</sup> Taras Shevchenko National University of Kyiv, Ukraine

<sup>3</sup> Ukrainian Hydrometeorological Institute, Kyiv, Ukraine

(skrynyk\_olesya@rambler.ru, Savchenkovaleria94@gmail.com, skrynyk@univ.kiev.ua)

## Abstract

Results of homogenization of climate data series collected in Ukraine are presented in the paper. We considered two data sets. First set is monthly mean air temperature data and second one is monthly sums of precipitation. In both cases, data were collected at 174 Ukrainian climatological stations, which are uniformly distributed on Ukrainian territory.

Homogenization of climate data series was performed by means of the MASH software. We used a quasi-automatic algorithm for MASH, which was tested, proved and used in CARPATCLIM project. After homogenization we obtained test statistics which proximate critical value.

Comparison of break points detected by MASH with metadata has shown that approximately 30 % of detected break points can be explained by metadata

## 1. INTRODUCTION

Collecting data is a very important stage of climatologic research because its quality always affects the result, and often becomes a great source of errors. According to the results presented by Venema et. al. (2014) we can conclude that there exists a variety of error sources - starting from those caused by external factors, and finishing with negligence during the digitizing process conducted manually. Thus, the very first thing after collecting data should be to insure quality control. To complete this task, we can use a number of methods and programs. According to the results of benchmarking conducted in the framework of COST Action HOME (Venema et. al, 2012), the MASH homogenization procedure is one of the best ones.

The purpose of this study is to analyze two data sets collected in Ukraine, and to prepare them for wide use in UHMI research projects. The MASH procedure was used to perform quality control and to homogenize long monthly air temperatures and monthly precipitation sum data sets collected in Ukraine during time period 1961-2010 (2009 for precipitation sums).

## 2. DATA AND METHODOLOGY

There are two steps in meteorological data quality control in Ukraine. First, the information passes an internal observation station quality control procedure. Second, the meteorological information passes through quality control performed by the Meteorological Division of Central Geophysical Observatory of the Ministry of Emergencies of Ukraine. Finally, the meteorological information is published in special tables, which insure that it is both qualitative and reliable. However, we should mention that checking homogenization of temperature and precipitation data using up-to-date homogenization procedures has not yet been performed.

Thus, we have under consideration two data sets. The first set is monthly mean air temperature data, and the second is monthly sums of precipitation. In both cases, data were collected at 174 Ukrainian climatological stations, which are uniformly distributed throughout Ukrainian territory (Fig. 1). The mean distance between stations is approximately 50 km in flatland areas, and 30 km in mountainous regions. The period of interest is from 1961 to 2010. The original time series did not have a lot of missing data - less than 1 % in every time series.

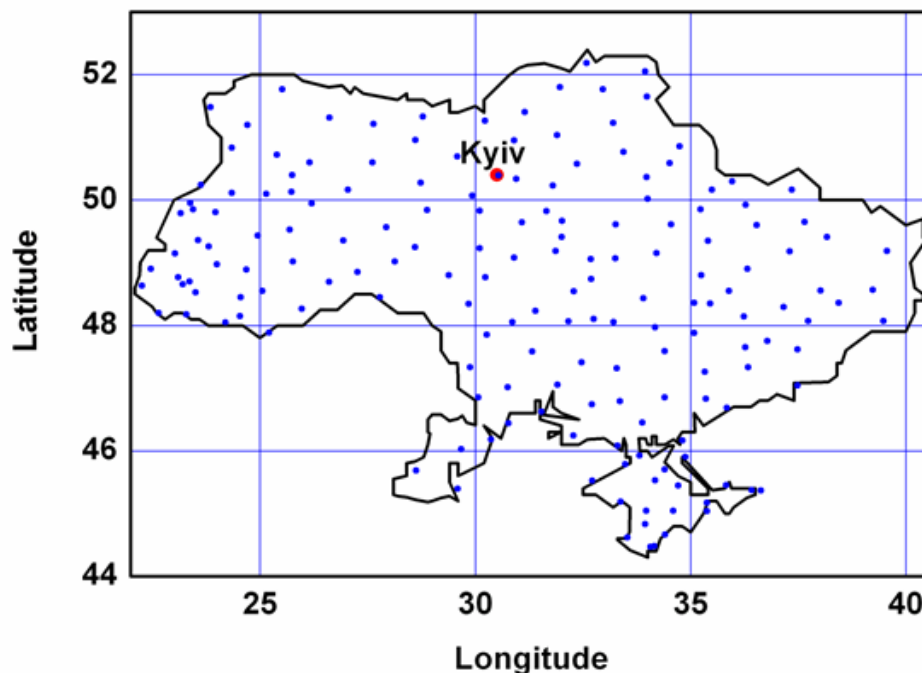


Fig. 1. Location of climatological stations used in the research

A historical description of Ukrainian stations is published in special issues (Climatologic handbook, USSR, 1968; Climatologic handbook, Ukraine, 2011), which serve as sources of possible break points (metadata). Most Ukrainian stations were relocated several times, sometimes at a very long distance (about 20 km). Furthermore, the measuring methodology for temperature was changed at all Ukrainian climatologic stations. In 1966, four-time observations were replaced with eight-time observations. We should note that some correction

for monthly air temperature for a period with four time observations was proposed and implemented in published tables

MASH software (Szentimrey, 1999) was applied according to the algorithm shared among participants of the CARPATCLIM project by T. Szentimrey. The algorithm is quasi-automatic. However, in our case, the use of the algorithm was insufficient, because test statistics (TS) still remained quite high. Therefore, according to the MASH Manual (Szentimrey, 2011), some additional steps were used.

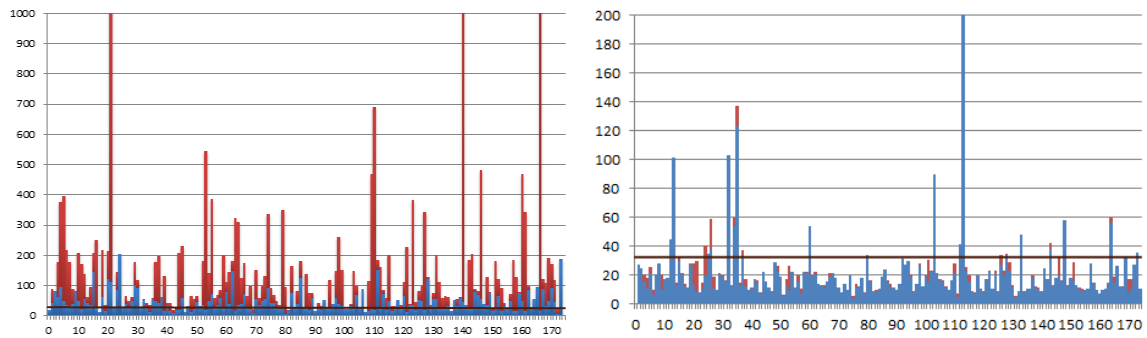
Firstly, the MASH procedure was used without metadata in order to assess the efficiency of this software in break points detection. Exactly these results are presented below in the text. After this, the MASH software was run second time with collected metadata. According to the recommendations of Venema et. al. (2012), the results of the second analysis can be considered more accurate. Therefore, these results will be used for further climatologic studies.

### 3. RESULTS

The inhomogeneity of the original air temperature time series was very high. For example, the average test statistics (TS) for yearly time series was equal to 301.33, which exceed the critical value (equal to 20.86) by more than 14 times. The TS for certain time series reached very high values (maximal TS was 25661.19) (Fig.2). After homogenization, we obtained the average TS 23.79, which seems to be acceptable (Fig. 3).

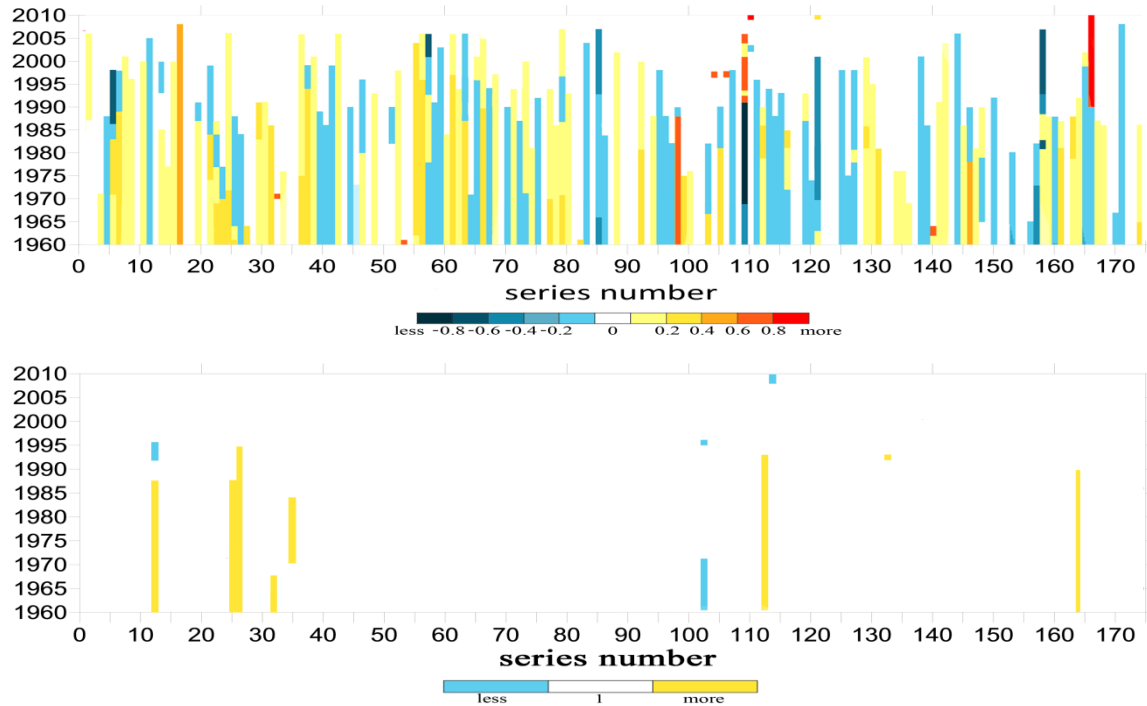


**Fig. 2. Example of inhomogeneity in air temperature data: top image - time series before homogenization process, the middle image - time series after homogenization process, bottom image - inhomogeneity**



**Fig. 3.** Test statistics (TS) of yearly air temperature (on the left) and yearly precipitation sum (on the right); red tables - TS before homogenization, blue tables - TS after homogenization, solid black line - critical value

The precipitation time series were much more qualitative. The average TS for yearly time series was equal to 22.29, which was less than the critical value. However, the TS for several time series were still very high (Fig. 3). This means that homogenization was necessary.



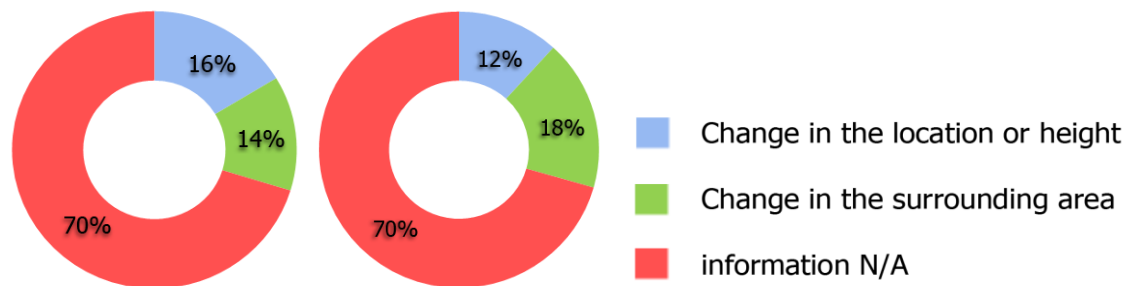
**Fig. 4.** Inhomogeneity in monthly air temperature (upper image) and monthly precipitation sum (bottom image)

According to the results, we constructed graphs of our time series inhomogeneity, based upon the MASHDRAW.BAT principle, but covering information from all stations during the entire period of time (Fig.4). It is clear from the observations the temperature time series include more inhomogeneities than the precipitation time series. It should also be noted that in the case of temperature, the percentage of positive and negative inhomogeneity was almost equal. Regarding the value, it should be mentioned that in 80% of cases, it was in a range from -0,2 to 0,2, and in 90% of the cases the range varied from -0,4 to 0,4. Little can be said about the inhomogeneity in the monthly precipitation sum. Only a few stations required

homogenization procedure, and in 90% of cases, the inhomogeneity values were lower than 0,3.

The MASH procedure reduced substantially the test statistics also in the case of monthly series. That can be seen from the files “Verisum”. Therefore, the homogenized time series (the homogenized data sets) can serve as a good base for further studies of current state of regional climate in Ukraine.

In both cases (temperature and precipitation), homogenization was performed without any metadata. In order to assess the efficiency of the homogenization software in break points detection, we collected metadata (possible break points) from historical descriptions of Ukrainian climatologic stations. Comparison of break points detected by MASH with metadata has shown that approximately 30 % of detected break points can be explained by metadata (Fig. 5).



**Fig. 5. The average percentage of break points which can be explained by metadata in monthly air temperature (to the left), and monthly precipitation sum (to the right) time series.**

#### 4. CONCLUSIONS

According to the results, MASH software detected many inhomogeneities, but most of them were less than 1 degree in absolute value. Comparison of available metadata and research results showed matches in almost 30% of all cases, which corresponds to a higher than average result.

The homogenized time series can serve as a good base for future studies of current state of regional climate in Ukraine, and also be used as a reference series for homogenization and quality control of another Ukrainian stations data.

## References

- Climatological handbook USSR, issue 10, Ukrainian SSR “History and physiographic description of meteorological stations, Kiev, 1968, 456 p. (in Russian)
- Climatological handbook “History and physiographic description of Ukrainian meteorological stations, Kiev, 2011, 462 p. (in Ukrainian)
- Szentimrey, T., 1999: Multiple Analysis of Series for Homogenization (MASH), Proceedings of the Second Seminar for Homogenization of Surface Climatological Data, Budapest, Hungary; WMO, WCDMP-No. 41, pp. 27-46.
- Szentimrey, T., 2011: Manual of homogenization software MASHv3.03, Hungarian Meteorological Service, p. 66.
- Venema, V.K.C., Mestre, O., Aguilar, E., Auer, I., Guijarro, J.A., Domonkos, P., Vertacnik, G., Szentimrey, T., Stepanek, P., Zahradnicek, P., Viarre, J., Müller-Westermeier, G., Lakatos, M., Williams, C.N., Menne, M.J., Lindau, R., Rasol, D., Rustemeier, E., Kolokythas, K., Marinova, T., Andresen, L., Acquaotta, F., Fratianni, S., Cheval, S., Klancar, M., Brunetti, M., Gruber, C., Prohom Duran, M., Likso, T., Esteban, P., and Brandsma, T. 2012: Benchmarking homogenization algorithms for monthly data, *Clim. Past*, 8, 89-115.
- Venema, V., Aguilar, E., Auchmann, R., Auer, I., Brandsma, T., Chimani, B., Gilabert, A., Mestre, O., Toreti, A., Vertacnik, G., Domonkos, P., 2014: [Inhomogeneities in daily data](#), 8<sup>th</sup> seminar for homogenization and quality control in climatological databases and 3<sup>rd</sup> conference on spatial interpolation techniques in climatology and meteorology, Budapest, Hungary, 12 – 16 May 2014

# HOMOGENIZATION PROCESS IN THE CLIMATE OF CARPATHIAN REGION PROJECT, VERIFICATION RESULTS

M. Lakatos<sup>1</sup>, T. Szentimrey<sup>1</sup>, Z. Bihari<sup>1</sup>, and S. Szalai<sup>2</sup>

<sup>1</sup>Hungarian Meteorological Service, Hungary

<sup>2</sup>Szent István University, Hungary  
(lakatos.m@met.hu)

## Abstract

This paper focuses on the homogenization activity was performed in the CARPATCLIM project. The main aim of CARPATCLIM was to create a daily harmonized gridded dataset during the period between 1961 and 2010 for the Carpathian Region. For ensuring the usage of largest possible station density the necessary work phases were implemented by counties with applying same methods and software. The common used methods and software in the project were the method MASH (Multiple Analysis of Series for Homogenization; *Szentimrey*) for homogenization, quality control, completion of the observed daily data series; and the method MISH (Meteorological Interpolation based on Surface Homogenized Data Basis; *Szentimrey and Bihari*) for gridding of homogenized daily data series. Besides the common software, the harmonization of the results across country borders was performed also by near border data exchange. The main steps of the homogenization process executed and the verification of the homogenization along with the quality control results are introduced in this paper.

## 1. INTRODUCTION

The main aim of CARPATCLIM project (*CARPATCLIM homepage*) was the spatial and temporal examination of the climate of the Carpathian Region using harmonized data and standard methodology. The consortium led by the Hungarian Meteorological Service (OMSZ) together with 10 partner organizations from 9 countries in the region was supported by the JRC to create a daily harmonized gridded dataset during the period between 1961 and 2010. The target area of the project partly includes the territory of Czech Republic, Slovakia, Poland, Ukraine, Romania, Serbia, Croatia, Austria and Hungary. 415 climate stations and 904 precipitation stations were used in the project to achieve the objectives. The final outcome of the CARPATCLIM is a  $\sim 10 \times 10$  km resolution homogenized and gridded dataset on daily scale for 13 basic meteorological variables and several climate indicators, 37 in total, on different time scales from 1961 to 2010.



## 2. METHODOLOGY

Uniform process of data homogenization was crucial due to the fact that significant differences might be occurred between the measurements and data handling of participant countries during the examined fifty-year-long period. The necessary work phases were implemented by country. The common used methods and software in the project were the method MASH (Multiple Analysis of Series for Homogenization; *Szentimrey*, 2011) for homogenization, quality control, completion of the observed daily data series; and the method MISH (Meteorological Interpolation based on Surface Homogenized Data Basis; *Szentimrey and Bihari*, 2007) for gridding of homogenized daily data series. The high quality of **times** series got through the commonly used MASH procedure are guaranteed by the excellent monthly benchmark results from the COST “HOME” Action (*Venema et al.*, 2012). Besides the common software, the harmonization of the results across country borders was performed by near border data exchange.

### 2.1. Main features of MASH

The original MASH (*Szentimrey*, 1999) procedure was developed for homogenization of monthly series. The present version: MASHv3.03 (*Szentimrey*, 2011) has been expanded for daily series as well. The main features of the applied procedure in CARPATCLIM are summarized here.

The MASHv3.03 software consists of two parts.

Part 1: Quality control, missing data completion, and homogenization of monthly series:

- Relative homogeneity test procedure.
- Step by step procedure: the role of series (candidate or reference series) changes step by step in the course of the procedure.
- Additive (e.g., temperature) or multiplicative (e.g., precipitation) model can be used depending on the climate elements.
- Providing the homogeneity of the seasonal and annual series as well.
- Metadata (probable dates of break points) can be used automatically.
- Homogenization and quality control (QC) results can be evaluated on the basis of verification tables generated automatically during the procedure.

Part 2: Homogenization of daily series:

- Based on the detected monthly inhomogeneities.
- Including quality control (QC) and missing data completion for daily data. The quality control results can be evaluated by test tables generated automatically during the procedure.

These attributes are favourable to achieve the project goals in CARPATCLIM. The time resolution of variables is daily, the upgraded version of MASHv3.03 is able to homogenize and control these daily data as well.

### 3. HOMOGENIZATION PROCESS IN THE CARPATCLIM

According to the project specification, the elements listed in *Table 1* have to be homogenized in the period of 1961–2010. The chosen homogenization model is depending on the distribution of given element. Additive model is used except in case of precipitation and wind speed, where the appropriate model is multiplicative.

*Table 1.* The set of meteorological variables in daily temporal resolution were homogenized (JRC, 2010)

VARIABLE	DESCRIPTION	UNITS
TA	2 m mean daily air temperature	°C
TMIN	Minimum air temperature from 18:00 to 06:00	°C
TMAX	Maximum air temperature from 06:00 to 18:00	°C
P	Accumulated total precipitation from 06:00 to 06:00	mm
DD	10 m wind direction	0°-360°
VV	10 m horizontal wind speed	m/s
SUNSHINE	Sunshine duration	hours
CC	Cloud cover	tenths
RGLOBAL	Global radiation	MJ/m <sup>2</sup>
RH	Relative humidity	%
PVAPOUR	Surface vapour pressure	hPa
PAIR	Surface air pressure	hPa

The main steps of homogenization in CARPATCLIM were as follows.

1. Near border data exchange before homogenization.
2. Homogenization, quality control, completion of the daily data series on national level by using near border data series.
3. Near border data exchange after homogenization.

MASH is an automatically working software. The test results of the homogenization and quality control (e.g., detected errors, degree of inhomogeneity of the series system, number of break points, estimated corrections, and certain verification results) are documented in automatically generated tables during the homogenization process.

### **3.1. Steps of creation of the homogenized station data series in CARPATCLIM**

#### **I. Compilation of the raw station data series of each country.**

1. Selection of the stations (with the help of spherical coordinates:  $\varphi$ ,  $\lambda$ ).
2. Collecting the daily station data series (missing data are allowed) and the metadata per countries. Exchange of the near border raw data series and the existing metadata between the neighboring countries.

#### **II. Homogenization, quality control, data completion of the station data series by MASH v3.03 on national level, using near border data.**

3. Derivation of monthly station data series from the daily station data series collected in step I.2. Homogenization, quality control, data completion of the monthly station data series. Metadata (probable dates of break points) can be used automatically.
4. Daily station data series (step I.2): homogenization, quality control, data completion. This procedure is based on the results of step II.1.
5. Exchange of the near border homogenized data for cross-border harmonization and for gridding (Module 2 of the project: modeling, interpolation).
6. Evaluation of the verification results of the homogenization and quality control. Controlling of the cross-border harmonization of the data series. Note that further cross-border harmonization is achieved after the modeling part of the gridding procedure in Module 2.

Summary of the main steps of homogenization of daily data series with quality control and missing data completion in CARPATCLIM are as follows:

- Monthly series derivation from daily series.
- MASH homogenization procedure for monthly series, estimation of monthly inhomogeneities. (Metadata can be used automatically.)
- Smooth estimation of daily inhomogeneities on the basis of estimated monthly inhomogeneities.
- Automatic correction of daily series.
- Automatic quality control (QC) of homogenized daily data.
- Automatic missing daily data completion.
- Monthly series derivation from the homogenized, quality controlled, and completed daily data.

- Test of homogeneity for the new monthly series with using the automatic verification results.

### **3.2. Verification of the homogenization**

This chapter is an overview of the evaluation of the implemented homogenization process. Validation is an essential part of the process, to make sure that the data quality increased as a result of homogenization. Hence a verification part is integrated into the MASH system for interpretation of the outcomes, it makes the evaluation of the different phases of the homogenization possible from the initial to the final stage. The basic conception of the verification test is that the confidence in the homogenization may be increased by the joint comparative mathematical examination of the original and the homogenized data series.

Two types of outcomes of the MASH software can be separated. The first type of output is the files containing the homogenized, controlled, and completed series, inhomogeneity series, detected breaks, and detected errors. The second type of output is the files containing the test results and verification tables in order to evaluate the homogenization. The verification tables contain the test statistics values before and after homogenization, measures to characterize the modification of series, the spatial representativity of the station network, and the evaluation of metadata. The quality control results for the daily data are also included.

The verification procedure based on hypothesis test results. The null hypothesis is that examined series are homogeneous. The test statistics can be compared to the critical value before and after homogenization. The critical values belong to different significance levels are built in the MASH software (it is 20.86 on the 0.05 significance level in our case). The homogenization is successful if the test statistics after homogenization is low. The theoretical background and more details of the derivation of the verification statistics can be found in MASH manual (*Szentimrey, 1999*).

The test statistics before (TS<sub>b</sub>) and after homogenization (TS<sub>a</sub>) and characteristics of the modified series are presented in this paper. Annual statistics are examined here; though all of them are produced automatically on the monthly and seasonal scales altogether. *Tables 2 to 4* contain the average measures for maximum and minimum temperatures and precipitation for each of the station systems and the QC results alike. Number of the partners in the header lines is as follows: Hungary and Croatia with their jointly handled dataset (1), Serbia (2), Romania (3), Ukraine (4), Slovakia (5), Poland (6), Czech Republic (7). The representativity is about 50 km for climate stations and 25 km for precipitation stations, respectively. Participants have contributed to the project with data of 415 climate stations and 904 precipitation stations in all.

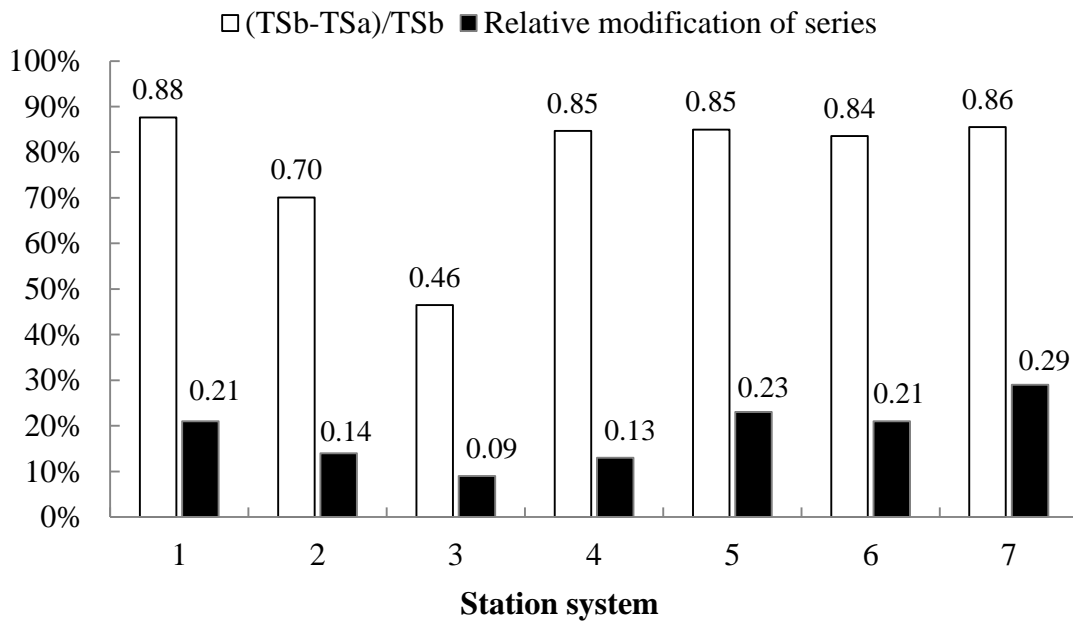
The TSa has to be near to the critical value or much less than the TSb if the homogenization is acceptable. Moreover, the measures of the relative modification are expected to be in accordance with the relative change of the test statistics:  $(TSb-TSa)/TSb$ . The applied statistics for the measure of the relative modification is in fact the ratio of the RMSE (root mean square error) and the standard deviation. If the significant modification of series induces weak decreasing in the degree of inhomogeneity, overdrawing the series is unnecessary and erroneous. *Tables 2–4* containing the summary statistics and the complementary diagrams in *Figs. 1–3* support the evaluation of homogenization.

The degree of inhomogeneity of the raw minimum temperatures (*Table 3.*) is substantially higher for Serbia (2) and much higher for the Hungarian and Croatian (1) dataset than in case of the maximum temperatures. The relative modification (42%) for the Hungarian and Croatian (1) series is achieved the most, although the largest improvement (*Fig. 2.*). The Serbian (2) system has been upgraded in the same rate by less relative modification. The Slovakian (5) system is near to homogeneous after processing. Relative changes of the test statistics are small in the Romanian (3) and Ukrainian (4) series, in accordance with the low value of relative modification. At the Czech Republic (7), the degree of homogeneity increased with relatively high modification. It can be found that MASH reduced the inhomogeneity of all systems, but less than in the case of maximum temperatures. The QC results relating the minimum temperatures show that the number of erroneous data per station is the largest in the Ukrainian (4) system. The Romanian (3) and Ukrainian (4) series contained more than 400 (°C) negative error and almost 100 (°C) positive errors in the data. The smallest correction has to be performed in the Czech (7) system, although it is a minor system with 18 stations.

Summary results of quality control and the homogenization performed in the project can be followed up and reported based on these tables. Verification statistics can be added to the homogenized series as the newly created metadata.

*Table 2.* Average test statistics and quality control (QC) results for maximum temperature

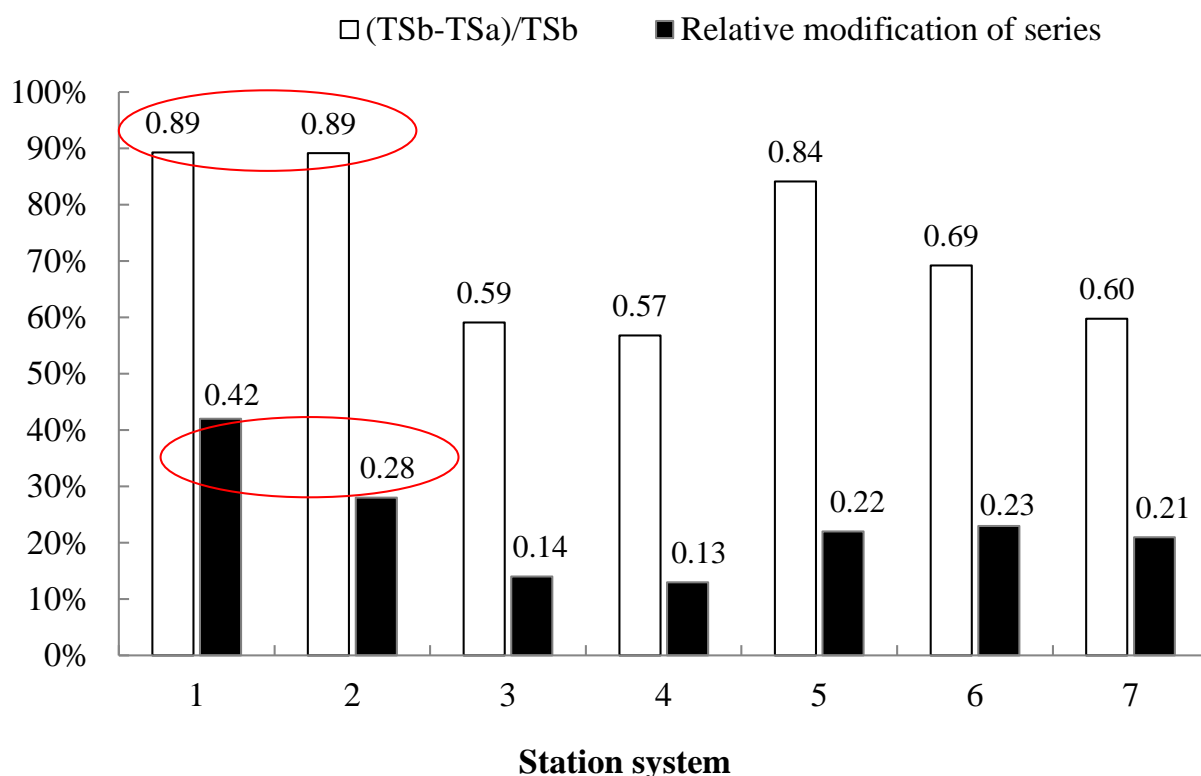
Station Sytem	1	2	3	4	5	6	7
Number of stations	68	39	140	53	59	38	18
TS after homog. (TSa)	23.6	55.7	39.0	23.7	26.4	24.8	26.7
TS before homog. (TSb)	190.7	186.2	72.9	154.0	175.6	150.6	184.3
Relative modification (%)	21	14	9	13	23	21	29
Total number of errors	6307	3811	10241	5444	4542	3288	1400
Maximal positive error (°C)	10.9	13.5	996.6	107.7	11.3	22.7	10.4
Minimal negative error (°C)	-2.3	-7.5	-21.0	-22.0	-14.5	-26.3	-6.2



**Fig. 1. Verification results for maximum temperature.**

*Table 3 Average test statistics and quality control (QC) results for minimum temperature*

Station System	1	2	3	4	5	6	7
Number of stations	68	39	140	53	59	38	18
TS after homog. (TSa)	24.3	52.5	52.5	51.9	28.5	43.5	37.8
TS before homog. (TSb)	227.5	484.7	128.3	120.3	179.7	141.3	93.9
Relative modification (%)	42	28	14	13	22	23	21
Total number of errors	4110	2161	6689	4111	3197	2592	375
Maximal positive error (°C)	23.7	11.8	95.1	79.3	14.9	15.9	0.7
Minimal negative error (°C)	-9.7	-8.0	-416.6	-417.6	-9.9	-10.0	-1.1



**Fig. 2. Verification results for minimum temperature**

Analyzing the precipitation results, we have to take into consideration that the MASH procedure carefully detects the break points. Lower inhomogeneity arose for the precipitation series than for temperatures (*Table 4*). During the homogenization, all of the networks became more homogeneous; nevertheless, the modification was precautionous. The test statistics indicates that the Polish (6) system was the most inhomogeneous, and the improvement is also little afterward, although the similar relative modification caused higher improvement than in the Romanian (3) system (*Fig. 3.*). The Slovakian (5) dataset passed through the most advance, at the expense of remarkable modifications of the series comparing to the others. Resulting from the QC numerous errors were detected, about in the rate of the amount of contributed stations. The amplitude of the errors in several systems is higher towards extremely heavy precipitations.

Table 4. Average test statistics and quality control (QC) results for precipitation

Station sytem	1	2	3	4	5	6	7
Number of stations	233	114	182	57	165	102	51
TS after homog. (TSa)	21.6	31.27	28.09	25.61	21.89	38.97	35.53
TS before homog. (TSb)	27.93	34.73	31.88	28.98	38.17	46.29	39.77
Relative modification (%)	4	5	6	3	10	5	4
Total number of errors	1531	672	975	313	803	408	223
Maximal positive error (mm)	71.94	230.27	10.27	179.46	94.29	93.36	60.38
Minimal negative error (mm)	23.24	-36.87	-1.52	-5.68	-59.46	-25.47	-11.41

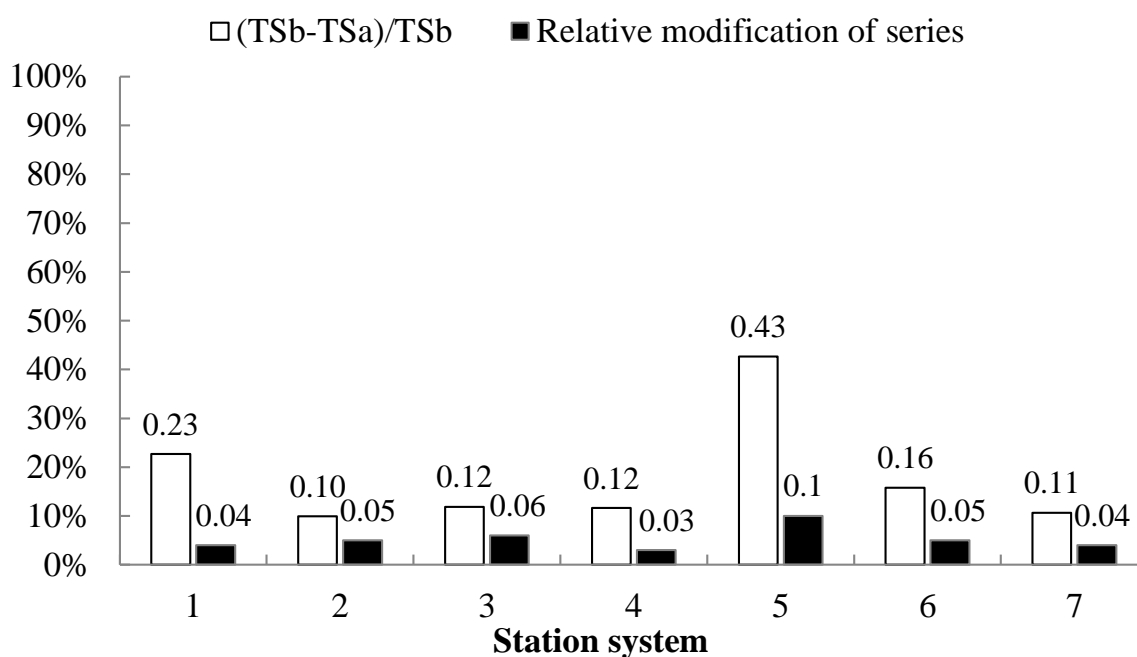


Fig. 3. Verification results for precipitation

Verification results for all the 12 elements can be followed up in the project deliverables related to the issues of the homogenization process (*DI.12*). The data rescue and digitization



activity in Module 1, and the data homogenization and QC performed by applying MASH procedure guarantee the availability of the high quality daily time series for the basic climate elements in the Carpathian region in the period of 1961–2010.

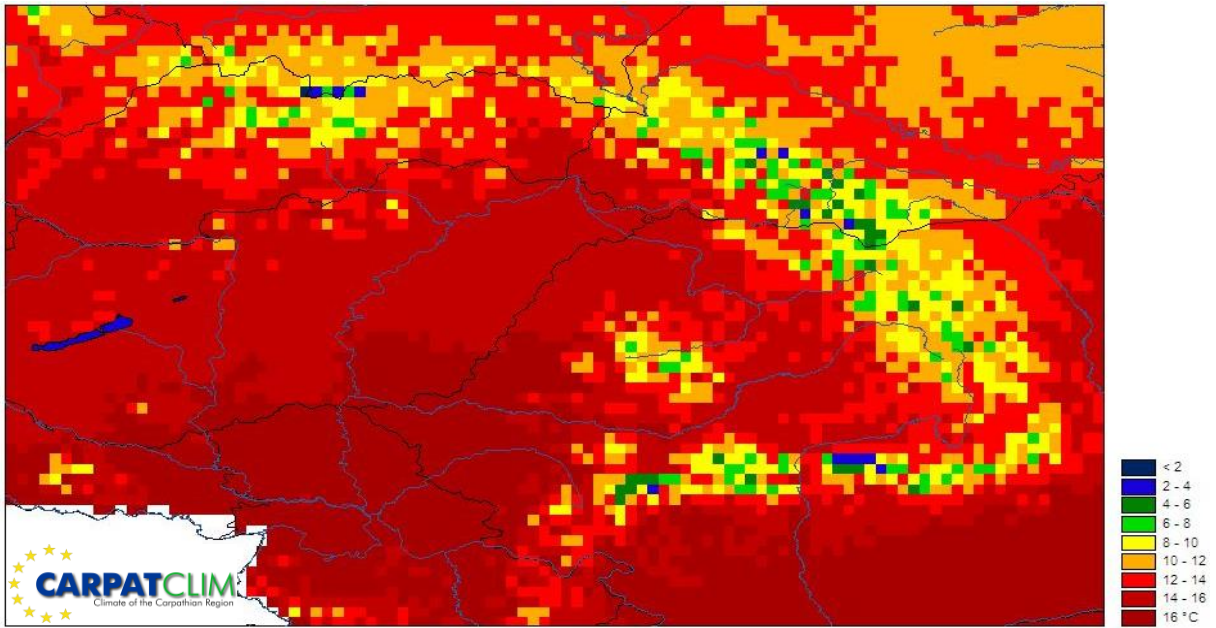
### **3.3. Interpolation of homogenized data**

The final outcome of the CarpatClim tender service is a  $\sim 10 \times 10$  km resolution gridded dataset on daily scale for elements listed in *Table 1*. Interpolation of the homogenized time series is carried out by applying the MISH (Meteorological interpolation based on surface homogenized data basis; *Szentimrey and Bihari, 2007*) method. The MISH method was developed for interpolation of meteorological data, and an adequate mathematical background was also developed (*Szentimrey et al., 2011*) for the purpose of efficient use of all the valuable meteorological and auxiliary model information. The main difference between MISH and the usual geostatistical interpolation methods is the application of the meteorological data series for modeling. In geostatistics (*Cressie, 1991*), the sample for modeling is only the predictor data, which is a single realization in time, while in meteorology there are data series, i.e., a sample in time and space as well.

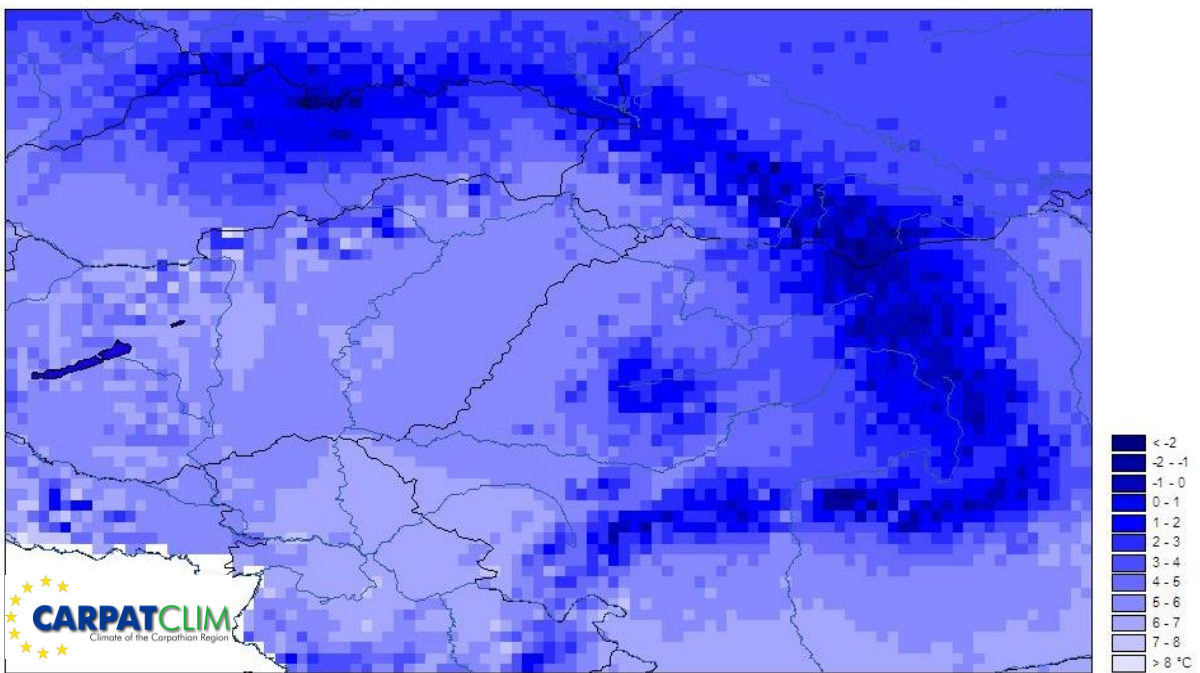
### **3.4. Homogenized, harmonized and interpolated averages**

The cross border harmonization is essential in the project to avoid breaks at the boundaries on climate maps based on the gridded data. It can be ensured by the changes of the homogenized series across the borders as it was in case of the raw data exchange. Test statistics of the cross border harmonization are detailed in a publically available deliverable D2.5 (*D2.5*). The gridding of the harmonized series was executed by countries by applying MISH, and the merging of the separate but harmonized grid parts followed up in the end.

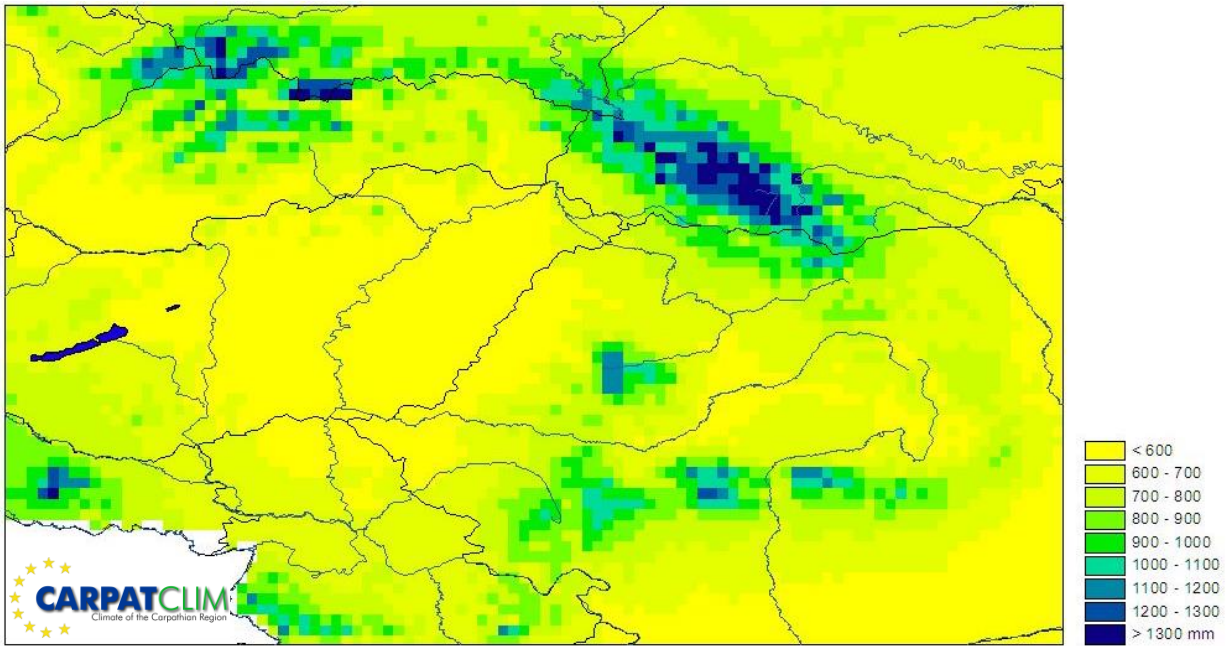
The harmonized averages for the time interval of the project: 1961-2010 are presented on maps in Figs 4-14.



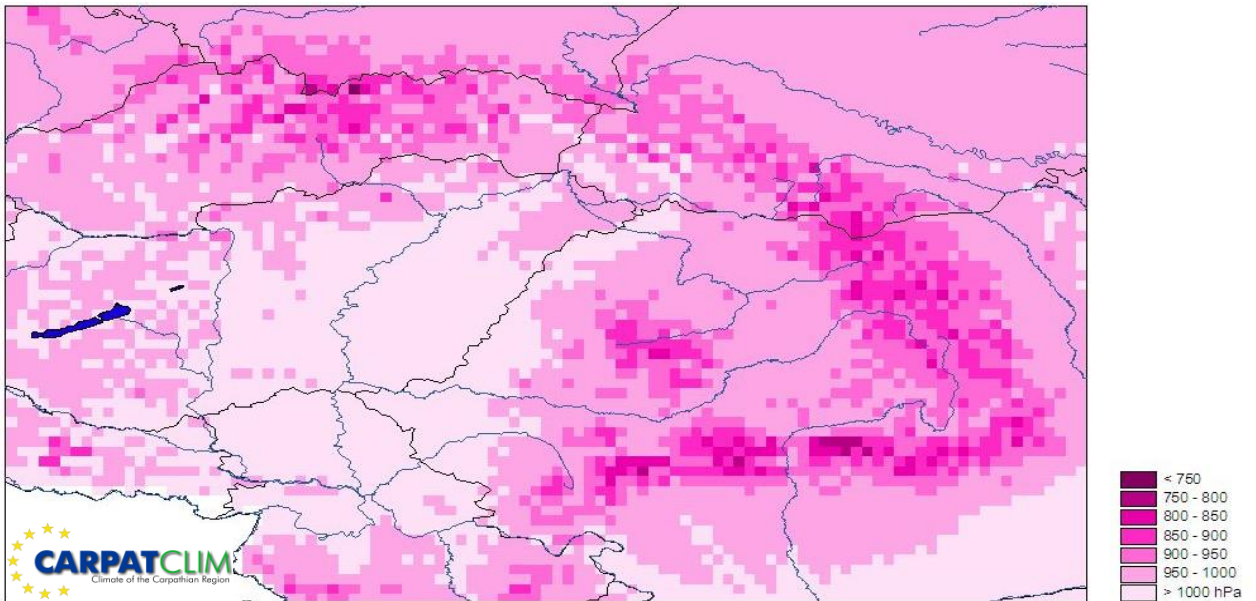
**Fig. 4. Average maximum temperature in the period of 1961-2010 for Carpathian Region**



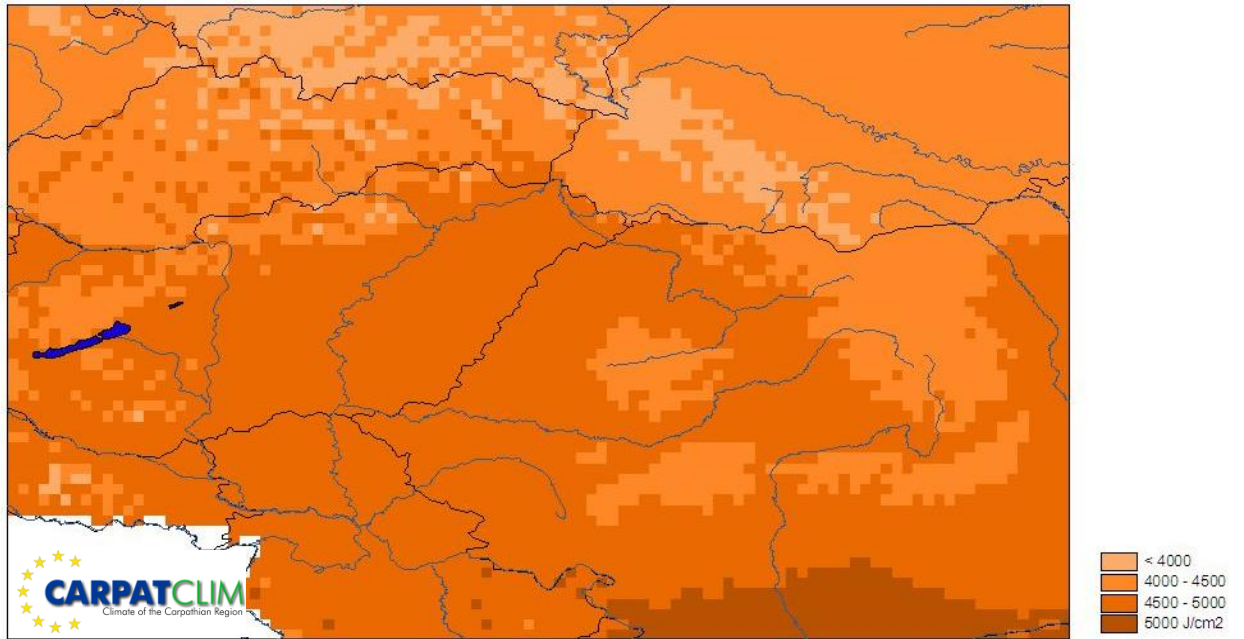
**Fig. 5. Average minimum temperature in the period of 1961-2010 for Carpathian Region**



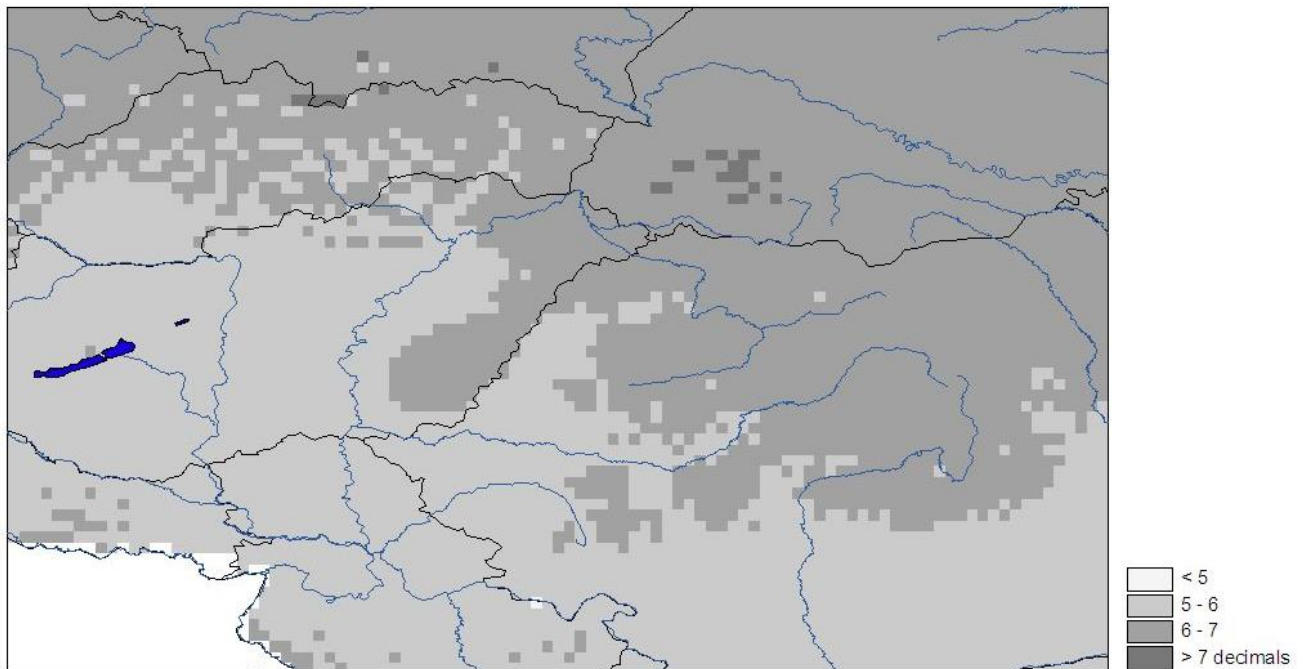
**Fig. 6. Annual mean precipitation in the period of 1961-2010 for Carpathian Region**



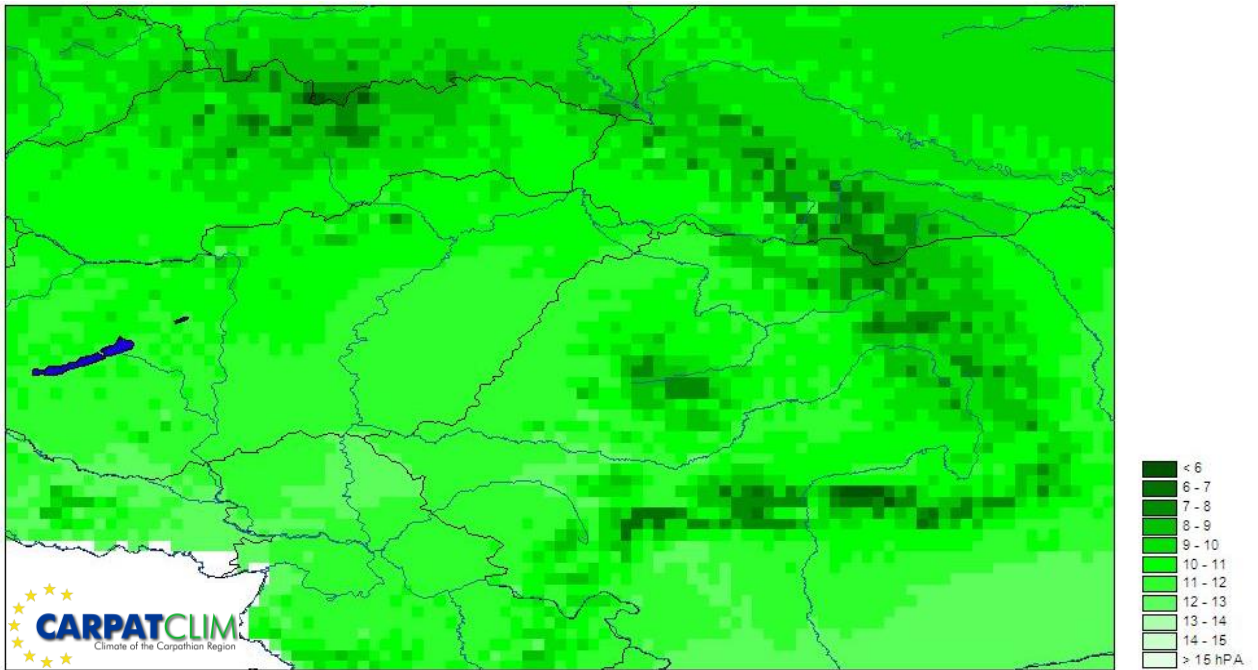
**Fig. 7. Average surface air pressure in the period of 1961-2010 for Carpathian Region**



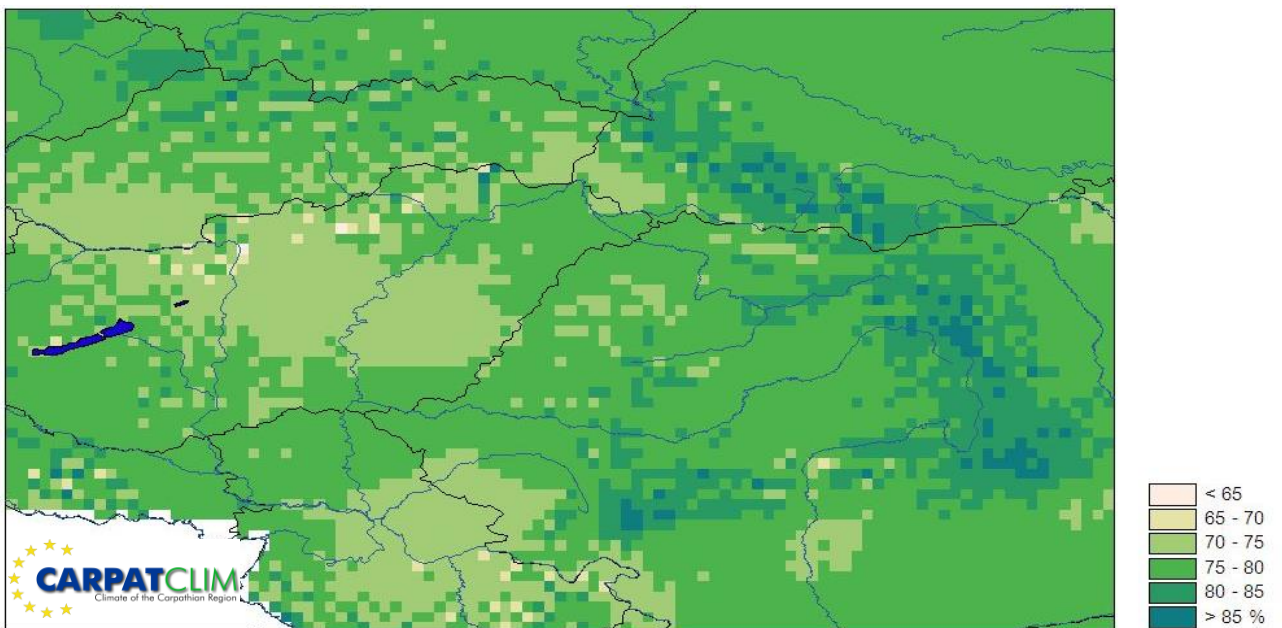
**Fig. 8. Yearly average global radiation in the period of 1961-2010 for Carpathian Region**



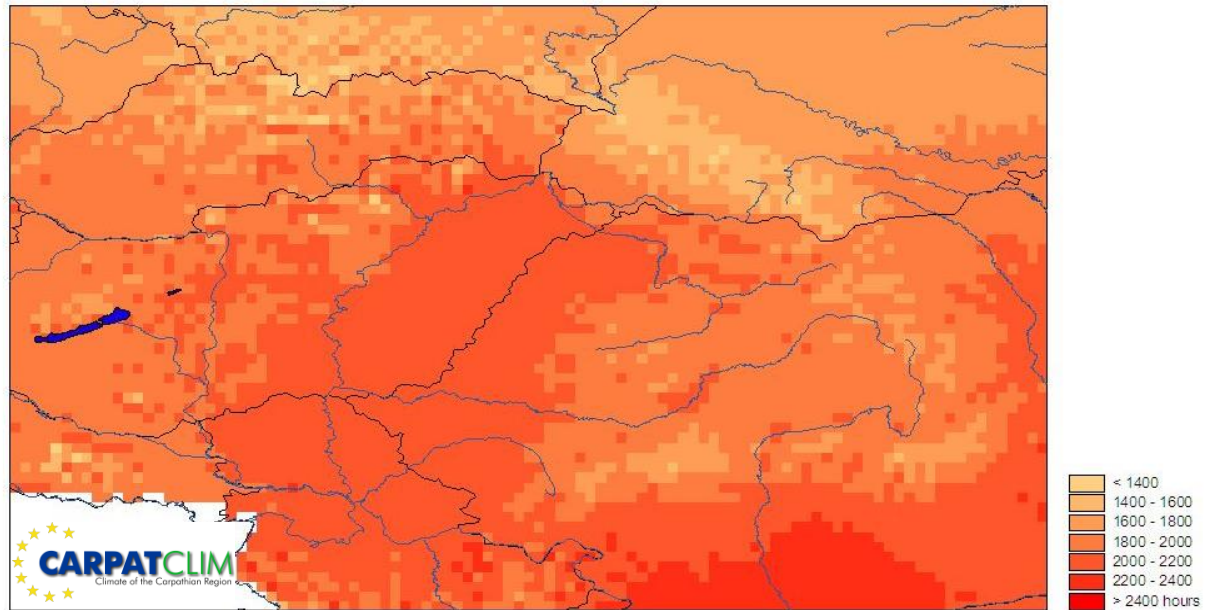
**Fig. 9. Average cloudiness in the period of 1961-2010 for Carpathian Region**



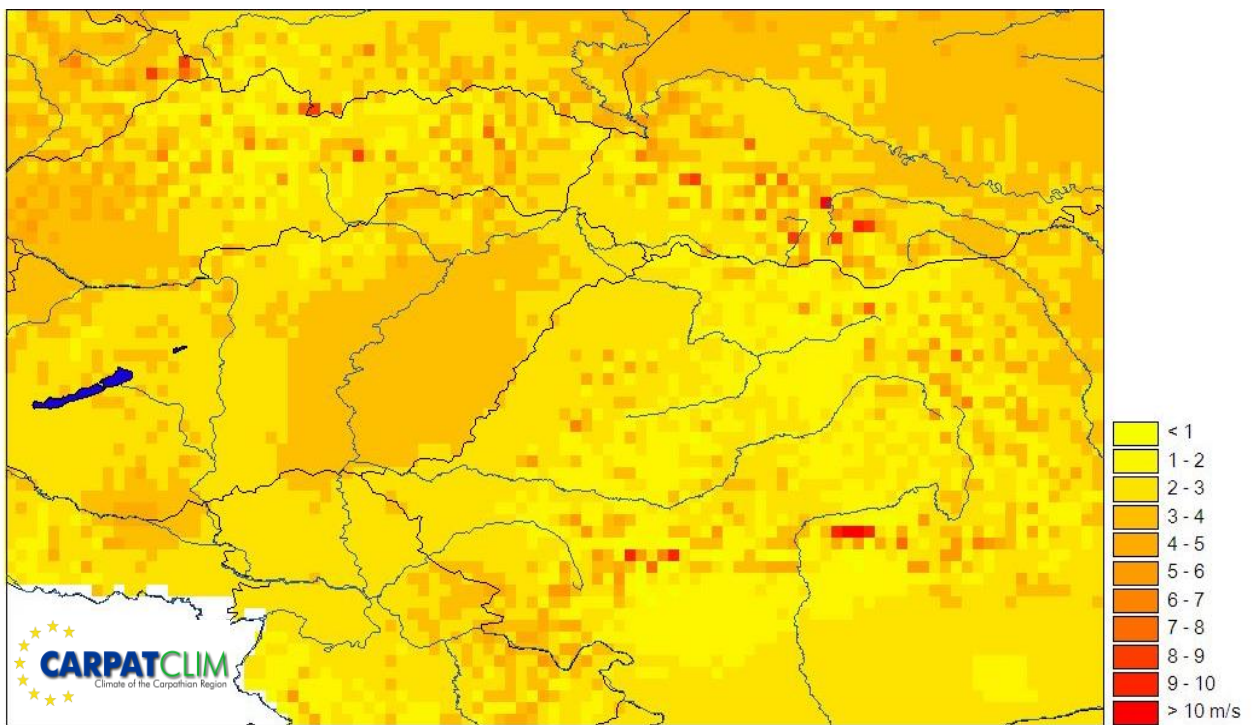
**Fig. 10. Average vapour pressure in the period of 1961-2010 for Carpathian Region**



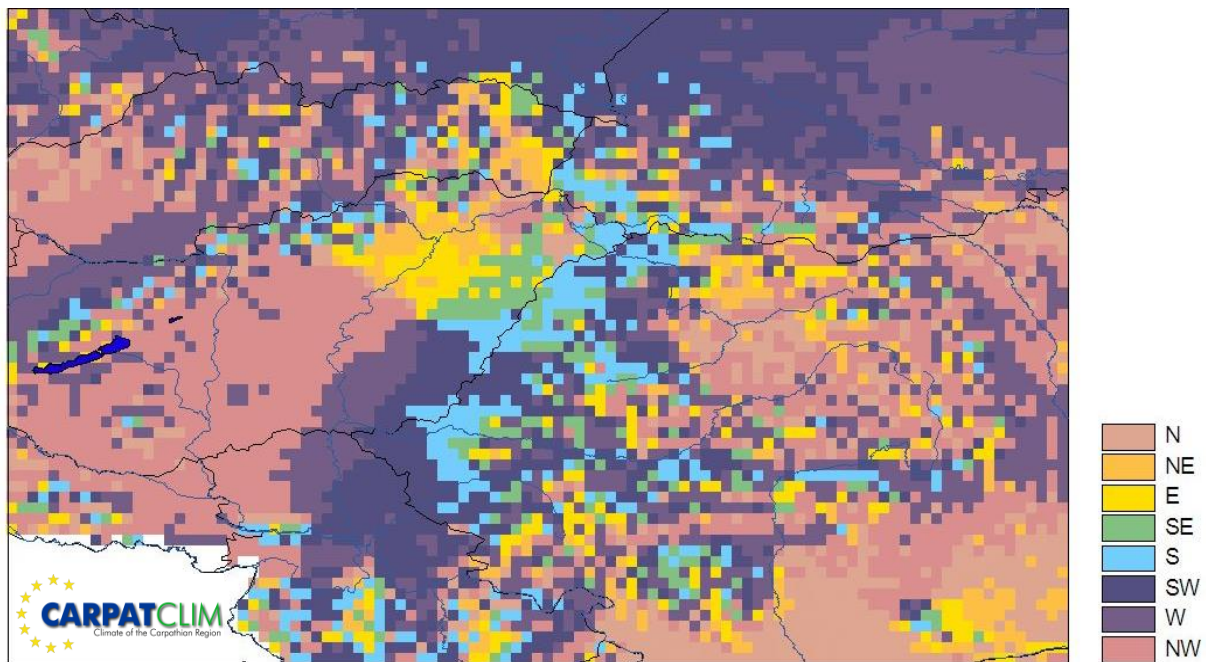
**Fig. 11. Average relative humidity in the period of 1961-2010 for Carpathian Region**



**Fig. 12. Average sunshine duration in the period of 1961-2010 for Carpathian Region**



**Fig. 13. Average wind speed on 10 m in the period of 1961-2010 for Carpathian Region**



**Fig. 14. Average wind direction on 10 m in the period of 1961-2010 for Carpathian Region**

#### 4. CONCLUSION

The CARPATCLIM project is a well-accomplished cooperation for applying a single homogenization method in a region fragmented by boundaries and a pioneer work for countervailing against differences in measuring practice and strict data policies. The high quality of times series got through the commonly used MASH procedure are guaranteed by the excellent monthly benchmark results from the COST “HOME” Action. The Climate of the Carpathian Region Project contributes to the availability of a set of homogeneous and spatially representative data to prepare climate change studies relevant in the region. The final outcome of the CARPATCLIM are the quality controlled, homogenized, in-situ daily time series and gridded data per country and the whole region as well, including a metadata catalogue with the documentation of the existing homogenized datasets. The daily grids with the metadata will be freely accessible for scientific purposes (*CARPATCLIM homepage*).

#### References

CARPATCLIM homepage: <http://www.carpatclim-eu.org/pages/home/>

Cressie, N., 1991: Statistics for Spatial Data. Wiley, New York.

D1.12: Final report on quality control and data homogenization measures applied per country, including QC protocols and measures to determine the achieved increase in data quality. <http://www.carpatclim-eu.org/pages/deliverables/>

D 2.5: Report with final results of the data harmonization procedures applied, including all protocols, per country. <http://www.carpatclim-eu.org/pages/deliverables/>

- JRC, 2010: Climate of the Carpathian Region. Technical Specifications (Contract Notice OJEU 2010/S 110-166082 dated 9 June 2010).
- <http://desert.jrc.ec.europa.eu/action/php/index.php?action=view&id=550>
- Lakatos, M., Szentimrey, T., and Bihari, Z., 2010: Application of gridded daily data series for calculation of extreme temperature and precipitation indices in Hungary. *Időjárás* 115, 99–109.
- Szentimrey, T., 1999: Multiple Analysis of Series for Homogenization (MASH). Proceedings of the Second Seminar for Homogenization of Surface Climatological Data, Budapest, Hungary; WMO, WCDMP-No. 41, 27–46.
- Szentimrey, T. and Bihari, Z., 2007: Mathematical background of the spatial interpolation methods and the software MISH (Meteorological Interpolation based on Surface Homogenized Data Basis). Proceedings from the Conference on Spatial Interpolation in Climatology and Meteorology, Budapest, Hungary, 2004, COST Action 719, COST Office, 17–27.
- Szentimrey, T., 2011: Manual of homogenization software MASHv3.03. Hungarian Meteorological Service.
- Szentimrey, T., Bihari, Z., Lakatos, M., and Szalai, S., 2011: Mathematical, methodological questions concerning the spatial interpolation of climate elements. Proceedings from the Second Conference on Spatial Interpolation in Climatology and Meteorology, Budapest, Hungary, 2009, *Időjárás* 115, 1–2, 1–11.
- UNEP, 2007: Carpathians Environment Outlook. Geneva. <http://www.grid.unep.ch>.
- Venema, V., Mestre, O., Aguilar, E., Auer, I., Guijarro, J.A., Domonkos, P., Vertacnik, G., Szentimrey, T., Štěpánek, P., Zahradnicek, P., Viarre, J., Müller-Westermeier, G., Lakatos, M., Williams, C.N., Menne, M., Lindau, R., Rasol, D., Rustemeier, E., Kolokythas, K., Marinova, T., Andresen, L., Acquaotta, F., Fratianni, S., Cheval, S., Klancar, M., Brunetti, M., Gruber, C., Duran, M.P., Likso, T., Esteban, P. and Brandsma, T., 2012: Benchmarking monthly homogenization algorithms. *Climate of the Past* 8, 89–115.



# BIASES AND CORRECTIONS OF WIND SPEED TIME SERIES

Csilla Péliné Németh<sup>1</sup>, Judit Bartholy<sup>2</sup>, Rita Pongrácz<sup>2</sup>,  
Tamás Szentimrey<sup>3</sup>, Kornélia Radics<sup>3</sup>

<sup>1</sup> Geoinformation Service of the Hungarian Defence Forces, Szilágyi Erzsébet fasor 7-9.,  
H-1024 Budapest, Hungary

<sup>2</sup> Department of Meteorology, Eötvös Loránd University, Pázmány Péter sétány 1/A, H-1117  
Budapest, Hungary

<sup>3</sup> National Meteorological Service, Kitaibel Pál utca 1., H-1024 Budapest, Hungary

## Abstract

Global climate change affects regional atmospheric circulation, and thus, these changes may modify regional weather patterns such as regional storm tracks, main wind directions and wind speed values of a specific region (Emeis 2013). Therefore, it is essential to learn both the present wind conditions of the Carpathian Region and the regional effects of global warming as deeply as possible. Reliable estimation of mean and extreme parameters of wind climate contributes to assess future conditions and to adapt to the changing climate.

To improve future wind assessments, different wind speed time series are checked, compared, and corrected. Firstly, wind speed and wind gust from Hungarian synoptic data sets (1975–2012) are homogenized using Multiple Analysis of Series for Homogenization (MASH) application (Szentimrey 1999) developed at the Hungarian Meteorological Service. Quality of measured data sets improved after homogenization, and missing data were filled automatically by the software. Secondly, interdependence of different time series are estimated by comparing measured and ERA Interim reanalysis wind data series. Our results showed that spatial difference cannot be reproduced by homogeneous gridded reanalysis data unlike in case of station measurements.

Both average and extreme values of homogenized station and grid point data sets were analyzed. We concluded that the seasonal variability is low, and biases of reanalysis data are smaller in summer than in winter. The high (i.e., 0.9 and 0.99) percentiles' values are underestimated in the reanalysis data series for most of the analyzed grid points. Due to significant differences between distributions of measured and reanalysis data sets, wind speed and wind gust extreme value analysis of present wind climate are calculated from homogeneous data sets.

Keywords: regional wind climate, homogenization, MASH, ERA Interim, extremes

## 1. INTRODUCTION

It is essential to learn the present state of wind climate and the regional changes of wind field due to global climate changes for drawing correct conclusions and estimating future consequences. Estimation of different wind climate parameters contributes to better understanding of regional environmental effects; moreover, it helps adaptation for changing climate. We aim to calculate mean and extreme wind parameters from reliable and quality controlled data series, therefore wind speed and wind gust data series of the Hungarian

synoptic station network are homogenized. This station network has been developed and installed by the Hungarian Meteorological Service (HMS) taking into account suggestions of the World Meteorological Organisation (WMO 2011). Because of the last decades' developments of measurement and communication technologies, the wind observing network has changed several times, which is quite usual. The most significant change was automation – i.e., change traditional measuring instruments into automated measuring systems – during 1995–1996. This major change introduced large variations in the climate signal, and caused inhomogeneities in the data sets. In fact, long instrumental records are very rarely homogeneous because of the changing surroundings of measuring sites (new buildings, vegetation growth, etc.). To avoid misinterpretation due to this inhomogeneity, the available time series can be divided into subsets. For instance, we used two subsets in case of previous wind climate analysis (Péliné et al. 2011) using wind data originating from traditional (1975–1994) and automated (1997–2012) measuring systems.

In addition to automation other causes may also lead to inhomogeneities such as substitution or relocation of weather stations, changing anemometer type or aging of the instruments, changes in measuring height, surroundings (e.g., urbanization), surface coverage, and roughness. Therefore, documentation of metadata is a crucial issue during any kind of meteorological measurement.

The above-mentioned changes could result in inhomogeneities, which cannot be explained by climatological reasons. Brake points in the data sets coincide with change in the probability distribution function of the measurements. These inhomogeneities must be detected and removed before further analyses. For this purpose mathematical methods are widely used, one of them is the Multiple Analysis of Series for Homogenization, MASH v3.03 (Szentimrey, 1999, 2011) developed in HMS. This technique is used here for homogenization of available daily wind speed time series between 1975 and 2012 for records of 19 Hungarian synoptic stations (Péliné et al. 2014). As a result, quality of the measured datasets improved significantly and reliability of datasets enhanced after filling the missing data.

Homogeneity of reanalysis data sets is also tested with MASH application and it has been proved that gridded data sets are homogeneous. In addition, measured and reanalysis data series are compared with estimation of Weibull parameters calculated from wind speed distributions.

## **2. HOMOGENIZATION WITH MASH APPLICATION**

A homogeneous climatological time series can be defined as time series where variability is only caused by changes in weather and climate (Aguilar et al. 2003). To decide whether or not a long time series is homogeneous, there are different detection and correction methods available for possible use. These methods are all based on mathematical formulation and climatological experience, however, their performances are different. Objective comparison of these existing methods was carried out in the framework of a scientific programme COST Action HOME ES0601: Advances in Homogenization Methods of Climate Series: an

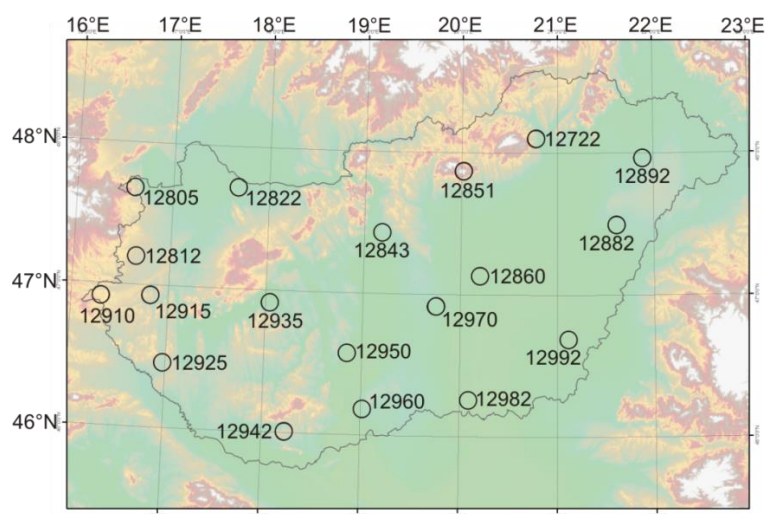
integrated approach (HOME 2011; Szentimrey 2013). The HOME tests concluded that MASH was one of the most successful methods (Domonkos et al. 2012, Domonkos 2013, Venema et al. 2012), that is why we used it in this study.

MASH application is a relative homogeneity test procedure (Szentimrey 1999). This tool consists of mathematical formulation, climatological station information (metadata), and software development for automation. Application does not assume that the reference series are homogeneous. The candidate series is chosen from the available time series (for example daily wind speed data), and the remaining series are considered as reference series. As running the application, the role of series changes step by step during the procedure. Depending on the climatic element, additive (for temperature) or multiplicative (for precipitation or wind speed) models can be used.

It is possible to homogenize monthly, seasonal, or annual time series. The daily inhomogeneities can be derived from the monthly ones (Szentimrey 2008). The application provides automatically the probable dates of break points for further usage, and the homogenized, completed and quality controlled time series. Although MASH is able to use metadata (for example the date of relocation) during the break point detection, it was not used during this work.

In this study, daily wind speed data sets for 19 stations (Fig. 1) were derived from at least 8 hourly wind speed data a day. Before calculating daily wind speed, hourly data was quality controlled and corrected. Metadata of stations is summarized in Table 1. Data are available from 1975 to 2012 at most stations. At station Paks (No. 15), measurements started only on May 1, 1979. Altogether more than one year is missing at Zalaegerszeg (No. 11) during 1993 and 1994. It is also important to note that 50 days are missing at Kecskemét (No. 17) in 2009.

A multiplicative model was applied for homogenization of daily wind speed data using the 0.05 significance level. Original series can be affected by climate change, inhomogeneity, and noise effect (Szentimrey 2011).



**Fig. 1. Hungarian stations used at MASH application for homogenization.**

Table 1. Metadata of Hungarian stations (in 2012) used at MASH application for homogenization

No.	WMO	Station name	Lon [° E]	Lat [° N]	Altitude [m]	Anemometer elevation [m]	Missing data [%]
1	12772	Miskolc	20.77	48.10	232.8	16.25	0.0
2	12805	Sopron	16.60	47.68	233.8	18.40	< 0.1
3	12812	Szombathely	16.65	47.20	201.1	10.56	< 0.1
4	12822	Győr	17.67	47.71	116.7	11.16	0.0
5	12843	Budapest Lőrinc	19.18	47.43	139.1	14.68	< 0.1
6	12851	Kékestető	20.02	47.87	1011.3	25.07	< 0.1
7	12860	Szolnok	20.13	47.17	108.1	10.40	< 0.1
8	12882	Debrecen	21.61	47.49	107.6	10.23	0.1
9	12892	Nyíregyháza	21.89	47.96	142.1	15.98	0.2
10	12910	Szentgotthárd	16.31	46.91	311.7	16.61	0.1
11	12915	Zalaegerszeg	16.81	46.93	240.1	10.40	3.3
12	12925	Nagykanizsa	16.97	46.46	139.8	13.69	0.1
13	12935	Siófok	18.04	46.91	108.2	15.10	0.0
14	12942	Pécs	18.23	46.01	202.8	10.55	0.0
15	12950	Paks	18.85	46.57	97.2	9.80	11.4
16	12960	Baja	19.02	46.18	113.0	10.30	0.1
17	12970	Kecskemét	19.75	46.91	114.0	10.40	0.4
18	12982	Szeged	20.09	46.26	81.8	12.25	< 0.1
19	12992	Békéscsaba	21.11	46.68	86.2	6.50	< 0.1

Today, input data of climate models are the widely used gridded reanalysis fields generated from measured and observed data. Climate model runs driven by reanalysis fields are essential, and provide important knowledge for modern climate research. However, the question arises how reliable are different reanalysis data sets for estimation of wind climate parameters and validation of climate models. Global reanalysis data sets (e.g., ERA Interim) are used in our study that was provided by the European Centre for Medium-Range Weather Forecasts (ECMWF) for researchers and climate modelers (Berrisford et al. 2009).

ERA Interim is remarkably improved compared to the earlier ERA-40 reanalysis data sets (1957–2002) due to data assimilation methods and inclusion of more types of observations, e.g., satellite measurements (Berrisford et al. 2009). In our study, datasets of wind

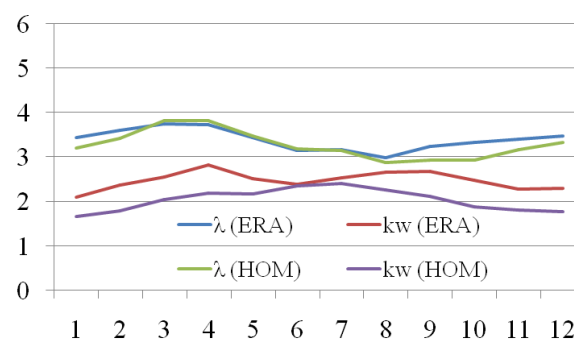
components with fine resolution for the Carpathian Basin (45°–49.5°N and 15°–24°E) are analyzed for 1979–2012.

Homogeneity of 10-meter daily average wind speed of 190 grid points of ERA Interim data sets is checked with MASH 3.03 software (Szentimrey 2011) for the Carpathian Basin between 1979 and 2012. Results of homogenization proved that these gridded data series are homogeneous. Values of the applied test statistics for characterization of inhomogeneity of time series were almost unchanged before and after homogenization and remained under the critical value (20.57; significant level: 0.05) at 72% of grid points. Values of yearly relative estimated inhomogeneity and yearly relative modification of time series differed from zero altogether 15% of the grid points.

### 3. COMPARISON OF DATA SERIES

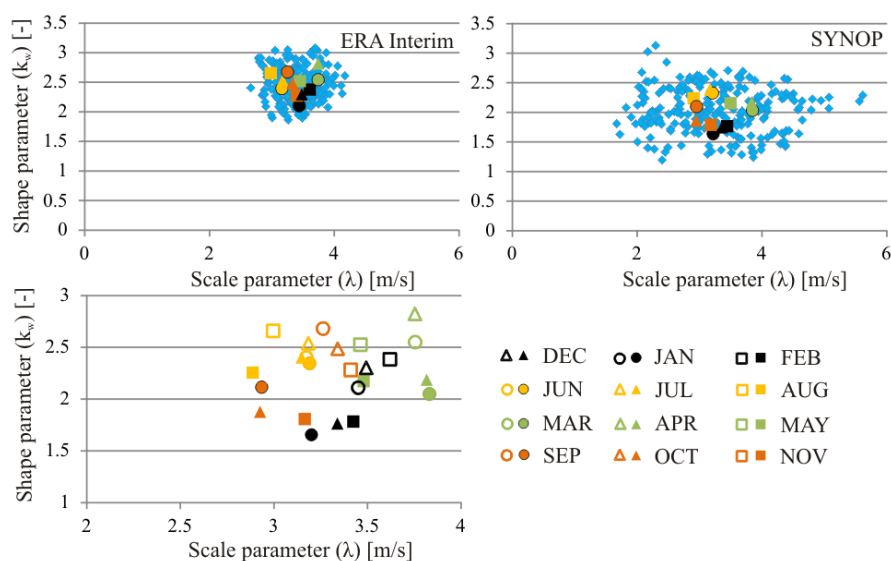
Weibull distributions are fitted in order to compare reanalysis and measured data series. Shape parameter of Weibull distribution ( $k_w$ ) describes frequencies of larger wind speeds. The larger the value of  $k_w$ , the smaller the variability of wind speed. Increasing scale parameter ( $\lambda$ ) when constant shape parameter is assumed occurs as an elongation of probability density function (pdf) along the abscissa with decrease and right-shift of the maxima of pdf (Wilks 2006). Variability of scale parameter is smaller in ERA Interim grid points (3.06–3.83) compared to the synoptic stations (2.13–4.51). Values of Weibull shape parameters of the reanalysis grid points are between 2.10 and 2.65, which are larger than what is found in case of the stations data (1.38–2.16). This overestimation of Weibull shape parameters reduces the variability of wind climate and the probability of extreme wind speed (Rodrigo et al. 2013).

The main disadvantage of homogeneous gridded reanalysis data series is that spatial difference cannot be reproduced by reanalysis data unlike in case of station measurements. Monthly scale parameters of both station and gridded data averages are close in spring and summer when regional differences are compensated (Fig. 2). The monthly average shape parameters are almost equal in June, however, in all the other months overestimations are found at ERA Interim grid points.



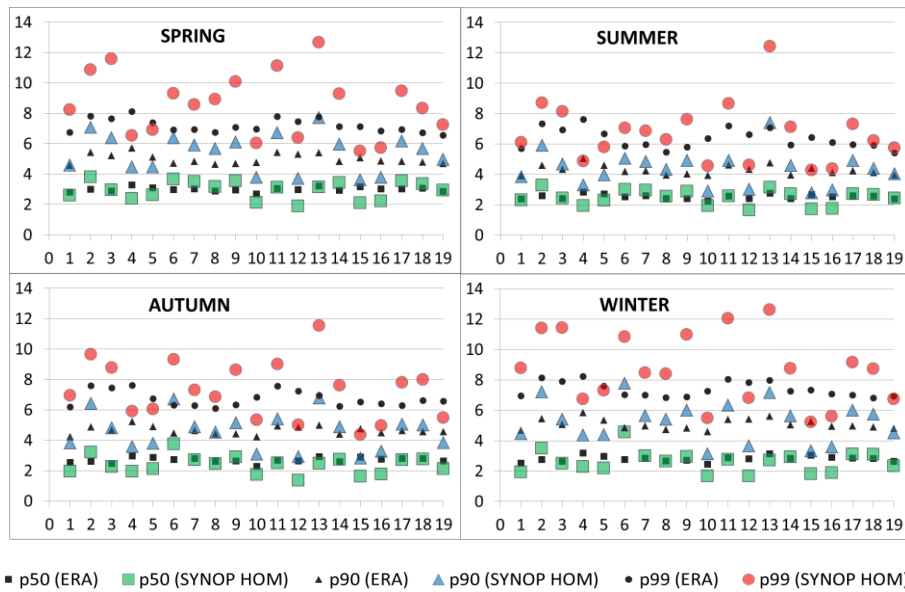
**Fig. 2. Monthly station and gridpoint averages of scale ( $\lambda$ ) and shape ( $k_w$ ) parameters calculated from fitted Weibull distributions of daily wind speed.**

Shape parameters [dimensionless] are shown in Fig. 3 as a function of scale parameter [m/s] of fitted Weibull distributions. Smaller shape parameter can occur in winter due to cyclone activity. Higher scale parameter was found in spring when both the value and the variability of monthly average wind speed are the largest. Average station shape parameters are generally overestimated by the average gridpoint shape parameters, the similar is valid for scale parameter. The only exception occurs in springtime when average station scale parameters are underestimated by the average gridpoint scale parameters. Because the scale parameter depends on wind speed, that is why the wind speed is overestimated, except in spring. The smallest differences (biases) of calculated parameters are observed in June and July.



**Fig. 3. Parameters of Weibull distribution fitted daily wind speed data series of grid points (left up) and stations (right up) in every month (blue) and in different seasons (winter – black, spring – green, summer – yellow, autumn – brown). Monthly grid (unfilled) and station (filled) averages are plotted in the lower diagram.**

Both average and extreme values of homogenized station and grid point data sets were analyzed (Fig. 4), for comparison the nearest grid point of each station is selected. Regarding yearly percentile values, average bias of median of 19 grid points is +11% (minimum: -24%; maximum: +64%). Generally, higher percentile (0.90 and 0.99) values are smaller at reanalysis grid points than at the stations, except near Győr, Budapest, Szentgotthárd, Nagykanizsa, Paks and Baja (Nos. 4, 5, 10, 12, 15, and 16, respectively). Average bias of the 0.90 percentiles of ERA Interim grid points is 2.5% (between -31% and +53%), and the 0.99 percentile's bias is -7.5% (between -39% and +40%).



**Fig. 4. Different seasonal percentiles (square – 0.5, triangle – 0.9, circle – 0.99) values [m/s] for 19 homogenized synoptic stations (colour) and ERA Interim grid points (black) calculated from 34-year time series (1979–2012)**

#### 4. CONCLUSIONS

As a summary of our study, we conclude that discrepancies between station datasets and gridded reanalysis datasets can be explained by the following possible causes. (1) Reanalysis data sets are basically created using limited number of measured data of regular meteorological stations. Hungary, similarly to other countries, reports its synoptic data in a standard code format to ECMWF according to the international cooperation of the member states. However, the number of stations of which data are transmitted fluctuates from year to year. There is only five stations – Miskolc, Budapest, Debrecen, Pécs, Szeged – from which time series are shared continuously starting from the 1979 till today. Since wind varies quite much spatially, when only a few stations with low spatial representation of a specific area are used for interpolation purposes, they result in relatively large errors. In contrast, estimation (interpolation) of atmospheric pressure for Hungary could be successful using only 5-7 Hungarian stations' data, because spatial representation of a specific station is much better in case of atmospheric pressure than wind. (2) Source observational data of reanalysis data sets were not homogenized before interpolation to the grid points, so inhomogeneities of predictor data may decrease correctness of reanalysis data. (3) Regarding gridding interpolation of meteorological elements, it is needed to interpolate not only spatially (i.e., as a GIS problem) but both spatially and temporally in order to create reliable time series in the grid points utilizing climatic knowledge. Reanalysis process is a data assimilation problem using variation analysis and supposing that the background field (where it exists) is equal to the analysis field. This approximation could lead to further sources of biases.

Researchers should make every possible effort to use quality controlled, homogenized and reliable data sets in order to complete reliable wind climatological analysis and to estimate potential renewable wind energy sources.

## Acknowledgement

We are grateful to the Hungarian Meteorological Service and the Geoinformation Service of the Hungarian Defence Forces for the wind data of the Hungarian synoptic meteorological stations and ERA Interim reanalysis data sets. This work was partially supported by the European Union and the European Social Fund through the project FuturICT.hu (grant no.: TÁMOP-4.2.2.C-11/1/KONV-2012-0013). Research leading to this paper has been supported by the Hungarian National Science Research Foundation under grant K-78125.

## References

- Aguilar, E., Auer, I., Brunet, M., Peterson, T.C., and Wieringa, J., 2003: Guidelines on climate metadata and homogenization. WCDMP-No. 53, WMO-TD No. 1186. World Meteorological Organization, Geneva.
- Berrisford, P., Dee, D. P., Fielding, K., Fuentes, M., Kallberg, P., Kobayashi, S., Uppala, S. M., 2009: The ERA-Interim Archive. ERA Report Series No. 1. ECMWF: Reading, UK.
- Emeis, S., 2013: Wind energy meteorology. Atmospheric Physics for Wind Power Generation, Springer, DOI 10.1007/978-3-642-30523-8
- Domonkos, P., Venema, V., and Mestre, O., 2012: Efficiencies of homogenization methods: our present knowledge and its limitation. In Proceedings of the 7th Seminar for Homogenization and Quality Control in Climatological Databases in press, [www.c3.uv.cat/publicacions/publicacions2012.html](http://www.c3.uv.cat/publicacions/publicacions2012.html)
- Domonkos, P., 2013: Measuring performances of homogenization methods. *Időjárás* 117, 91–112
- Freitas, L., Gonzalez Pereira, M., Caramelo, M., Mendes, M., and Nunes, L., 2013: Homogeneity of monthly air temperature in Portugal with HOMER and MASH. *Időjárás* 117, 69–90
- HOME, 2011: Homepage of the COST Action ES0601 - Advances in Homogenisation Methods of Climate Series: an Integrated Approach (HOME), <http://www.homogenisation.org>.
- Péliné N. Cs., Radics K., Bartholy J., 2011: Seasonal variability of Hungarian wind climate. *Acta Silvatica & Lignaria Hungarica*, 7, 39-48.
- Péliné N. Cs., Bartholy J., Pongrácz R., 2014: Homogenization of Hungarian daily wind speed data series. *Időjárás*, Vol. 118, No. 2, 119-132.
- Rodrigo, J. S., Buchlin, J., van Beeck, J., Lenaerts, J. T. M., van den Brooke, M. R., 2013: Evaluation of the Antarctic surface wind climate from ERA reanalyses and RACMO2/ANT simulations based on automatic weather stations. *Climate Dynamics* 40: p353-376., Springer, DOI 10.1007/s00382-012-1252-2
- Szentimrey, T., 1999: Multiple Analysis of Series for Homogenization (MASH). Proceedings of the Second Seminar for Homogenization of Surface Climatological Data, Budapest, Hungary; WMO, WCDMP-No. 41, 27-46.
- Szentimrey, T., 2008: Development of MASH homogenization procedure for daily data, Proceedings of the Fifth Seminar for Homogenization and Quality Control in Climatological Databases, Budapest, Hungary, 2006; WCDMP-No. 68, WMO-TD No. 1434, 116–125.
- Szentimrey, T., 2011: Manual of homogenization software MASHv3.03. Hungarian Meteorological Service, Budapest.
- Szentimrey, T., 2013: Theoretical questions of daily data homogenization, Special Issue of the COST-ES0601 (HOME) ACTION, *Időjárás* 117, 113–122.
- Venema, V., Mestre, O., Aguilar, E., Auer, I., Guijarro, J.A., Domonkos, P., Vertacnik, G., Szentimrey, T., Stepanek, P., Zahradnicek, P., Viarre, J., Müller-Westermeier, G. Lakatos, M., Williams, C.N., Menne, M., Lindau, R., Rasol, D., Rustemeier, E., Kolokythas, K., Marinova, T., Andresen, L., Acquaotta, F., Fratianni,



S., Cheval, S., Klancar, M., Brunetti, M., Gruber, C., Duran, M.P., Likso, T., Esteban, P., and Brandsma, T., 2012: Benchmarking monthly homogenization algorithms. *Climate of the Past* 8, 89–115.

Wilks, D. S., 2006: *Statistical methods in the atmospheric sciences*. International Geophysics Series, Volume 91, Department of Earth and Atmospheric Sciences, Cornell University, Elsevier

World Meteorological Organization, 2011: *Guide to Climatological Practices*, WMO/No. 100, Geneva.

# PROGRAMME

**Budapest, Hungary**  
**12-16 May 2014**

**Venue:**

The Headquarters of the Hungarian Meteorological Service (1 Kitaibel Pál street, Budapest)

Homogenization sessions oral and posters: 12 May Monday-14 May Wednesday  
Interpolation session oral and posters: 15 May Thursday  
Software session: 16 May Friday  
EUMETNET DARE (Data Recovery and Rescue) Team on 13 May Tuesday afternoon

## MONDAY, 12 MAY

8:30 – 9:00 Registration

9:00 – 12:00

Opening addresses by  
President of OMSZ  
Delegate of WMO  
Organizers

Introductory Presentations

Hechler, P., Baddour, O.: Elements of sustained data management solutions for climate

10:00 – 10:30 coffee break

Szentimrey, T., Lakatos, M., Bihari, Z.: Mathematical questions of homogenization and quality control

Lindau, R., Venema, V.: On the reliability of using the maximum explained variance as criterion for optimum segmentations in homogenization algorithms

12:00 – 14:00 Lunch break

14:00 – 17:00 Homogenization and quality control of monthly data

Coll, J., Curley, M., Walsh, S., Sweeney, J.: Homogenising Ireland's monthly precipitation records - an application of HOME-R and statistical exploration protocols to the station network

Curley, M., Walsh, S.: Homogenisation of Monthly Maximum and Minimum Air Temperatures in Ireland

Dubuisson, B., Gibelin, A-L., Jourdain, S., Deaux, N., Laval, L.: Reliable long term series for analysing climate change at Météo-France

15:30 – 16:00 coffee break

Domonkos, P.: The ACMANT2 software package

Yosef, Y.: Homogenization of monthly temperature series in Israel - an integrated approach for optimal break-points detection

18:00 – Welcome party

(Hungarian Meteorological Service, 1 Kitaibel Pál street, Budapest)

## **TUESDAY, 13 MAY**

9:00 – 12:30 Homogenization and quality control of monthly data

Willett, K., Venema, V., Williams, C., Aguilar, E., Lopardo, G., Jolliffe, I., Alexander, L., Vincent, L., Lund, R., Menne, M., Thorne, P., Auchmann, R., Warren, R., Bronnimann, S., Thorarinsdottir, T., Easterbrook, S., Gallagher, C.: Homogenisation algorithm skill testing with synthetic global benchmarks for the International Surface Temperature Initiative

Luhunga, P., M., Mutayoba, M., Ng'ongolo, H., K.: Homogeneity of monthly mean air temperature of the United Republic of Tanzania with HOMER

Zahradníček, P., Rasol, D., Cindrić, K., Štěpánek, P.: Homogenization of monthly precipitation time series in Croatia

10:30 – 11:00 coffee break

Lijuan, C., Ping, Z., Zhongwei, Y., Jones, P., Yani, Z., Yu, Y., Guoli, T.: Instrumental Temperature Series in Eastern and Central China Back to the 19th Century

Dunn, R.: Identifying Homogeneous sub-periods in HadISD

Elfadli, K., Brunet, M.: The WMO/MEDARE Initiative: bringing and developing high-quality historical Mediterranean climate datasets into the 21st century

12:30 – 14:00 Lunch break

14:00 – 17:00 EUMETNET DARE (Data Recovery and Rescue) Expert Team meeting (open for everybody)

15:30 – 16:00 coffee break

**WEDNESDAY, 14 MAY**

9:00 – 12:00 Homogenization and quality control of monthly data

Tayyar, A.: Climate data in Jordan

Djamel, B.: Homogenization of the pluviometric series and the climatic variability in the Northeast region of Algeria

Casabella, N., González-Rouco, J., F., Navarro, J., Hidalgo, A., Lucio-Eceiza, E., E., Conte, J., L., Aguilar, E.: Homogeneity of monthly wind speed time series in the Northeast of the Iberian Peninsula

10:30 – 11:00 coffee break

Guijarro, J., A.: Homogenization of Spanish mean wind speed monthly series

Lucio-Eceiza, E., E., González-Rouco, J., F., Navarro, J., Hidalgo, Á., Jiménez, P., A., García-Bustamante, E., Casabella, N., Conte, J., Beltrami, H.: Quality control of a surface wind observations database for north eastern north America

12:00 – 13:30 Lunch break

13:30 – 16:30 Homogenization and quality control of daily data

Legg, T.: Comparison of daily sunshine duration recorded by Campbell-Stokes and Kipp & Zonen sensors

Venema, V., Aguilar, E., Auchmann, R., Auer, I., Brandsma, T., Chimani, B., Gilibert, A., Mestre, O., Toreti, A., Vertacnik, G., Domonkos, P.: Inhomogeneities in daily data

Acquaotta F., Fratianni, S., Venema, V.: Comparison study of two independent precipitation networks on daily and monthly scale in Piedmont, Italy

15:00 – 15:30 coffee break

Warren, R.: Benchmarking the Performance of Daily Temperature Homogenisation Algorithms

Yuan, F., Tang, G., Wang, X., L., Wan, H., Lijuan, C.: Quality Control and Homogenization of China's 6-hourly Surface Pressure Data

19:00 – Seminar banquet

(Venue: Hungarian Meteorological Service, 1 Kitaibel Pál street, Budapest; 19:00)

(Location of the restaurant: Kaltenberg Étterem, Kinizsi street 30-36, Budapest; 19:30)

**THURSDAY, 15 MAY**

9:00 – 12:00 Spatial Interpolation, Homogenization and Gridding

Szentimrey, T., Bihari, Z., Lakatos, M.: Mathematical questions of spatial interpolation of climate variables

Bertrand, C.: Creation of a 30 years-long high resolution homogenized solar radiation data set over the Benelux

Journée, M.: Gridding of precipitation and air temperature observations in Belgium

10:30 – 11:00 coffee break

Wypych, A., Ustrnul, Z., Henek, E.: Meteorological hazard maps – methodological approach

Petrović, P., Simić, G., Kordić, I.: Practical Aspects of Raw, Homogenized and Gridded Daily Precipitation Datasets

12:00 – 14:00 Lunch break

14:00 – 17:00 Presentations connected with CARPATCLIM project

Skrynyk, O., Savchenko, V., Radchenko, R., Skrynyk, O.: Homogenization of monthly air temperature and monthly precipitation sum data sets collected in Ukraine

Birsan, M-V., Dumitrescu, A.: Homogenization and gridding of the Romanian climatic dataset using the MASH and MISH software packages

Szalai, S., Bihari, Z., Lakatos, M., Szentimrey, T.: The CARPATCLIM (Climate of Carpathian Region) project

15:30 – 16:00 coffee break

Lakatos, M., Szentimrey, T., Bihari, Z., Szalai, S.: Homogenization in CARPATCLIM (Climate of Carpathian Region) project

Bihari, Z., Szentimrey, T., Lakatos, M., Szalai, S.: Gridding in CARPATCLIM (Climate of Carpathian Region) project

**FRIDAY, 16 MAY**

9:00 – 12:00 Software Presentations

Szentimrey, T.: Software MASH (Multiple Analysis of Series for Homogenization)

Stepanek, P.: Software AnClim for tutorial of statistical methods in climatology, including homogenization and ProClimDB for processing of climatological datasets

10:20 – 10:50 coffee break

Domonkos, P.: Software ACMANT2

Szentimrey, T.: Software MISH (Meteorological Interpolation based on Surface Homogenized Data Basis)

## LIST OF PARTICIPANTS

### ALGERIA

BOUCHERF DJAMEL  
National Meteorological Office Algeria  
d.boucherf@meteo.dz

### AUSTRIA

INGEBORG AUER  
Central Institute for Meteorology and  
Geodynamics  
ingeborg.auer@zamg.ac.at

BARBARA CHIMANI  
Central Institute for Meteorology and  
Geodynamics  
barbara.chimani@zamg.ac.at

### BELGIUM

CEDRIC BERTRAND  
Royal Meteorological Institute of Belgium  
cedric@meteo.be  
cedric.bertrand@meteo.be

MICHEL JOURNEE  
Royal Meteorological Institute of Belgium  
michel.journee@meteo.be

### CHINA

FANG YUAN  
National Meteorological Informational  
Center  
yuan-fang-1984@hotmail.com

LIJUAN CAO  
National Meteorological Information  
Center  
caolj@cma.gov.cn

### CROATIA

DUBRAVKA RASOL

Meteorological and Hydrological Service,  
Croatia  
rasol@cirus.dhz.hr

### CZECH REPUBLIC

VÍT KVĚTOŇ  
Czech Hydrometeorological Institute  
vit.kveton@chmi.cz

PETR STEPANEK  
Global Change Research Centre AS CR, v.  
v. i.  
stepanek.p@czechglobe.cz

### ESTONIA

KAIRI VINT  
Estonian Environment Agency  
kairi.vint@envir.ee

### FINLAND

ANNA FREY  
Finnish Meteorological Institute,  
Observation Services  
anna.frey@fmi.fi

### FRANCE

ANNE-LAURE GIBELIN  
Météo-France  
anne-laure.gibelin@meteo.fr  
BRIGITTE DUBUISSON  
Météo-France  
brigitte.dubuisson@meteo.fr

### MACEDONIA

ALEKSANDAR PRODANOV  
Hydrometeorological Service of  
Macedonia  
aprodanov@meteo.gov.mk

## **GERMANY**

KARSTEN FRIEDRICH  
Deutscher Wetterdienst  
karsten.friedrich@dwd.de

RALF LINDAU  
Meteorological Institute of University  
Bonn  
rlindau@uni-bonn.de

VICTOR VENEMA  
Meteorological Institute of University  
Bonn  
Victor.Venema@uni-bonn.de

## **GREECE**

ANNA MAMARA  
Hellenic National Meteorological Service  
annamamara@yahoo.gr

## **HUNGARY**

TAMÁS SZENTIMREY  
Hungarian Meteorological Service  
szentimrey.t@met.hu

ZITA BIHARI  
Hungarian Meteorological Service  
bihari.z@met.hu

MÓNICA LAKATOS  
Hungarian Meteorological Service  
lakatos.m@met.hu

SÁNDOR SZALAI  
Szent István University  
Szalai.Sandor@mkk.szie.hu

TAMÁS KOVÁCS  
Hungarian Meteorological Service  
kovacs.t@met.hu

ENIKŐ VINCZE  
Hungarian Meteorological Service  
vincze.e@met.hu

CSILLA PÉLINÉ NÉMETH  
Geoinformation Service of the Hungarian  
Defence Forces  
pelinenemeth.csilla@mhtehi.gov.hu

## **IRELAND**

JOHN COLL  
Irish Climate Analysis and Research Unit  
john.coll@nuim.ie

MARY CURLEY  
Met Éireann  
mary.curley@met.ie

## **ISRAEL**

YIZHAK YOSEF  
Israel Meteorological Service Climatology  
Department  
yosefy@ims.gov.il

## **ITALY**

FIORELLA ACQUAOTTA  
University of Turin, Earth Science  
Department, NatRisk  
fiorella.acquaotta@gmail.com

## **JORDAN**

AHMAD MAH'D MOH'D TAYYAR  
Jordan Meteorological Department  
tayarcom@yahoo.com

## **LIBYA**

KHALID ELFADLI IBRAHIM  
Libyan National Meteorological Centre  
kelfadli@yahoo.com

## **MONTENEGRO**

MIRJANA SPALEVIC  
Institute of Hydrometeorology and  
Seismology of Montenegro  
mirjana.spalevic@meteo.co.me

## **MOROCCO**

EL GUELAI FATIMA ZOHRA  
Moroccan Meteorological Service  
faty.elguelai@gmail.com

## **POLAND**



AGNIESZKA WYPYCH  
Institute of Geography and Spatial  
Management, Jagiellonian University  
agnieszka.wypych@uj.edu.pl

## **ROMANIA**

MARIUS-VICTOR BIRSAN  
Meteo Romania (National Meteorological  
Administration)  
marius.birsan@gmail.com

## **SERBIA**

GORDANA SIMIĆ  
Republic Hydrometeorological Service of  
Serbia  
gordana.simic@hidmet.gov.rs

IVANA KORDIĆ  
Republic Hydrometeorological Service of  
Serbia  
ivana.kordic@hidmet.gov.rs

PREDRAG PETROVIĆ  
Republic Hydrometeorological Service of  
Serbia  
predrag.petrovic@hidmet.gov.rs

## **SLOVAKIA**

OLIVER BOCHNÍČEK  
Slovak Hydrometeorological Institute  
oliver.bochnicek@shmu.sk

PETER KAJABA  
Slovak Hydrometeorological Institute  
peter.kajaba@shmu.sk

## **SPAIN**

DHAIS PEÑA  
University of Saragossa  
dpang@unizar.es

ETOR EMANUEL LUCIO-ECEIZA  
Universidad Complutense Madrid  
eelucio@fis.ucm.es

JOSÉ A. GUIJARRO

AEMET (Spanish State Meteorological  
Agency)  
jguijarrop@aemet.es

NURIA CASABELLA  
CIEMAT (Centro de Investigaciones  
Energéticas, Medioambientales y  
Tecnológicas) & UCM (University  
Complutense of Madrid)  
nucasabe@ucm.es

PÉTER DOMONKOS  
Centre for Climate Change (C3),  
University Rovira i Virgili, Tortosa, Spain  
peter.domonkos@urv.cat

ENRIC AGUILAR  
CENTER FOR CLIMATE CHANGE, C3,  
URV  
enric.aguilar@urv.cat

## **SWITZERLAND**

RENATE AUCHMANN  
Institute of Geography, University of Bern  
renate.auchmann@giub.unibe.ch

## **TANZANIA**

PHILBERT MODEST LUHUNGA  
Tanzania Meteorological Agency (TMA)  
philuhunga@yahoo.com

## **TUNISIA**

MELIKA NAFFATIA  
Institut National de la Météorologie  
melika@meteo.tn

## **UNITED KINGDOM**

RACHEL WARREN  
College of Engineering, Maths and  
Physical Sciences, University of Exeter  
rw307@exeter.ac.uk

ROBERT DUNN  
Met Office Hadley Centre  
robert.dunn@metoffice.gov.uk

TIM LEGG  
Met Office  
tim.legg@metoffice.gov.uk

## **UKRAINE**

VALERIYA SAVCHENKO  
Taras Shevchenko National University of Kyiv  
savchenkovaleria94@gmail.com

## **WMO**

PEER HECHLER  
Data Management Applications Division  
pechler@wmo.int

For more information, please contact:

**World Meteorological Organization**

**Observing and Information Systems Department**

Tel.: +41 (0) 22 730 82 68 – Fax: +41 (0) 22 730 80 21

E-mail: [wcdmp@wmo.int](mailto:wcdmp@wmo.int)

7 bis, avenue de la Paix – P.O. Box 2300 – CH 1211 Geneva 2 – Switzerland

**[www.wmo.int](http://www.wmo.int)**