

# Homogenization of Temperature Series via Pairwise Comparisons

Matthew J. Menne\* and Claude N. Williams, Jr.  
NOAA/National Climatic Data Center  
Asheville, North Carolina 28801

*Submitted to the Journal of Climate*  
*September 2007*

\*Corresponding author address:  
Dr. Matthew Menne  
NOAA/National Climatic Data Center  
151 Patton Avenue  
Asheville, North Carolina 28801  
Tel. 828-271-4449  
Fax. 828-271-4328  
E-Mail: [Matthew.Menne@noaa.gov](mailto:Matthew.Menne@noaa.gov)

## **Abstract**

An automated homogenization algorithm based on the pairwise comparison of temperature series is described. The algorithm works by forming a matrix of pairwise difference series between serial temperature values from a network of stations. The pairwise difference series are then evaluated for undocumented changepoints and the station series responsible for the breaks are identified automatically. The algorithm can also incorporate station history information to help improve changepoint detection. When the magnitude of a documented or undocumented changepoint is determined to be statistically significant, an adjustment is made for the target series. The algorithm is shown to be robust and efficient at detecting undocumented step changes under a variety of step and trend inhomogeneities. In addition, the pairwise comparison method implemented by the algorithm is shown to yield a lower false alarm rate for undocumented changepoint detection relative to the more common use of a reference series.

## 1. Introduction

Discontinuities in a climate series can be induced by virtually any change in instrumentation or observation practice. The relocation, replacement or recalibration of an instrument, for example, can lead to an abrupt shift in time-ordered observations that is unrelated to any real change in climate. Likewise, gradual alterations to the land use or land cover surrounding a measurement site might induce a “creeping” change (Carretero et al. 1998; Karl et al. 1988) that could limit the degree to which observations are representative of a particular region. Such artifacts in the climate record ultimately confound attempts to quantify climate variability and change (Thorne et al. 2005a). Unfortunately, changes to the circumstances behind a series of climate observations are practically inevitable at some point during the period of record. For this reason, testing for artificial discontinuities or “inhomogeneities” is an essential component of climate analysis. Often, the test results can be used to adjust a series so that it more closely reflects only variations in weather and climate.

Numerous approaches have been employed to detect discontinuities in climate series (Peterson et al. 1998a), and comparison studies have recently proliferated (e.g., Ducré-Robitaille et al. 2003; DeGaetano 2006; Reeves et al. 2007). The goal of this work is to describe an automated homogenization algorithm that builds on the most efficient changepoint detection techniques using a holistic design approach. For example, the algorithm relies upon a pairwise comparison of temperature series in order to reliably distinguish artificial changes from true climate variability, even when the changes are undocumented (Caussinus and Mestre 2004). Consequently, the procedure detects inhomogeneities regardless of whether there is *a priori* knowledge of the date or

circumstances of a change in the status of observations (Lund and Reeves, 2002). In addition, the algorithm employs a recursive testing strategy to resolve multiple undocumented changepoints within a single time series (Menne and Williams, 2005). Lastly, the procedure explicitly looks for both abrupt “jumps” as well as local, unrepresentative trends in the temperature series (DeGaetano, 2006).

Additional background on the design considerations for constructing this “pairwise” homogenization algorithm is provided in section 2. In section 3, the specific components of the algorithm based on these considerations are described. In section 4, an assessment of the algorithm’s skill at changepoint detection and how this skill compares to previous studies is provided by means of simulated temperature series. Because of recent interest in land-use change and its impact on the temperature record (e.g. Peterson and Owen 2005; Kalnay et al. 2006; Parker 2006; Pielke et al. 2007), the algorithm was also applied to historical temperature data from the U.S. Cooperative Observing Network to assess the frequency of local, non representative trends as discussed in section 5. Some concluding remarks are offered in section 6.

## **2. Design considerations for the pairwise algorithm**

### *a. Relative changepoint testing*

Conrad and Pollak (1962) describe the concept of relative homogeneity as follows: *a climatological series is relatively homogeneous with respect to a synchronous series at another place if the temperature differences (or precipitation ratios) of pairs of homologous averages constitute a series of random numbers that satisfies the law of errors.* The assumption is that similar variations in climate occur at nearby locations

because of the spatial correlation inherent to meteorological fields (Livezey and Chen 1983). A statistically significant and persistent violation of relative homogeneity is presumed to be artificial or, at least, to have origins other than the background variations in weather and climate.

Relative homogeneity testing is conducted primarily to distinguish artificial breaks from real climate variability, but it also provides a number of additional benefits. For instance, there is a greater likelihood that a sequence of temperature differences will satisfy the assumption of independently and identically distributed Gaussian errors relative to the original raw series. Likewise, when two monthly temperature series are highly correlated, the variance of their differences will be lower relative to the original series. The reduction in variance improves the power of changepoint detection.

To carry out relative homogeneity testing, a reference series is commonly constructed by averaging values from locations near the target site whose observations are in question (Karl and Williams 1987; Alexandersson and Moberg 1997; Vincent 1998). Unfortunately, the homogeneity of the reference series cannot be taken for granted since undocumented changepoints may be present in any one of the averaged series (Hanssen-Bauer and Førland 1994; Menne and Williams 2005). Strategies for reducing changepoint-attribution errors have included assessing the homogeneity of the reference series itself (e.g., McCarthy et al. 2007) and building a reference from previously adjusted series (e.g., Gonzalez-Ruoco et al. 2001). Unfortunately, conducting a separate assessment of reference series homogeneity fails to exploit the enhanced sensitivity of relative homogeneity testing, and many small amplitude changepoints may go undetected in the reference series only to be later attributed to the target series.

Similar problems may arise when adjusted data are used to build a reference series because artifacts from the original imperfect reference series can be transferred to the adjusted data.

Alternatively, relative homogeneity testing can be implemented via a pairwise comparison of individual climate series (Jones et al. 1986; Slonosky et al 1999; Menne and Duchon 2001; Caussinus and Mestre 2004). In pairwise testing, the cause of undocumented changepoints can be traced more directly, that is, without first testing the reference series or assuming it is homogeneous. Unfortunately, implementing pairwise testing has previously required a manual review of the results. For example, Jones et al. (1986) conducted an arduous station-by-station homogenization by manually determining the cause of changepoints in paired difference series. Caussinus and Mestre (2004) computed the locations of changepoints in difference series automatically, but still deferred to an analyst to attribute the cause. In contrast, an automated approach was developed for the pairwise algorithm as described in section 3.

*b. Distinction between documented and undocumented changepoints*

In the absence of station history records, the date of an inhomogeneity must be treated as an unknown parameter. In such cases, a systematic search through all values in a series is required to identify the dates of statistically significant discontinuities. The systematic nature of the search necessitates the use of a more conservative set of critical values relative to the standard values which are appropriate for testing the significance of known changes to observation practice (Lund and Reeves 2002). This means that tests for undocumented changepoints are less sensitive than comparable tests for documented

changes. It follows that to maximize the power of changepoint detection, station histories should be exploited whenever possible.

The strategy used by the pairwise algorithm is to first identify all evidence of changepoints using the less sensitive tests for undocumented changepoints. Where possible, the results are then combined with information about documented changes whose impact may go undetected by these less sensitive tests. An important benefit of this approach is that all possible changepoints are identified before estimates of their magnitude are made.

*c. Resolving multiple undocumented changepoints*

While the issue of accurately resolving multiple undocumented changepoints remains an active area of statistical research (Reeves et al., 2007), two approaches are in operative use (Menne and Williams 2005). The first, more common approach uses a recursive testing procedure (e.g., Vincent 1998) to overcome the “at most one changepoint” assumption behind most hypothesis tests for undocumented changepoints. The second approach relies on a penalty function to constrain the number of changepoints resolved through an optimization routine used to maximize the contrast between sequential mean levels of a series (e.g., Caussinus and Mestre 2004).

A recursive testing approach is used in the pairwise algorithm for two reasons. First, the approach is associated with a low probability of false changepoint detection without requiring an analyst to interpret the results (cf. Caussinus and Mestre 2004). Second, Menne and Williams (2005) noted that when the recursive hypothesis test method is carried out using a semi-hierarchical splitting algorithm (Hawkins 1976), the

power of changepoint detection can be comparable to that of optimal algorithms.

Recursive testing is based on a hierarchic, binary segmentation of the test series whereby a series is split at the location where the test statistic reaches a maximum, i.e., the point at which the separation between the mean before and after the breakpoint is greatest. Then, the sub-sequences on either side of the first split are likewise evaluated, and the process is repeated recursively until the magnitude of the statistic does not exceed the chosen significance level in any remaining sub-sequences (or the sample size in a segment is too small to test). A semi-hierarchic implementation of this method means that each splitting step is followed by a merging step to test whether a split chosen at an earlier stage has lost its importance after subsequent breakpoints are identified, thereby more closely approximating an optimal solution.

*d. Impact of local, unrepresentative trends*

Ideally, a changepoint detection method should be able to differentiate trend changes from step changes. In practice, however, many of the commonly used undocumented changepoint tests are not robust to the presence of trends in the test data since they are based solely on comparing the means of two sequential intervals. Use of such tests in the presence of trends can lead to falsely detected step changes as well as to inaccurate estimates of the magnitude of a shift when it occurs within a general trend (DeGaetano 2006; Pielke et al. 2007). Conversely, methods that directly account for both step changes and trend changes (e.g., Vincent 1998; Lund and Reeves 2002; Wang 2003) are characterized by much lower powers of detection than the simpler difference in means tests.

While no one test clearly outperforms others under all circumstances, the

Standard Normal Homogeneity Test (SNHT; Alexandersson, 1986) has been shown to have superior accuracy in identifying the position of a step change under a wide variety of step and trend inhomogeneity scenarios relative to other commonly used methods (DeGaetano 2006; Reeves et al. 2007). For this reason, the pairwise algorithm uses the SNHT along with a verification process that identifies the form of the apparent changepoint (e.g, step change, step change within a trend, etc.). The pairwise testing procedure is similar to the Vincent (1998) and Reeves et al. (2007) forward and backward regression methods, respectively, but is more easily adaptable to a recursive testing approach for resolving multiple undocumented changepoints, and at the same time retains the higher power of detection of the SNHT.

### **3. Description of the pairwise algorithm**

The pairwise algorithm begins by selecting neighbors for each target series. A matrix of difference series is then formed between the target series and a number of its neighbors. Each difference series is evaluated for the presence of undocumented changepoints, and the station series responsible for the breaks is identified through an iterative attribution process. The impacts of station history events, if available, are evaluated when adjustments are calculated for all breaks in the series as described below.

#### *a. Selection of neighbors and formulation of difference series*

The pairwise algorithm starts by finding the 100 nearest neighbors for each temperature station within a network of stations. These neighboring stations are then ranked according to their correlation with the target. First differences of monthly anomalies are used to calculate the correlation coefficients in order to minimize the

impact of artificial shifts in determining the correlation (Peterson et al., 1998b). A series must simply be positively correlated with the target series to be eligible as a neighbor.

From all eligible neighbors, the set used for the pairwise analysis is selected using a two-step process. First, an account is made of the years and months for which both the target and its 40 most highly correlated neighbors report monthly mean maximum and minimum temperature data. Then, beginning with the 41<sup>st</sup> most highly correlated neighbor, the algorithm assesses whether an additional neighbor adds any data for the years and months which have fewer than seven viable neighbors. If the neighbor in question provides records for such data sparse periods, it replaces the least correlated of the original 40 with the new neighbor provided that the addition does not remove data for other data sparse periods. This process ensures that, whenever possible, at least seven neighbors are available at all times during the target station's period of record (the rationale for attempting to make at least seven target-neighbor comparisons is provided in section 4).

Next, time series of differences,  $\{D_i\}$ , are then formed between all target-neighbor monthly temperature series. As an example, Fig. 1 depicts  $\{D_i\}$  series formed between mean monthly maximum temperature anomalies from Chula Vista, California and nine highly correlated neighbor series. The reduction in variance of the  $\{D_i\}$  series relative to the original target series is clearly evident. The variety of overlapping periods and relative breakpoints between the records from Chula Vista and its neighbors is common in surface temperature records.

#### *b. Identification of undocumented changepoints*

Once the matrix of difference series has been formed, the SNHT is used to

identify undocumented changepoints in each  $\{D_t\}$  using the semi-hierarchical splitting algorithm and a 5% significance level ( $\alpha = 0.05$ ). The SNHT evaluates the null hypothesis ( $H_0$ ) that the  $\{D_t\}$  series has a constant mean against the alternative hypothesis ( $H_A$ ) that there is an undocumented step change on date  $c$ . To account for the possibility of multiple changepoints, the difference series is assumed to consist of  $K$  segments, each bounded by two changepoints ( $c_{k-1}$  and  $c_k$ ). In the pairwise algorithm, SNHT takes the form

$$H_0 : \{D_t\} \rightarrow N(\mu_k, \sigma^2), \quad c_{k-1} + 1 \leq t \leq c_k \quad (1)$$

$$H_A : \begin{cases} \{D_t\} \rightarrow N(\mu_1, \sigma^2), & c_{k-1} + 1 \leq t \leq c \\ \{D_t\} \rightarrow N(\mu_2, \sigma^2), & c + 1 \leq t \leq c_k \end{cases} \quad (2)$$

where  $N(\mu, \sigma^2)$  refers to a random variable with a mean  $\mu$  and variance  $\sigma^2$  and  $\mu_1 \neq \mu_2$ .

For convenience we define  $c_0 = 1$  and  $c_K = n$ , the total number of values in the  $\{D_t\}$  series. The un-subscripted  $c$  in (2) refers to the assignment of an undocumented changepoint between two previously established changepoints ( $c_{k-1}$  and  $c_k$ ) as the semi-hierarchical splitting algorithm iterates through the succession of splitting and merging steps, ultimately converging on a solution of  $K$  segments bounded by  $K-1$  changepoints.

*c. Classification of breakpoints identified by the SNHT test*

The result of step (b) is a set of  $K-1$  apparent changepoints for each  $\{D_t\}$  series.

Because the SNHT assumes that each series is of the form

$$\{D_t\} = \mu_k + \varepsilon_t, \quad c_{k-1} + 1 \leq t \leq c_k, \text{ and } k = 1, K, \quad (3)$$

the next step determines whether this piecewise stationary model is justified for each changepoint. The determination is made by fitting a hierarchy of potential models for all

segments centered on each  $k^{\text{th}}$  breakpoint. The five models, M1 through M5, are described in Table 1 (after Reeves et al. 2007). The model that minimizes the Bayesian Information Criterion (BIC; Schwarz, 1978) is selected as the best representation for each changepoint.

Procedurally, the BIC is calculated by fitting M1 through M5 to every segment  $c_{k-1} + 1$  to  $c_{k+1}$  for all  $k=1, K$ . The BIC is defined as

$$\text{BIC}(p) = -2 \log(L) + \log(n')p, \quad (4)$$

where  $p$  is the number of parameters required to fit the model,  $n'$  is the number of data points in the segment  $c_{k-1} + 1$  to  $c_{k+1}$ , and  $L$  is the likelihood of the model in question.

For the models listed in Table 1,

$$-2 \log(L) = n' \log(\text{SSE} / n'), \quad (5)$$

where SSE refers to the sum of squared errors for the particular model fit.

In some cases, one or more of the original  $K - 1$  changepoints may be eliminated from the solution for a particular  $\{D_t\}$  series. For example, if the true model between the values  $c_{k-1} + 1$  and  $c_{k+1}$  is a constant increasing trend (M2), the SNHT may have identified an apparent jump in the middle of the trend interval whereas the BIC is likely to be lower for M2 than for any of the other four models. In such a case, the false changepoint is removed from the original undocumented changepoint solution. Alternatively, the use of the BIC may determine that the  $\{D_t\}$  segment between  $c_{k-1} + 1$  and  $c_{k+1}$  more appropriately follows M4 (step change within a constant trend) or M5 (a step change separated by different trends). If so, there is evidence of a relative trend between the two series and the magnitude of the step change,  $\Delta$ , required in subsequent

steps (e) and (f) should be calculated using the higher dimension models (M4, M5) to avoid calculating a biased estimate of the step.

*d. Attribution of changepoints*

Given that breaks in a difference series will be induced by discontinuities in either of the original temperature series, the next step is to identify the series responsible for a particular discontinuity. To begin, a matrix of change dates by station is formed, and all changepoint dates detected in the  $\{D_i\}$  series are temporarily assigned to both of the original series used to form the differences. Specifically, a count is incremented for the date of change,  $i$ , each time a station,  $j$ , is implicated by a break in one of its difference series. The resulting  $I \times J$  matrix of change dates by station is then “un-confounded” by systematically identifying those stations that are common to numerous difference series with the same date of change. More specifically, the station/date with the highest overall changepoint count is identified. This station is then tagged as the perpetrator, that is, as the cause of the breaks on the date with the highest breakpoint count. The corresponding count on that particular change date is then decremented for all of the perpetrator’s neighbors, and the process is repeated using the updated change-date tallies. The procedure continues recursively until no station/change date count is greater than one for any station-date in the period of record.

*e. Assignment of undocumented changepoint dates*

Although undocumented changepoints are assigned to a perpetrating series in step (d), the date of an undocumented changepoint returned by the SNHT is subject to some

sampling variability. As illustrated in Fig. 2, the degree of this sampling variability is a function of the magnitude of changepoint, with larger changepoints associated with more precise estimates of the date of change. This means that testing a group of target-neighbor difference series often leads to a range of undocumented changepoint dates clustered around the time of the actual change.

To determine which change dates likely refer to the same discontinuity, an interim estimate of step-change magnitude is necessary. The estimate is calculated using the most appropriate change model (M3, M4 or M5) according to the BIC and is used to calculate confidence limits for the change dates returned by the SNHT. The cluster of dates falling within overlapping confidence limits is then conflated to a single change date at the target in one of two ways: 1) it is assigned to the date of a known event in the target station's history that occurs within the confidence limits for the change dates; or 2) it is assigned to the most common changepoint date that falls within the confidence limits, which means that the discontinuity appears to be truly undocumented.

*f. Calculation of adjustments*

Steps (a) through (e) are necessary simply to identify undocumented changepoints in all temperature series. In many applications, however, station histories also may be available, which might provide additional information regarding possible discontinuities. When available, the dates of documented events should be combined with evidence of undocumented changepoints because the impact of documented events may be too subtle for the tests for undocumented changepoints to detect. Adjustments can then be calculated for all undocumented and documented changepoints.

Adjustments are determined by calculating multiple estimates of relative changepoint magnitude using segments from neighboring series that are homogeneous for at least 24 months before and after the target changepoint. In addition, when two changepoints occur within 24 months in the target series, an adjustment is made for their combined effect. The range of pairwise estimates for a particular target step change is considered to be a measure of the confidence with which the magnitude of the discontinuity can be estimated. As in step (e), the step model found to be most appropriate (i.e., M3, M4, or M5) according to the BIC is used to calculate a final estimate of the step change for each relevant  $\{D_t\}$  segment. At least three separate pairwise estimates of step-change magnitude are required for each target changepoint because the distribution of estimates is used to determine the significance of the adjustment. Moreover, since the distribution of step-change estimates is not necessarily symmetric, the median estimate is used to adjust the target series.

Statistical significance is determined by comparing the median estimate of step-change magnitude to the 5<sup>th</sup> percentile (median > 0) or to the 95<sup>th</sup> percentile (median < 0) of all estimates, subject to an initial outlier check. Because fewer than 20 estimates may be available for any given changepoint, a multiple of the difference between the median and the first quartile ( $Q_1$ ) or between the median and third quartile ( $Q_3$ ) serves as an estimate of the 5<sup>th</sup> or 95<sup>th</sup> percentile, respectively. A factor of 2.5 is used because it approximates a one-tailed test at the 5% ( $\alpha=0.05$ ) significance. When the median and the tail of the distribution closest to zero are of the same sign (i.e., median -  $Q_1*2.5$  or median +  $Q_3*2.5$ ), the step-change is considered to be significant, and an adjustment is

made to the target series. This significance test is similar to the Tukey outlier test, but allows for asymmetry in the distribution of estimates.

*g. Example of changepoint detection and adjustment*

Application of the pairwise algorithm to the group of series shown in Fig. 1 revealed two significant changepoints in Chula Vista maximum temperatures, both of which were associated with documented station moves, first on 1 January 1982 and then again on 25 April 1985. Difference series between the pairwise-adjusted mean monthly maximum temperatures for Chula Vista and its neighbors are shown in Fig. 3, which suggests that the algorithm has removed the major step inhomogeneities from all series in the group.

#### **4. Evaluation of the algorithm**

To evaluate the performance of the pairwise algorithm more generally, temperature series were simulated under a number of trend and step change scenarios. The simulations were designed to test the skill of changepoint detection as well as to facilitate comparison of the results to previous investigations regarding the use of a reference series and the identification of the type of changepoint.

*a. Evaluation under monthly temperature simulations*

The performance of the pairwise algorithm was first evaluated using two different sets of simulated monthly temperature anomalies. One set was comprised of series with step changes while the second set contained series with both trend and step

inhomogeneities. Both sets consisted of 1000 groups of 21 correlated “red noise” series generated as described in Menne and Williams (2005; hereafter MW05). The average correlation between each series within a group was about 0.7. For all series the mean ( $\mu$ ) was zero and the standard deviation ( $\sigma$ ) was one; the number of values in each series ( $n$ ) was equal to 1200, the equivalent of 100 years of monthly means.

A random number of step changes was imposed on each series at random dates. The number of steps per series varied normally about an average of five, with as few as zero and as many as ten. The magnitude of each step change was also assigned randomly by sampling from the standard normal distribution, which means that about two-thirds of the imposed steps were equal to one  $\sigma$  or less. As discussed in MW05, the standard normal distribution is a good proxy for the distribution of known impacts to U.S. temperature series (Karl and Williams, 1987). All imposed step changes were treated as undocumented, and 10 neighbors were identified by the pairwise algorithm for all 21 series in the groups.

In the “monthly steps and trends” simulations, a trend inhomogeneity was added to roughly a quarter of the simulated series. The magnitude of this trend was varied randomly from  $0.001\sigma/\text{month}$  up to about  $0.18\sigma/\text{month}$ , while the trend interval varied randomly from 2 months up to the full period of record. In general, the trend inhomogeneity did not initiate with a step change, although steps frequently occurred randomly within the intervals of a creeping inhomogeneity.

Fig. 4 illustrates the impact of random step-only changes on the first group of simulated series. Prior to imposing step changes, the true trend in each series was zero. After imposing step changes, the trends ranged from  $-7.62\sigma/\text{century}$  to  $+4.34\sigma/\text{century}$ .

The pairwise algorithm correctly identified 34 of the 43 imposed step changes. Of the nine step changes not identified, six had a magnitude of less than  $0.3\sigma$ , which is below the sensitivity of most tests for undocumented changepoints (DeGaetano, 2006; Ducré-Robitaille et al. 2003). Furthermore, the largest undetected changepoint ( $+0.696\sigma$ ) was preceded 10 time steps earlier by another undetected changepoint of  $-0.451\sigma$ , i.e., the two changepoints essentially masked one another. The overall effectiveness of the pairwise adjustments is evident in Fig. 5, which depicts the ten series after homogenization by the pairwise algorithm. Note that changepoints have been adjusted relative to the latest mean level in each series, the convention in climate data homogenization. In general, the adjusted series all have trends much closer to the true “climate” trend of zero.

Table 2 more generally summarizes the detection skill of the pairwise approach for both the step-only and the step/trend change scenarios. The hit rate (the ratio of the number of changepoints correctly identified relative to the total number imposed) is roughly 67% for both scenarios. The false alarm rate (the ratio of falsely detected changepoints to the total number detected) is 6.77% for the step-only scenario (only slightly higher than the expected type I error rate at the  $\alpha=0.05$  significance level) and 19.65% for the step/trend scenario. The increase in false alarms when trend inhomogeneities are present occurs for two main reasons. First, the beginning or end of a trend inhomogeneity is often identified as a step change by the pairwise algorithm. Second, short interval trends of about 24 months or less tend to be virtually indistinguishable from step changes and are therefore adjusted as an abrupt change. Indeed, the largest magnitude false alarms under the steps-and-trend inhomogeneity

simulations result from short-interval, but large magnitude trend inhomogeneities that are approximated by a step change.

Histograms indicating the magnitude of hits, misses and false alarms for the step-only and step/trend simulations are shown in Fig. 6 and 7, respectively. In both cases, changes in excess of  $0.5\sigma$  are readily detected, and most misses are generally less than  $0.5\sigma$ . The number of false alarms is also generally small, suggesting that they will have little impact on the homogenized trends for the simulated series.

Regarding the series trends, two measures of error are provided in Table 2. The first is the root mean square error (RMSE) for a trend calculated using the unadjusted series and the second is the RMSE for trends calculated using the adjusted series. As shown in the table, the pairwise homogenization process greatly reduces the error associated with the calculation of the true background climate trend. Table 2 also indicates that the RMSE for changepoint estimates in series with trends about as good as in the series with no trend inhomogeneities, which suggests that the model identification is reasonably successful at identifying step changes that occur within local trends. A more thorough assessment of changepoint type identification is provided in section 4c.

*b. Pairwise versus reference series changepoint detection skill*

The use of a reference series is the most widely employed approach to relative changepoint detection, and MW05 evaluated the implications of such an approach for undocumented changepoint detection. The pairwise approach was therefore evaluated using the same simulations and scenarios as in MW05 to directly compare its skill of undocumented changepoint detection against the reference series approach. Table 3 depicts the seven scenarios evaluated in MW05. Each case was comprised of 1000

groups of six correlated series (one target and five neighbors) with  $n = 100$  values. Of the three reference series formulations evaluated by MW05, the one based on a correlated weighted average of the five neighbors (Alexandersson and Moberg, 1997) is compared here. As in the pairwise algorithm, the SNHT was used to test the target minus weighted-average reference  $\{D_t\}$  series ( $\alpha=0.05$ ). All changepoints detected in the  $\{D_t\}$  series were attributed to the target series to test the consequences of assuming reference series homogeneity.

Table 4 summarizes the pairwise and reference series detection skill for the MW05 target series. Two statistics are presented for each case: the FAR (previously described) and the CorRect Changepoint (CRC) power statistic (Reeves et al. 2007), which is the percentage of time that (a) the changepoint date in the target series was selected within  $\pm 2$  time steps of the correct date or (b) the target was correctly identified as homogeneous. Basically, the CRC is synonymous with hit rate except that it also credits the number of times that the target series was successfully identified as homogeneous.

In general, the pairwise algorithm has a much higher success rate in identifying homogeneous target series than the reference series approach as indicated by the higher CRC percentages for cases 1, 3, and 5. This is true when the neighbor series are themselves homogeneous as in Cases 1 and 5, but especially when all the neighbors have changepoints as in Case 3, which cause numerous inhomogeneities in the reference series. More generally, Table 4 indicates the degree to which the pairwise approach limits the number of false alarms whenever the neighboring series are impacted by

undocumented changepoints as evidenced by the low FAR for Cases 3, 4, and 6 relative to the reference series approach.

As shown in Fig. 8, the pairwise hit rate meets or exceeds that of the reference series approach when there are at least seven viable neighbors available at all times during a target station's history. (This is the foundation for the number of neighbors selected for comparison as described in section 3a). The relatively steep increase in the power of detection as the number of comparisons increases illustrates an advantage of pairwise testing, namely, that there are multiple chances to detect a changepoint in any particular target series. If the SNHT misses a changepoint in one target-neighbor difference series, or if it misidentifies the date, there are a number of additional chances to test for the same undocumented break. The chances are not completely independent, however, because any two  $\{D_i\}$  series with a common target will have an expected correlation of 0.5 (Menne and Duchon 2001). Moreover, the power of pairwise detection can be further improved by increasing the sample size between changepoints, which can be achieved by testing serial monthly values rather than annual or seasonal averages. This accounts for the higher hit rate in the "monthly" simulations, i.e., 67% (Table 2) compared to the rate of a little less than 50% shown in Fig. 8 when 10 neighbors are available.

*c. Skill in identifying the type of changepoint*

The magnitude of a step change will not be accurately estimated if the type of changepoint has been misidentified. Consequently, the skill of the pairwise algorithm in classifying changepoint type was assessed for the range of models in Table 1. As in section 4b, a set of 1000 groups of target and neighbor series with  $n = 100$  values were

used for each scenario. In this case different magnitudes of trend and step parameters, i.e.,  $c, \Delta, \beta, \beta_1, \beta_2$ , were imposed on the target series as shown in Table 5; the five neighbor series, in contrast, were always homogeneous (M1). The magnitudes of the parameters imposed on the target series were the same as those used by Reeves et al. (2007; hereafter R07), although only a portion of the results are summarized here.

A comparison of the CRC's in Table 5 for the  $\Delta=1\sigma$  simulations indicates that the pairwise algorithm correctly identified more than 85% of these step changes regardless of whether the target series followed M3, M4 or M5. Moreover, the algorithm also correctly identified more than 85% of the M2 (constant trend) target series as homogeneous (no steps). On the other hand, there is more variability in the skill of classifying the type of changepoint as indicated by the correct type percentages shown in bold. The percentages indicate that the algorithm had somewhat less success in classifying M4- and M5-type changepoints relative to M3-type changepoints and series that follow M2.

Under the M2 scenarios, the pairwise algorithm correctly classifies more than 85% of the  $\{D_t\}$  series when  $\beta$  is greater than 0.01 (a slope yielding a change of  $1\sigma$  in 100 time steps), but less than 50% when  $\beta=0.005$  (a change of  $0.5\sigma$  in 100 time steps). The reason for the difference is that the BIC does not always distinguish a sloped line from a flat line when  $\beta$  is small. This kind of misclassification, however, does not impact the CRC since there is no step change assigned to the target. On the other hand, when  $\beta$  is larger, the SNHT tends to partition the  $\{D_t\}$  trend into one or more step-type changes. The BIC correctly reclassifies most of these breaks as M2, but also cannot always distinguish a trend (M2) from a step change (M3, M4 or M5). Consequently, the

pairwise algorithm classifies only 91% of M2 target series as homogeneous (no step) when  $\beta=0.01$  and 86.9% when  $\beta=0.02$ . The impact of this type of misclassification is to inadvertently remove some of the unique target series trend as a step adjustment, thereby bringing the target series more in line with the regional background climate trend captured by the neighbors (DeGaetano, 2006; Pielke et al. 2007).

For target series under M3 (step change with no trend), the overall power of detection is a function of the magnitude of the step, as shown in previous investigations (e.g. DeGaetano 2006). In the pairwise algorithm, most ( $> 88\%$ ) of the  $\{D_t\}$  series with a step change of  $1\sigma$  or greater were correctly identified as M3, and the CRC exceeds 90% in such cases. On the other hand, many (about 45%) of the  $0.5\sigma$  magnitude step changes are misclassified as a trend change (M2).

When the target series follows M4 (step change within a constant trend), the pairwise CRC varies between 85 and 90% for the  $1\sigma$  step changes, close the M3 rate. However, in the M4 simulations, the algorithm frequently (about 80% of the time when  $\beta=0.005$ ) misclassifies the  $\{D_t\}$  series as M3, especially when  $\beta$  is small. This type of misclassification leads to a biased estimate of the magnitude of the jump by aliasing the unique target trend on to the estimate of the step change. Much like a false alarm when the target follows M2, the biased estimate would bring the adjusted target more in agreement with the background trend captured by the neighbors (DeGaetano 2006; Pielke et al. 2007).

Under M5, the target series has a step change within a trend change, but there is also a change in trend coincident with the step. In this scenario, the CRCs are comparable to the M4 simulations, but in this case, the pairwise algorithm tended to

misclassify the  $\{D_t\}$  series as M3 or M4 in roughly equal proportions. Consequently, some of the target series trends would be aliased onto the estimate of the M5 step changes, as in the case of the M4 target series simulations.

Overall, the results in Table 5 are consistent with the changepoint type identification capabilities of the generalized methods investigated by R07, namely, that it is more challenging to classify M4 and M5 type changepoints. As shown in R07, the lower identification skill occurs even when changepoint tests specifically designed for these types of change are used, i.e., Wang (2003) for M4 and Lund and Reeves (2002) for M5. Nevertheless, from Table 5 and results (not shown) based on directly testing a target series as in R07, it appears that the pairwise approach (SNHT plus BIC) has comparable skill at model identification compared to the methods evaluated by R07. The advantage of the pairwise approach is that the SNHT's superior power of detection is exploited.

The skill of identifying changepoint type, like the power of detection, can also be improved by increasing the sample size of the test series, i.e., by testing serial monthly series. For example, the percentage of correctly identified M4 difference series is about 70% at  $\beta=0.02$  when  $n=240$  and  $c=120$  versus 50% for  $n=100$  and  $c=50$ . Similarly, when  $\beta_1=0.01$  and  $\beta_2=0.03$  under M5, the percentage of series correctly identified increases to 75% for  $n=240$  versus about 50% for  $n=100$ . In addition, the skill of changepoint detection and identification increases with increasing correlation between series, which reduces the variance of the  $\{D_t\}$  series. As noted by DeGaetano (2006), the correlation between temperature series in the United States is typically higher than in the simulations used here.

## 5. Application to U.S. temperature series

A number of recent studies have focused on the impact of land-use change on the temperature record (e.g. Peterson and Owen 2005; Kalnay et al. 2006; Parker 2006; Pielke et al. 2007), yet no general assessment of the frequency of the various types of changepoints in observed temperature series has been conducted. For this reason, the pairwise algorithm was applied to monthly temperature series from the U.S. Cooperative Observer (Coop) Network in order to assess relative frequency of changepoint types, including local trends, in U.S. temperature records. Monthly mean maximum and minimum values from over 7000 stations covering the period 1895 to 2006 were used, although the specific period of record varied from station to station.

An analysis of the more than 1,000,000  $\{D_i\}$  series segments used to calculate the adjustments for all Coop temperature series indicates that about 50% of the step changes follow M3 (step change with no trend), while approximately 40% follow M5 (step change accompanied by a trend change) and about 10% follow M4 (step change within a general trend). In other words, trend inhomogeneities are as widespread as step changes in the Coop network.

To evaluate the pairwise adjustments for these changes, the adjusted series for a commonly used subset of the Coop network, i.e., the U.S. Historical Climatology Network (HCN; Easterling et al. 1996) were manually inspected. In brief, this entailed graphing each HCN series and its Coop neighbors as in Fig. 3, then subjectively deeming the adjusted series as plausible or implausible. This subjective evaluation revealed that roughly 15 to 20% of the adjusted series exhibited physically unrealistic trends that were clearly inconsistent with neighboring stations. The minimum temperature series at Cheesman, Colorado is an extreme example. As shown in Fig. 9, a saw tooth pattern is

evident in the  $\{D_t\}$  series formed between the Cheesman series and its neighbors. The increasing difference between Cheesman and surrounding stations (particularly after 1980) sometimes exceeded 4°C in five years, a relative change that was easily classified as M5 (step change with a trend change) by the pairwise algorithm. The consequence of adjusting the series using M5 (i.e., removal of the step and retention of the trend) is shown in Fig. 10. The result is clearly unrealistic.

Given that preserving local trends (i.e., trend inhomogeneities) can often result in undesirable adjusted series, the pairwise algorithm was modified to employ the more commonly used M-3 adjustment for all step changes (DeGaetano 2006). (Note that M3, M4, and M5 were still employed to detect step changes). The impact of the M3-only approach on the Cheesman series is also shown in Fig. 10. Although the saw tooth signature remains in the adjusted data, the trend at Cheesman using the M3 adjustments is clearly in sync with the average of trends in surrounding series. A similar visual inspection of all HCN temperature series suggests that an M3-only adjustment approach works well for all situations in which there is evidence of a step change because any associated trend inhomogeneity is consistently aliased onto estimates of the step change in a way that favors the background climate signal.

The same result occurs when M3 alone is used to adjust the simulated series in the “monthly steps and trends” simulations, as shown in Table 6. From a comparison of the RMSE for the adjusted trends in Tables 2 and 6, it is evident that using M3 for all step change adjustments removes the impact of most trend inhomogeneities since the error for the adjusted trends is roughly the same for the steps only and steps and trends simulations. Still, while the temperature series that result using the M3 only adjustments

arguably approximate the best theoretical climate series for each location, the local trend signal is nevertheless aliased out of the original series, thus limiting the use of the adjusted series in some attribution studies of observed temperature change. Ultimately, a better solution would be to remove trend inhomogeneities via trend adjustments and step inhomogeneities via step adjustments. Unfortunately, unlike step changes that occur at the same time within a group of target/neighbor  $\{D_i\}$  series, a trend inhomogeneity at a given target station may begin and end at different times with respect to each of its neighbors. This makes identifying the true interval of trend inhomogeneity more difficult than detecting step changes, and is beyond the scope of this paper. Nevertheless, the pairwise algorithm provides a major first step in identifying intervals of trend inhomogeneities.

## **6. Conclusion**

Our evaluation of the pairwise algorithm suggests that it is a robust, reliable and accurate approach to detecting step-type inhomogeneities under a wide variety of circumstances. Relative to the more traditional use of a climate reference series, a pairwise approach to undocumented changepoint detection reduces the number of false alarms in general and is particularly successful at identifying homogeneous series. In addition, unlike the reference approach, there are no requirements for a group of series to have a common base period. As a result, the estimation of step-change magnitude is not confined to the shortest homogeneous interval within a group of neighboring series. In this regard, the pairwise method is similar to the graph theory approach used by Christy

et al. (2006) except that the pairwise algorithm makes no attempt to compare climate series that do not overlap in time.

Moreover, since each climate series is paired with a unique set of neighboring series in the algorithm, it is possible to determine whether more than one nearby station series shares a particular change date because both stations will have been implicated multiple times on or about the same date. This property of the algorithm is important when a widespread and near simultaneous change in observation practice occurs in a network. Such a situation arose in the U.S. Cooperative Network when liquid in glass thermometers were replaced with electronic thermistors at roughly two-thirds of sites during the mid- and late 1980s (Quayle et al. 1991; Hubbard and Lin 2006). Of course, if a change is implemented on exactly the same date at all stations, relative homogeneity testing will not be effective.

Results from applying the pairwise algorithm to observed temperature series suggest that there is widespread evidence of relative trends between series in the U.S. surface temperature record. Although there is some interest in preserving such trend inhomogeneities for land use/land change impact studies (e.g., Pielke et al 2007), the results of this analysis indicate that physically implausible trends can result when all trend inhomogeneities are preserved. On the other hand, if the goal is to produce an accurate estimate of the background climate signal, all step changes can nevertheless be removed using the step only model. While this necessarily leads to the aliasing of any trend inhomogeneity onto the estimate of the step change, a reliable estimate of the background climate signal is obtained.

## **Acknowledgements**

Special thanks to Russell Vose, Imke Durre and Tamara Houston for constructive comments on earlier drafts for this manuscript. Thanks also to Xiaolan Wang and Lucie Vincent for helpful comments on the pairwise methodology. Partial support for this work was provided by the Office of Biological and Environmental Research, U.S. Department of Energy (Grant number DE-AI02-96ER62276).

## References

- Alexandersson, H.: 1986: A homogeneity test applied to precipitation data. *J. Climatol.*, **6**, 661-675.
- Alexandersson, H., and A. Moberg, 1997: Homogenization of Swedish temperature data. Part I: Homogeneity test for linear trends. *Int. J. Climatol.*, **17**, 25-34.
- Carretero J.C., M. Gomez, I. Lozano, A.R. de Elvira, O. Serrano, K. Iden, M. Reistad, H. Reichardt, V. Kharin, M. Stolley, H. von Storch, H. Gunther, A. Pfizenmayer, W. Rosenthal, M. Stawarz, T. Schmith, E. Kaas, T. Li, H. Alexandersson, J. Beersma, E. Bouws, G. Komen, K. Rider, R. Flather, J. Smith, W. Bijl, J. de Ronde, M. Miletus, E. Bauer, H. Schmidt, H. Langenberg (The WASA Group), 1998: Changing waves and storms in the northeast Atlantic? *Bull. Amer. Meteor. Soc.*, **79**, 741-760.
- Caussinus, H., and O. Mestre, 2004: Detection and correction of artificial shifts in climate series. *J. Royal Stat. Soc, Series C*, **53**, 405-425.
- Christy J.R., W.B. Norris, K. Redmond, and K.P. Gallo, 2006: Methodology and results of calculating central California surface temperature trends: Evidence of human-induced climate change? *J. Climate*, **19**, 548-563.
- Conrad, V., and L.W. Pollack, 1962: *Methods in Climatology*. Harvard University Press, 459 pp.
- DeGaetano, A.T., 2006: Attributes of several methods for detecting discontinuities in mean temperature series. *J. Climate*, **19**, 838-853.
- Ducré-Robitaille, J.-F., L.A. Vincent and G. Boulet, 2003: Comparison of techniques for detection of discontinuities in temperature series. *Int. J. Climatol.*, **23**, 1087-1101.
- Easterling, D. R., T. R. Karl, E.H. Mason, P. Y. Hughes, and D. P. Bowman. 1996. United States Historical Climatology Network (U.S. HCN) Monthly Temperature and Precipitation Data. ORNL/CDIAC-87, [NDP-019/R3](#). Carbon Dioxide Information Analysis Center, Oak Ridge National Laboratory, U.S. Department of Energy, Oak Ridge, Tennessee.
- Gonzalez-Rouco, J.F., J.L. Jimenez, V. Quesada, and F. Valero, 2001: Quality control and homogeneity of precipitation data in the southwest of Europe. *J. Climate*, **14**, 964-978.
- Hanssen-Bauer, I., and E.J. Førland, 1994: Homogenizing long Norwegian precipitation series. *J. Climate*, **7**, 1001-1013.

- Hawkins, D.M., 1976: Point estimation of the parameters of a piecewise regression model. *Appl. Statist.*, **25**, 51-57
- Hubbard, K.G., and X. Lin, 2006: Reexamination of instrument change effects in the U.S. Historical Climatology Network. *Geophys. Res. Lett.*, **33**, ISI:000239786500008
- Jones, P.D., S. S.C.B. Raper, P.M. Kelly, and T.M.L. Wigley, R.S. Bradley and H.F. Diaz, 1986: Northern Hemisphere Surface Air Temperature Variations: 1851-1984. *J. Appl. Meteor.*, **25**, 161-179.
- Kalnay, E., M. Cai, H. Li, and J. Tobin, 2006: Estimation of the impact of land-surface forcings on temperature trends in eastern United States. *J. Geophys. Res.*, **111**, D06106, doi:10.1029/2005JD006555, 2006
- Karl, T.R., H.F. Diaz, and G. Kukla, 1988: Urbanization: its detection and effect in the United States climate record, *J. Climate*, **1**, 1099-1123.
- Karl, T.R., and C.N. Williams Jr., 1987: An approach to adjusting climatological time series for discontinuous inhomogeneities. *J. Climate Appl. Meteor.*, **26**, 1744-1763.
- Livezey, R.E., and W.Y. Chen, 1983: Statistical field significance and its determination by Monte Carlo techniques. *Mon. Wea. Rev.*, **111**, 46-59.
- Lund, R.B., and J. Reeves, 2002: Detection of undocumented changepoints—a revision of the two-phase regression model. *J. Climate*, **17**, 2547-2554.
- McCarthy, M.P., H. A. Titchner, P. W. Thorne, S. F. Tett, L. Haimberger, and D. E. Parker, 2007: Assessing Bias and Uncertainty in the HadAT Adjusted Radiosonde Climate Record . *J. Climate*, in press.
- Menne, M.J., and C.E. Duchon, 2001: A method for monthly detection of inhomogeneities and errors in daily maximum and minimum temperatures. *J. Atmos. Oceanic Tech.*, **18**, 1136-1149.
- Menne, M.J., and C.N. Williams, Jr., 2005: Detection of undocumented changepoints using multiple test statistics and composite reference series. *J. Climate*, **18**, 4271-4286.
- Parker, D.E., 2006: A demonstration that large-scale warming is not urban. *J. Climate*, **19**, 2882-2895.
- Peterson, T.C., and coauthors, 1998a: Homogeneity adjustments of in situ atmospheric climate data: A review. *Int. J. Climatol.*, **18**, 1493-1517.
- Peterson, T.C., T.R. Karl, P.F. Jamason, R. Knight, and D.R. Easterling, 1998b: First difference method: Maximizing station density for the calculation of the long-term global temperature change. *J. Geophys. Res.*, **74**, 22967-22974.

- and T.W. Owen, 2005: Urban heat island assessment: Metadata are important. *J. Climate*, **18**, 2637–2646.
- Pielke, R.A., C. Davey, J. Angel, O. Bliss, N. Doesken, M. Cai, S. Fall, D. Niyogi, K. Gallo, R. Hale, K. Hubbard, X. Lin, H. Li, J. Nielsen-Gammon, and S. Raman, 2007: Documentation of bias associated with surface temperature measurement sites for climate change assessment. *Bull. Amer. Meteor. Soc.*, **88**, 913-928.
- Reeves, J., J. Chen, X.L. Wang, R. Lund, and Q.Q. Lu, 2007: A review and comparison of changepoint detection techniques for climate data. *J. Appl. Meteor. Climatol.*, **46**, 900-914.
- Quayle, R.G., D.R. Easterling, T.R. Karl and P.Y. Hughes, 1991: Effects of recent thermometer changes in the cooperative station network. *Bull. Amer. Meteor. Soc.*, **72**, 1718-1723.
- Schwarz, G., 1978: Estimating the dimension of a model. *Ann. Stat.*, **6**, 461-464.
- Slonosky, V.C., P.D. Jones and T.D. Davies, 1999: Homogenization techniques for European monthly mean surface pressure series. *J. Climate*, **12**, 2658-2672.
- Thorne, P. W., D. E. Parker, J. R. Christy, and C. A. Mears, 2005a: Uncertainties in climate trends: Lessons from upper-air temperature records, *Bull. Am. Meteor. Soc.*, **86**, 1437-1442.
- Thorne, P. W., D. E. Parker, S. F. B. Tett, P. D. Jones, M. McCarthy, H. Coleman, and P. Brohan, 2005b: Revisiting radiosonde upper air temperatures from 1958 to 2002. *J. Geophys. Res.*, **110**, D18105, doi:10.1029/2004JD005753
- Vincent, L.A., 1998: A technique for the identification of inhomogeneities in Canadian temperature series. *J. Climate*, **11**, 1094-1104.
- Wang, X.L., 2003: Comments on “Detection of undocumented changepoints: A revision of the two-phase model”. *J. Climate*, **16**, 3383-3385.

## List of Tables

Table 1. Hierarchy of changepoint models for a temperature difference series  $\{D\}$ , where the subscript  $t$  refers to the time step of the series (e.g., one month),  $\mu$  refers to the mean,  $\beta$  refers to the trend, and  $\varepsilon_t$  represents a random error term.

Table 2. Changepoint detection and magnitude estimation skill for monthly temperature. RMSE of  $\Delta$  and  $\beta$  expressed in standardized units ( $\sigma$ ). The RMSE of  $\beta$  is calculated with respect to the true trend of zero.

Table 3. Number of changepoints imposed on each target and/or neighbor series for various case studies. The cases are comprised of 1000 simulations of six correlated series with  $n=100$  as described in Menne and Williams (2005).

Table 4. Skill scores from the pairwise homogenization algorithm for the case studies described in Table 3. The subscripts “*pw*” and “*ref*” refer to the pairwise and reference series approaches, respectively.

Table 5. Changepoint detection and model identification results (in percent) for 1,000 sets of 5 target–neighbor difference ( $\{D_t\}$ ) series ( $n=100$ ). Parameters were added as indicated to the target series and  $c=50$  for the target simulated under M3, M4, and M5. The neighbor series always followed M1 (constant mean with no breaks). *CRC* refers to the pairwise algorithm’s detection results for the target series. The percentage of  $\{D_t\}$  identified correctly are given in bold.

Table 6. Changepoint detection and magnitude estimation skill for monthly temperature series using a constant mean model (M3) for all step change adjustments regardless of the identified type.

## List of Figures

Figure 1. Mean monthly maximum temperature anomalies for Chula Vista (Target) and differences between monthly temperature anomalies at Chula Vista and nine neighboring series (T-N1 through T-N9).

Figure 2. Histogram of the most likely changepoint date identified by the SNHT for 10,000 series with  $n=100$  and a step change,  $\Delta$  at position 50. The magnitude of  $\Delta$  was varied systematically from 0.2 to 4.0.

Figure 3. As in Fig. 1, following adjustments by the pairwise algorithm.

Figure 4. “Annual” averages of simulated monthly series with a random number of changepoints imposed at random times and with random magnitudes. The true trend in all 10 correlated series is zero. Simulations are treated as beginning in January, 1901 and ending December, 2000.

Figure 5. As in Fig. 4, after homogenization by the pairwise algorithm.

Figure 6. Changepoint detection results for the monthly “steps only” simulations.

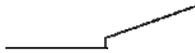
Figure 7. Changepoint detection results for the monthly steps and trends simulations.

Figure 8. Relationship between the hit rate (HR) and false alarm rate (FAR) for changepoints attributed to the target series as a function of the number of neighbors used to compute a composite reference series or in pairwise comparisons. Results are based on 1000 groups of series ( $n=100$ ) simulated under Case 6 (between 0 and 6 random changepoints added to the target and all neighbor series).

Figure 9. Mean monthly minimum temperature anomalies (in °C) for Cheesman, Colorado (Target) and differences between monthly temperature anomalies at Cheesman and 9 neighboring series (T-N1 to T-N9).

Figure 10. Differences between annual minimum temperatures at Cheezman, Colorado and 20 neighboring stations, and following adjustments for step changes using the most appropriate model determined by the pairwise algorithm (M3, M4 or M5) and using M3 only.

**Table 1. Hierarchy of changepoint models for a temperature difference series  $\{D\}$ , where the subscript  $t$  refers to the time step of the series (e.g., one month),  $\mu$  refers to the mean,  $\beta$  refers to the trend, and  $\varepsilon_t$  represents a random error term.**

<i>Model</i>	<i>Description</i>	<i>Schematic of Model</i>	<i>Number of parameters, <math>p</math>, required to fit model</i>
M1	$D_t = \mu + \varepsilon_t$		1
M2	$D_t = \mu + \beta t + \varepsilon_t$		2
M3	$D_t = \begin{cases} \mu_1 + \varepsilon_t, & t \leq c \\ \mu_2 + \varepsilon_t, & t > c \end{cases}$		3
M4	$D_t = \begin{cases} \mu_1 + \beta t + \varepsilon_t, & t \leq c \\ \mu_2 + \beta t + \varepsilon_t, & t > c \end{cases}$		4
M5	$D_t = \begin{cases} \mu_1 + \beta_1 t + \varepsilon_t, & t \leq c \\ \mu_2 + \beta_2 t + \varepsilon_t, & t > c \end{cases}$		5

**Table 2. Changepoint detection and magnitude estimation skill for monthly temperature. RMSE of  $\Delta$  and  $\beta$  expressed in standardized units ( $\sigma$ ). The RMSE of  $\beta$  is calculated with respect to the true trend of zero.**

<i>Case Study</i>	<i>Hit Rate (%)</i>	<i>FAR (%)</i>	<i>RMSE of <math>\Delta</math> (<math>\sigma</math>)</i>	<i>RMSE of <math>\beta</math> for unadjusted values (<math>\sigma</math>)</i>	<i>RMSE of <math>\beta</math> for adjusted values (<math>\sigma</math>)</i>
Monthly Data with Step Changes	67.11	6.77	0.284	2.455	0.401
Monthly Data with Step and Trend Changes	67.56	19.65	0.313	2.899	0.757

**Table 3. Number of changepoints imposed on each target and/or neighbor series for various case studies. The cases are comprised of 1000 simulations of six correlated series with  $n=100$  as described in Menne and Williams (2005).**

Scenario	Number of imposed changepoints	
	“Target” Series	Each “Neighbor” Series
Case 1 (null case)	0	0
Case 2	2	0
Case 3	0	2
Case 4	2	2
Case 5 (null case with missing values)	0	0
Case 6	0 to 6*	0 to 6*
Case 7	6**	0

\* the number of change points in each series is normally distributed about an average of 3.

\*\* changepoint position and magnitude are fixed as in Caussinus and Mestre (2004): +2.0 at  $c = 20$ , +2.0 at  $c = 40$ , -2.0 at  $c = 50$ , -2.0 at  $c = 70$ , +2.0 at  $c = 75$ , and +2.0 at  $c = 85$ .

**Table 4. Skill scores from the pairwise homogenization algorithm for the case studies described in Table 3. The subscripts “*pw*” and “*ref*” refer to the pairwise and reference series approaches, respectively.**

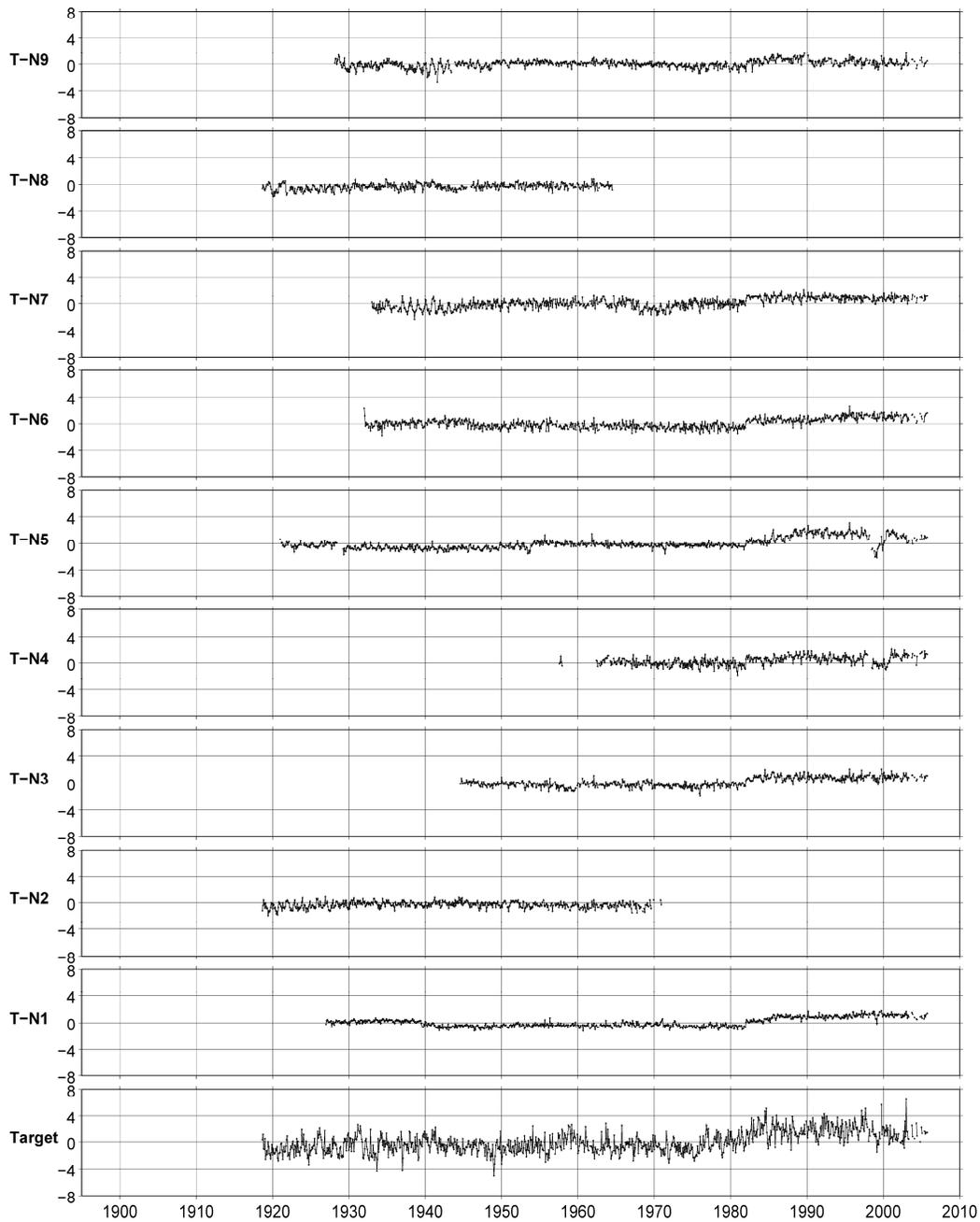
<b><i>CASE STUDY (and scenario description)</i></b>	<b><i>CRC<sub>pw</sub> (%)</i></b>	<b><i>FAR<sub>pw</sub> (%)</i></b>	<b><i>CRC<sub>ref</sub> (%)</i></b>	<b><i>FAR<sub>ref</sub> (%)</i></b>
CASE 1 (homogeneous target and neighbor series)	99.5	100.0	88.8	100.0
CASE 2 (two random changepoints in target; homogeneous neighbor series)	44.0	5.6	55.4	21.0
CASE 3 (homogeneous target series; two random changepoints in each neighbor series)	95.2	100.0	0.0	100.0
CASE 4 (two random change points in all series)	37.3	8.5	50.3	46.0
CASE 5 (homogeneous target and neighbor series with missing values)	100.0	<i>Undefined (zero false alarms)</i>	87.2	100.0
CASE 6 (up to six changepoints in all series)	31.6	7.0	45.4	41.0
CASE 7 (six changepoints in target [ $\Delta=2\sigma$ ]; homogeneous neighbors)	84.6	1.1	70.4	6.0

**Table 5. Changepoint detection and model identification results (in percent) for 1,000 sets of 5 target–neighbor difference ( $\{D_t\}$ ) series ( $n=100$ ). Parameters were added as indicated to the target series and  $c=50$  for the target simulated under M3, M4, and M5. The neighbor series always followed M1 (constant mean with no breaks). CRC refers to the pairwise algorithm’s detection results for the target series. The percentage of  $\{D_t\}$  identified correctly are given in bold.**

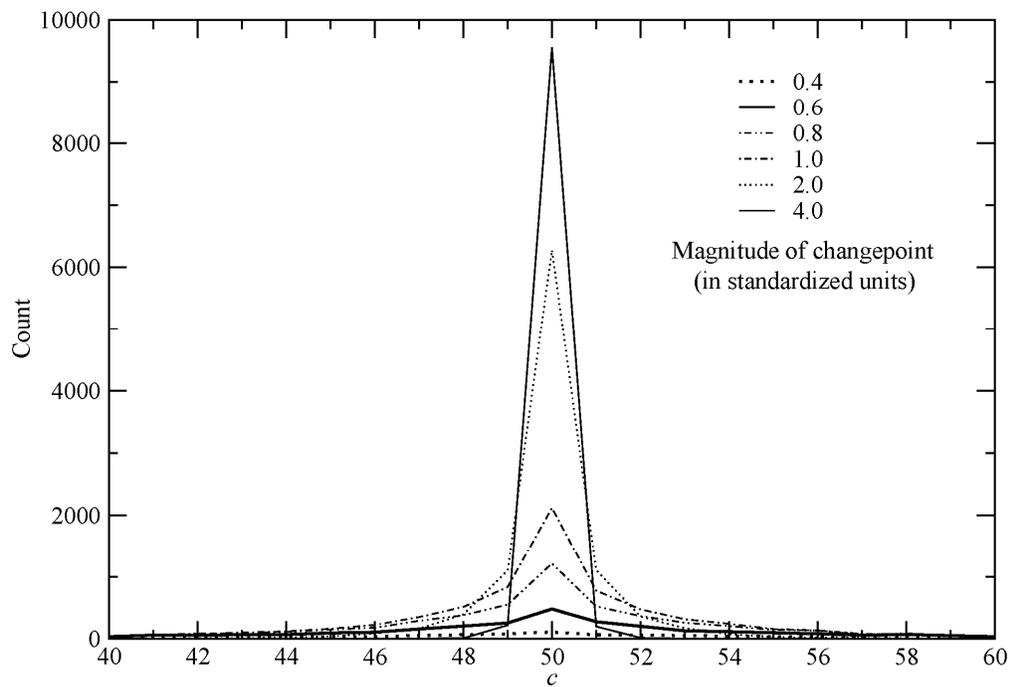
<i>Target Series Follows M2</i>								
	$\beta$		<i>M1</i>	<i>M2</i>	<i>M3</i>	<i>M4</i>	<i>M5</i>	<i>CRC</i>
	0.005		51.25	<b>44.36</b>	3.86	0.24	0.30	95.70
	0.010		2.45	<b>88.23</b>	7.30	0.72	1.31	91.00
	0.020		0.55	<b>85.75</b>	3.48	4.56	5.67	86.90
<i>Target Series Follows M3</i>								
$\Delta$			<i>M1</i>	<i>M2</i>	<i>M3</i>	<i>M4</i>	<i>M5</i>	<i>CRC</i>
0.5			11.71	45.16	<b>39.03</b>	1.66	2.44	30.30
1.0			0.06	4.83	<b>88.59</b>	1.82	4.70	90.90
2.0			0.00	0.10	<b>93.21</b>	1.85	4.85	99.90
<i>Target Series Follows M4</i>								
$\Delta$	$\beta$		<i>M1</i>	<i>M2</i>	<i>M3</i>	<i>M4</i>	<i>M5</i>	<i>CRC</i>
1.0	0.005		0.04	6.56	80.08	<b>5.24</b>	8.08	89.30
1.0	0.010		0.06	7.32	51.07	<b>24.56</b>	16.99	87.90
1.0	0.020		0.05	7.46	19.75	<b>52.87</b>	19.87	85.50
<i>Target Series Follows M5</i>								
$\Delta$	$\beta_1$	$\beta_2$	<i>M1</i>	<i>M2</i>	<i>M3</i>	<i>M4</i>	<i>M5</i>	<i>CRC</i>
1.0	0.010	0.015	0.11	8.16	34.84	33.17	<b>23.72</b>	87.11
1.0	0.010	0.020	0.07	8.45	25.17	33.51	<b>32.79</b>	86.20
1.0	0.010	0.030	0.07	7.47	19.34	23.22	<b>49.91</b>	85.70

**Table 6. Changepoint detection and magnitude estimation skill for monthly temperature series using a constant mean model (M3) for all step change adjustments regardless of the identified type.**

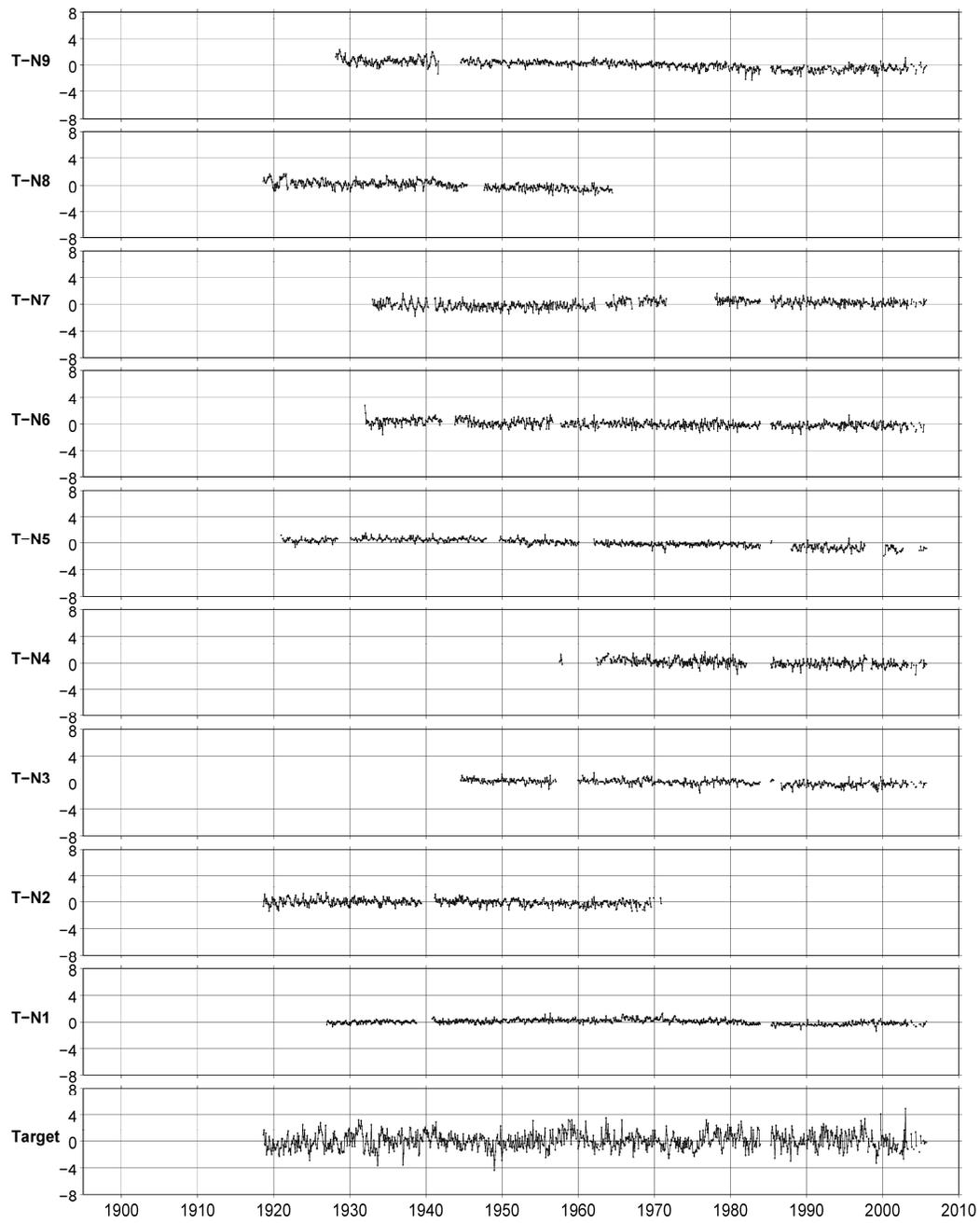
<i>Case Study</i>	<i>Hit Rate (%)</i>	<i>FAR (%)</i>	<i>RMSE of <math>\Delta</math> (<math>\sigma</math>)</i>	<i>RMSE of <math>\beta</math> for unadjusted values (<math>\sigma</math>)</i>	<i>RMSE of <math>\beta</math> for adjusted values (<math>\sigma</math>)</i>
Monthly Data with Step Changes	67.22	6.77	0.291	2.455	0.401
Monthly Data with Step and Trend Changes	67.58	20.14	0.349	2.899	0.488



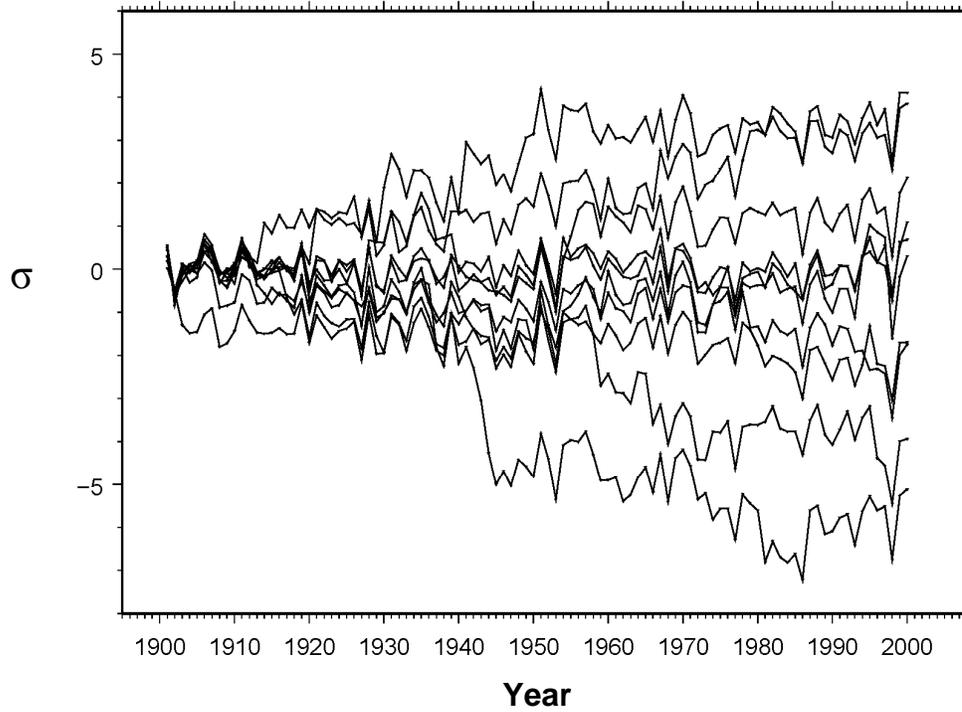
**Figure 1. Mean monthly maximum temperature anomalies for Chula Vista (Target) and differences between monthly temperature anomalies at Chula Vista and nine neighboring series (T-N1 through T-N9).**



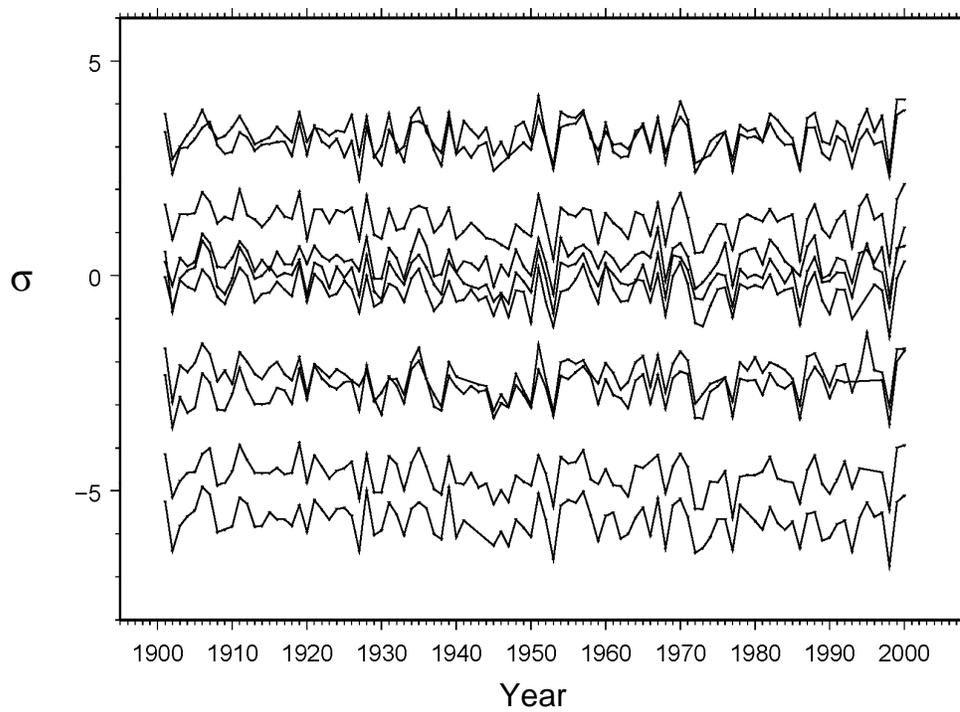
**Figure 2. Histogram of the most likely changepoint date identified by the SNHT for 10,000 series with  $n=100$  and a step change,  $\Delta$  at position 50. The magnitude of  $\Delta$  was varied systematically from 0.2 to 4.0.**



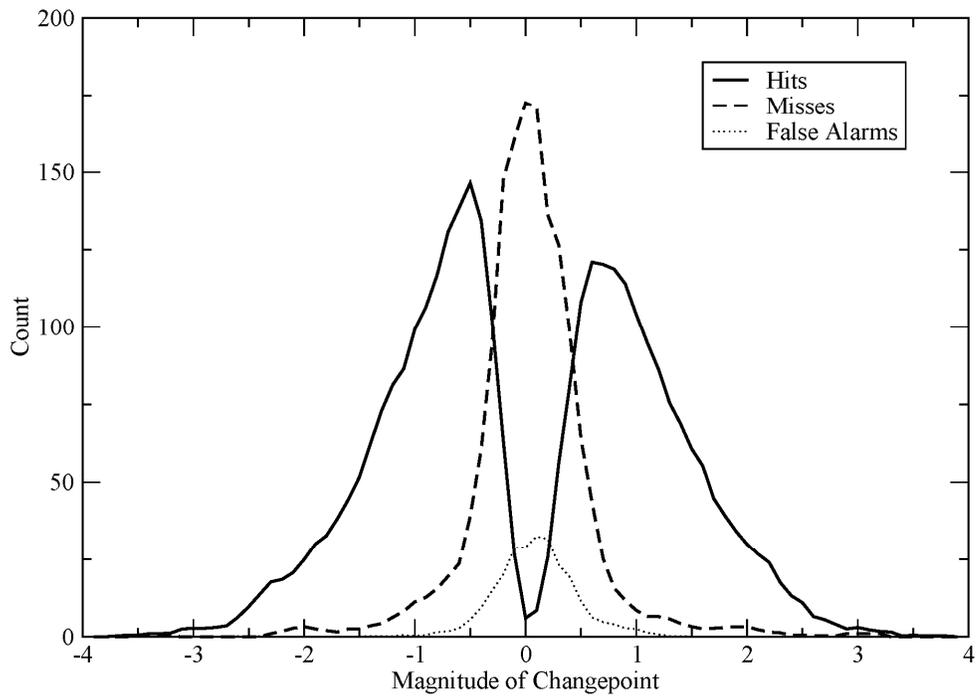
**Figure 3.** As in Fig. 1, following adjustments by the pairwise algorithm.



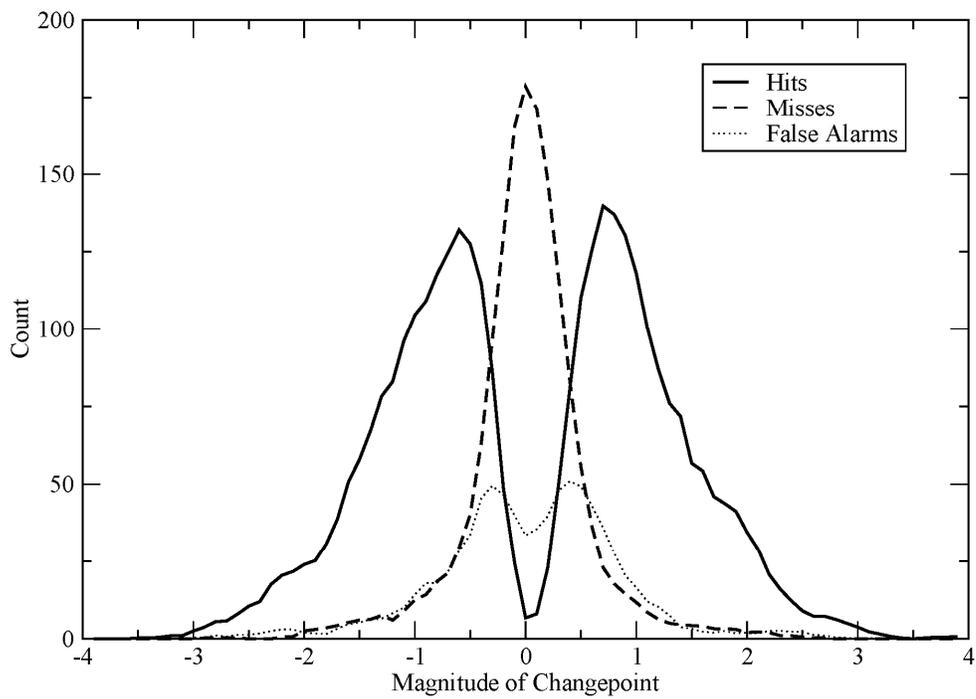
**Figure 4. “Annual” averages of simulated monthly series with a random number of changepoints imposed at random times and with random magnitudes. The true trend in all 10 correlated series is zero. Simulations are treated as beginning in January, 1901 and ending December, 2000.**



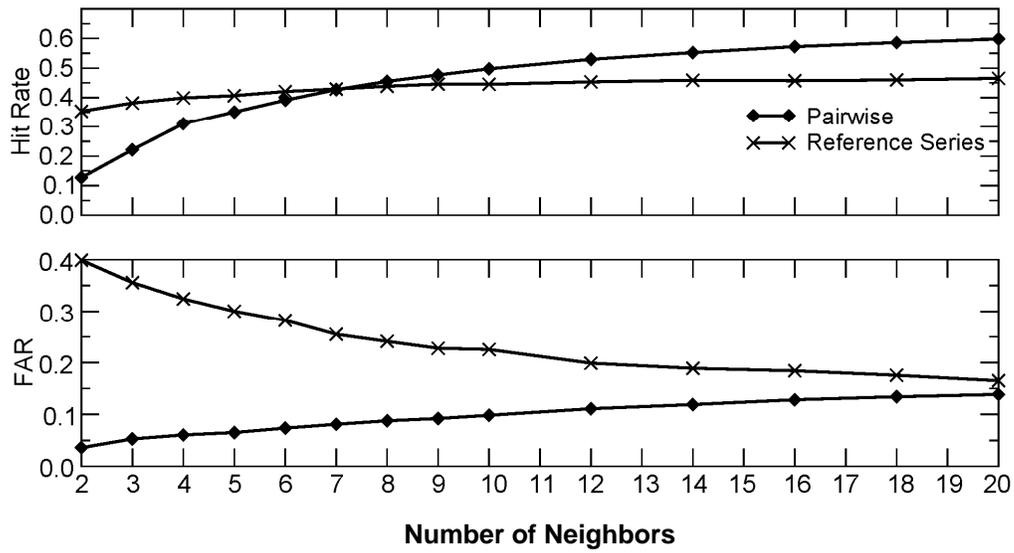
**Figure 5. As in Fig. 4, after homogenization by the pairwise algorithm.**



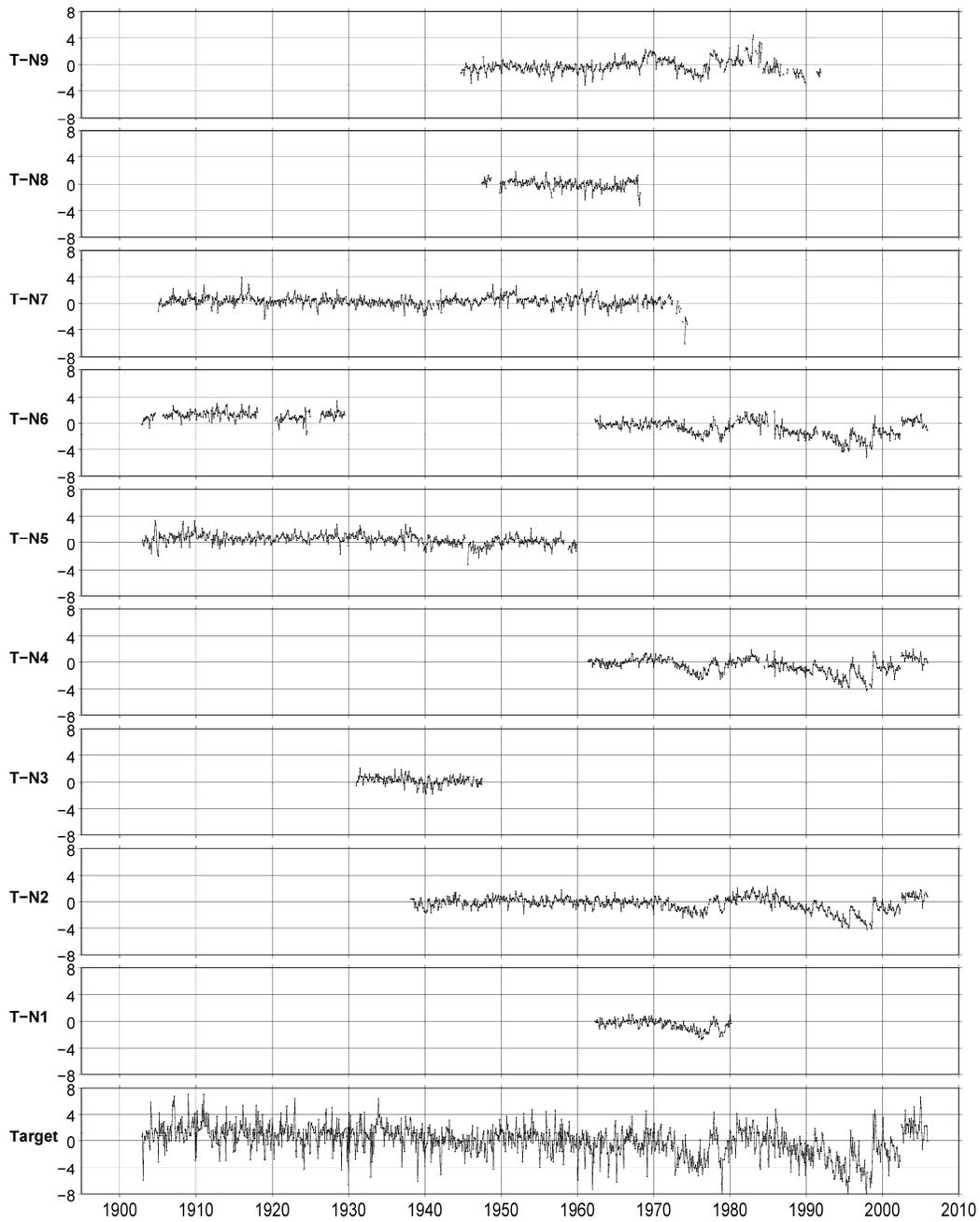
**Figure 6. Changepoint detection results for the monthly “steps only” simulations.**



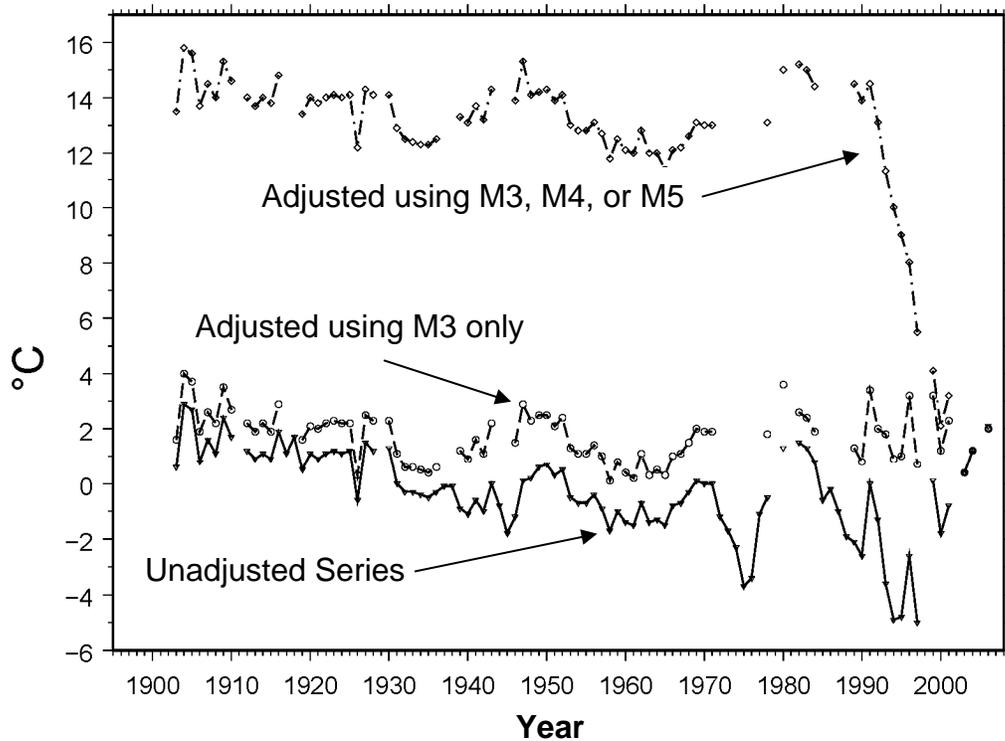
**Figure 7. Changepoint detection results for the monthly steps and trends simulations.**



**Figure 8. Relationship between the hit rate (HR) and false alarm rate (FAR) for changepoints attributed to the target series as a function of the number of neighbors used to compute a composite reference series or in pairwise comparisons. Results are based on 1000 groups of series ( $n=100$ ) simulated under Case 6 (between 0 and 6 random changepoints added to the target and all neighbor series).**



**Figure 9. Mean monthly minimum temperature anomalies (in °C) for Cheesman, Colorado (Target) and differences between monthly temperature anomalies at Cheesman and 9 neighboring series (T-N1 to T-N9).**



**Figure 10. Differences between annual minimum temperatures at Cheezman, Colorado and 20 neighboring stations (solid line) and following adjustments for step changes using the most appropriate model determined by the pairwise algorithm (M3, M4, or M5) and using only M3.**