

Article

## Reliability of Inference of Directed Climate Networks Using Conditional Mutual Information

Jaroslav Hlinka <sup>1,\*</sup>, David Hartman <sup>1</sup>, Martin Vejmelka <sup>1</sup>, Jakob Runge <sup>2,3</sup>, Norbert Marwan <sup>2</sup>, Jürgen Kurths <sup>2,3,4</sup> and Milan Paluš <sup>1</sup>

<sup>1</sup> Institute of Computer Science, Academy of Sciences of the Czech Republic, Pod vodarenskou vezi 2, 182 07, Prague 8, Czech Republic; E-Mails: hartman@cs.cas.cz (D.H.); vejmelka@cs.cas.cz (M.V.); mp@cs.cas.cz (M.P.)

<sup>2</sup> Potsdam Institute for Climate Impact Research (PIK), 14473 Potsdam, Germany; E-Mails: jakobrunge@gmail.com (J.R.); marwan@pik-potsdam.de (N.M.); kurths@pik-potsdam.de (J.K.)

<sup>3</sup> Department of Physics, Humboldt University, 12489 Berlin, Germany

<sup>4</sup> Institute for Complex Systems and Mathematical Biology, University of Aberdeen, Aberdeen AB24 3UE, United Kingdom

\* Author to whom correspondence should be addressed; E-Mail: hlinka@cs.cas.cz; Tel.: +420-266-053-808; Fax.: +420-286-585-789.

Received: 30 January 2013; in revised form: 11 May 2013 / Accepted: 14 May 2013 /

Published: 24 May 2013

---

**Abstract:** Across geosciences, many investigated phenomena relate to specific complex systems consisting of intricately intertwined interacting subsystems. Such dynamical complex systems can be represented by a directed graph, where each link denotes an existence of a causal relation, or information exchange between the nodes. For geophysical systems such as global climate, these relations are commonly not theoretically known but estimated from recorded data using causality analysis methods. These include bivariate nonlinear methods based on information theory and their linear counterpart. The trade-off between the valuable sensitivity of nonlinear methods to more general interactions and the potentially higher numerical reliability of linear methods may affect inference regarding structure and variability of climate networks. We investigate the reliability of directed climate networks detected by selected methods and parameter settings, using a stationarized model of dimensionality-reduced surface air temperature data from reanalysis of 60-year global climate records. Overall, all studied bivariate causality methods provided reproducible estimates of climate causality networks, with the linear approximation showing

higher reliability than the investigated nonlinear methods. On the example dataset, optimizing the investigated nonlinear methods with respect to reliability increased the similarity of the detected networks to their linear counterparts, supporting the particular hypothesis of the near-linearity of the surface air temperature reanalysis data.

**Keywords:** causality; climate; nonlinearity; transfer entropy; network; stability

---

## 1. Introduction

Across geosciences, many investigated phenomena relate to specific complex systems consisting of intricately intertwined interacting subsystems. These can be suitably represented as networks, an approach that is gaining increasing attention in the complex systems community [1,2]. The meaning of the existence of a link between nodes of a network depends on the area of application, but in many cases it is related to some form of information exchange between the nodes.

This approach has already been adopted for the analysis of various phenomena in the global climate system [3–7]. For a recent review discussing the progress, opportunities and challenges, see [8]. Typically, a graph is constructed by considering two locations linked by a connection, if there is an instantaneous dependence between the localized values of a variable of interest.

This dependence can be conveniently quantified by mutual information, an entropy-based general measure of statistical dependence that takes into account nonlinear contributions to the coupling. In practice, for reasons of theoretical and numerical simplicity, linear Pearson's correlation coefficient might be sufficient, although potentially neglecting the nonlinear contributions to interactions. In particular, while initial works by Donges *et al.* stressed the role of mutual information in detecting important features of global climate networks [9,10], more detailed recent work has shown that the differences between correlation and mutual information graphs are mostly (but not necessarily completely) spurious, such as due to natural and instrumental (related to data collection) nonstationarities of the data [11].

However, these methods do not allow to assess the directionality of the links and of the underlying information flow. This motivates the use of more sophisticated measures, known also as causality analysis methods. Indeed, it has been recently suggested that data-driven detection of climate causality networks could be used for deriving hypotheses about causal relationships between prominent modes of atmospheric variability [12,13].

The family of causality methods include linear approaches such as Granger causality analysis [14], as well as more general nonlinear methods. A prominent representative of nonlinear causality assessment is the conditional mutual information [15], especially in the form of transfer entropy [16].

Arguably, the nonlinear methods, due to their model-free nature, have the theoretical advantage of being sensitive to forms of interactions that linear methods may detect only partially or not at all. On the other hand, this advantage may be more than outweighed by a potentially lower precision. Depending on specific circumstances, this may adversely affect the reliability of the detection of network patterns.

Apart from uncertainty about the general network pattern, reliability is important when the focus lies on detecting changes in time, with the need to distinguish them from random variability of the estimates among different sections of time series under investigation—a task that is relevant in many areas of geoscience including climate research. In other words, before analyzing a complex dynamical system using network theory, a key initial question is that of the reliability of the network construction, and how it depends on the choice of a causality method and its parameters.

We study this question for a selection of standard causality methods, using a timely application in the study of a climate network and its variability. In particular, surface air temperature data from the NCEP/NCAR reanalysis dataset [17,18] is used. The original data contains more than 10,000 time series—a relatively dense grid covering the whole globe. For efficient computation and visualization of the results, it is convenient to reduce the dimensionality of the data. We use principal component analysis and select only components that have significantly high explained variance compared with the corresponding spatially independent but temporally dependent (*i.e.*, “colored”) random noise.

As the causality network construction reliability may crucially depend on the specific choice of causality estimator, we test different causality measures and their parametrization by quantifying the similarity of causality matrices reconstructed from independent realizations of a stationary model of data. These realizations are either independently generated, or they represent individual non-overlapping temporal windows of a single stationary realization. Optimal parameter choices of the applied nonlinear methods are detected, and the reliability of networks constructed using linear and nonlinear methods is compared. The latter method, *i.e.*, comparing networks reconstructed from temporal windows, allows also to assess the network variability on real data and to compare it with the variability observed in the analysis of stationary model time series.

## 2. Data and Methods

### 2.1. Causality Assessment Methods

#### 2.1.1. Granger Causality Analysis

A prominent method for assessing causality is Granger causality analysis, named after Sir Clive Granger, who proposed this approach to time series analysis in a classical paper [14]. However, the basic idea can be traced back to Wiener [19], who proposed that if the prediction of one time series can be improved by incorporating the knowledge of a second time series, then the latter can be said to have a causal influence on the former. This idea was formalized by Granger in the context of linear regression models. In the following, we outline the methods of assessment of Granger causality, following the description given in [20–22].

Consider two stochastic processes  $X_t$  and  $Y_t$  and assume they are jointly stationary. Let further the autoregressive representations of each process be:

$$X_t = \sum_{j=1}^{\infty} a_{1j} X_{t-j} + \epsilon_{1t}, \quad \text{var}(\epsilon_{1t}) = \Sigma_1 \tag{1}$$

$$Y_t = \sum_{j=1}^{\infty} d_{1j} Y_{t-j} + \eta_{1t}, \quad \text{var}(\eta_{1t}) = \Gamma_1 \tag{2}$$

and joint autoregressive representation

$$X_t = \sum_{j=1}^{\infty} a_{2j} X_{t-j} + \sum_{j=1}^{\infty} b_{2j} Y_{t-j} + \epsilon_{2t} \tag{3}$$

$$Y_t = \sum_{j=1}^{\infty} c_{2j} X_{t-j} + \sum_{j=1}^{\infty} d_{2j} Y_{t-j} + \eta_{2t} \tag{4}$$

where the covariance matrix of the noise terms is:

$$\Sigma = \text{Cov} \begin{pmatrix} \epsilon_{2t} \\ \eta_{2t} \end{pmatrix} = \begin{pmatrix} \Sigma_2 & \Lambda_2 \\ \Lambda_2 & \Gamma_2 \end{pmatrix} \tag{5}$$

The causal influence from  $Y$  to  $X$  is then quantified based on the decrease in the residual model variance when we include the past of  $Y$  in the model of  $X$ , *i.e.*, when we move from the independent model given by Equation (1) to the joint model given by Equation (3):

$$F_{Y \rightarrow X} = \ln \frac{\Sigma_1}{\Sigma_2} \tag{6}$$

Similarly, the causal influence from  $X$  to  $Y$  is defined as:

$$F_{X \rightarrow Y} = \ln \frac{\Gamma_1}{\Gamma_2} \tag{7}$$

Clearly, the causal influence defined in this way is always nonnegative.

The original introduction of the concept of statistical inference of causality [14] includes a third (potentially highly multivariate) process  $Z$ , representing all other intervening processes that should be controlled for in assessing the causality between  $X$  and  $Y$ . The bivariate (or “pairwise”) implementation of the estimator thus constitutes a computational simplification of the original process, for the sake of numerical stability as well as comparability with the bivariate transfer entropy (conditional mutual information) approach introduced later. See Section 4 for further discussion of related issues.

### 2.2. Estimation of GC

Practical estimation of Granger causality involves fitting the joint and univariate models described above to experimental data. While the theoretical framework is formulated in terms of infinite sums, the fitting procedure requires selection of the model order  $p$  for the models. For our report, we have selected  $p = 1$ , corresponding to links of lag 1, since this lag was also chosen in the nonlinear methods considered later. This is a common choice for Granger causality in literature and amounts to looking for links with lag 1 time unit.

### 2.3. Transfer Entropy

The information-theoretic analog to Granger causality is the concept of *transfer entropy* (TE [16]). TE can be defined in terms of *conditional mutual information* as shown below, closely following [15]. In particular, we can define that  $X$  causes  $Y$  if the knowledge of the past of  $X$  decreases the uncertainty about  $Y$  (above what the knowledge of past of  $Y$  and potentially all other relevant confounding variables already informs).

For two discrete random variables  $X, Y$  with sets of values  $\Xi$  and  $\Upsilon$  and probability distribution functions (PDFs)  $p(x), p(y)$  and joint PDF  $p(x, y)$ , the Shannon entropy  $H(X)$  is defined as

$$H(X) = - \sum_{x \in \Xi} p(x) \log p(x) \quad (8)$$

and the joint entropy  $H(X, Y)$  of  $X$  and  $Y$  as

$$H(X, Y) = - \sum_{x \in \Xi} \sum_{y \in \Upsilon} p(x, y) \log p(x, y) \quad (9)$$

The conditional entropy  $H(X|Y)$  of  $X$  given  $Y$  is

$$H(X|Y) = - \sum_{x \in \Xi} \sum_{y \in \Upsilon} p(x, y) \log p(x|y) \quad (10)$$

The amount of shared information contained in the variables  $X$  and  $Y$  is quantified by the mutual information  $I(X; Y)$  defined as

$$I(X; Y) = H(X) + H(Y) - H(X, Y) \quad (11)$$

The conditional mutual information  $I(X; Y|Z)$  of the variables  $X, Y$  given the variable  $Z$  is given as

$$I(X; Y|Z) = H(X|Z) + H(Y|Z) - H(X, Y|Z) \quad (12)$$

Entropy and mutual information are measured in bits if the base of the logarithms in their definitions is 2. It is straightforward to extend these definitions to more variables, and to continuous rather than discrete variables.

The transfer entropy from  $X$  to  $Y$  then corresponds to the conditional mutual information between  $X_t$  and  $Y_{t+1}$  conditional on  $Y_t$ :

$$T_{X \rightarrow Y} = I(X_t, Y_{t+1} | Y_t) \quad (13)$$

While the definition of these information-theoretic functionals is very general and elegant, the practical estimation faces challenges related to the problem of efficient estimation of the PDF from finite size samples. It is important to bear in mind the distinction between the quantities of the underlying stochastic process and their finite-sample estimators.

#### 2.4. Potential Causes of Observed Difference

Interestingly, it can be shown that for linear Gaussian processes, transfer entropy is equivalent to linear Granger causality, up to a multiplicative factor [23]:

$$\mathcal{T}_{X \rightarrow Y} = \frac{1}{2} \mathcal{F}_{X \rightarrow Y} \quad (14)$$

However, in practice, the estimates of transfer entropy and linear Granger causality may differ. There are principally two main reasons for this divergence between the results. Firstly, when the underlying process is not linear Gaussian, the true transfer entropy may differ from the true linear Granger causality corresponding to the linear approximation of the process. A second reason for divergence between sample estimates of transfer entropy and linear Granger causality, valid even for linear Gaussian processes, is the difference in the properties of the estimators of these two quantities, in particular bias and variance of the estimates.

#### 2.5. TE Estimation

There are many algorithms for the estimation of information-theoretical functionals that can be adapted to compute transfer entropy estimates. Two basic classes of nonparametric methods for the estimation of conditional mutual information are the *binning methods* and the *metric methods*. The former discretize the space into regions usually called bins or boxes—a robust example is the equiquantal method based on discretization of studied variables into  $Q$  equiquantal bins (EQQ [24]). In the latter methods, the probability distribution function estimation depends on the distances between the samples computed using some metric. An example of a metric method is the  $k$ -nearest neighbor (kNN [15,25]) algorithm. For more detail on methods of estimation of conditional mutual information and their comparison, see [15].

Note that both types of algorithms require setting an additional parameter. While some heuristic suggestions have been published in the literature, suitable values of these parameters may depend on specific aspects of the application including the character of the time series. For the purpose of this study, we use a range of parameter values and subsequently select the parameter values providing the most stable results for further comparison with linear methods, see below.

#### 2.6. Data

##### 2.6.1. Dataset

Data from the NCEP/NCAR reanalysis dataset [17] has been used. In particular, we utilize the time series  $x_i(t)$  of daily and monthly mean surface air temperature from January 1948 to December 2007 ( $T_d = 21,900$  and  $T_m = 720$  time points, respectively), sampled at latitudes  $\lambda_i$  and longitude  $\phi_i$  forming a regular grid with a step of  $\Delta\lambda = \Delta\phi = 2.5^\circ$ . The points located at the poles have been removed, giving a total of  $N = 10,224$  spatial sampling points.

### 2.6.2. Preprocessing

To minimize the bias introduced by periodic changes due to solar irradiation, the mean annual cycle has been removed to produce anomaly time series. The data were further standardized such that the time series at each grid point has unit variance. The time series are then scaled by the cosine of the latitude to account for grid points closer to the poles representing smaller areas and being closer together (thus biasing the correlation with respect to grid points farther apart). The poles are thus omitted entirely by effectively removing data for latitude  $\pm 90$ .

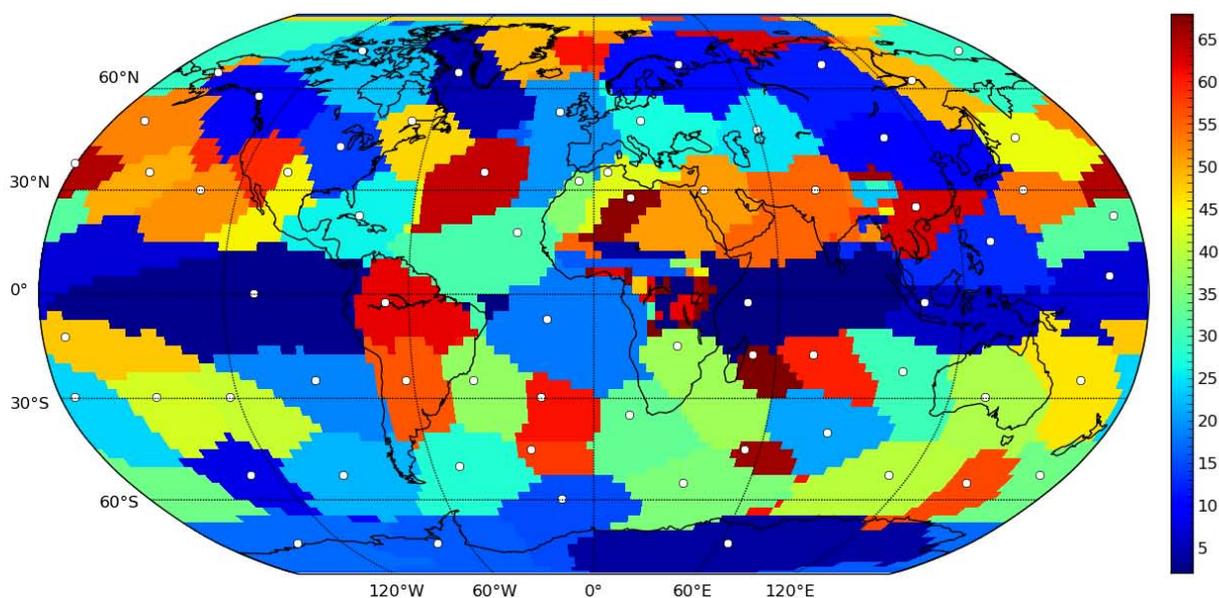
### 2.6.3. Computing the Components

First, the covariance matrix of the scaled time series obtained by preprocessing is computed. Note that this covariance matrix is equal to the correlation matrix, where each correlation is scaled by the inverse of the product of the cosines of the latitudes of the time series entering the correlation.

Next, the eigendecomposition of the covariance matrix is computed. The eigenvectors corresponding to genuine components are extracted (the estimation of the number of components is explained in the next paragraph). The eigenvectors are then rotated using the VARIMAX method [26].

The rotated eigenvectors are the resulting components. Each component is represented by a scalar field of intensities over the globe, and by a corresponding time series. See Figure 1 for a parcellation of the globe by the components. For each location, the color corresponding to the component with maximal intensity is used—due to good spatial localization and smoothness of the components, this leads to the parcellation of the globe into generally contiguous regions.

**Figure 1.** Location of areas dominated by specific components of the surface air temperature data using VARIMAX-rotated PCA decomposition. For each location, the color corresponding to the component with maximal intensity was used. White dots represent approximate centers of mass of the components, used in subsequent figures for visualization of the nodes of the networks.



#### 2.6.4. Estimating the Dimensionality of the Data

To reduce the dimensionality, only a subset of the components is selected for further analysis. The main idea rests in determining significant components by comparing the eigenvalues computed from the original data to eigenvalues computed from a control dataset corresponding to the null hypothesis of uncoupled time series with the same temporal structure as the original data. To accomplish this, the time series in the control datasets are generated as realizations of autoregressive (AR) models that are fitted to each time series independently. The dimension of the AR process is estimated for each time series separately using the Bayesian Information Criterion [27].

This model is used to generate 10,000 realizations in the control dataset. The eigendecomposition of each realization is computed and aggregated, providing a distribution for each eigenvalue (1st, 2nd, ...) under the above null hypothesis. Finding the significant eigenvalues then reduces to a multiple comparison problem, which we resolved using the False Discovery Rate (FDR) technique [28], leading to the identification of 67 components.

For computational reasons, the decomposition was carried out on the monthly data and the resulting component weights were used to extract daily time series from the corresponding preprocessed daily data (anomalization, standardization, cosine transform). The method thus provides full-resolution component localization on the 10224-point grid while also yielding a high-resolution time series associated with each component. Indeed, carrying out the decomposition directly on the daily data might have provided a slightly different set of components, as the decomposition would also take into account high-frequency variability. The current decomposition has both conceptual and practical motivation. Conceptually, it provides information about the high-frequency behavior of the variability modes detected in the monthly data (a timescale more commonly used for decomposition in the climate community). The practical reason is that the decomposition of the full  $10224 \times 21900$  matrix of daily data is computationally demanding, particularly in combination with the bootstrap testing procedure described above.

#### 2.7. Network Construction

Formally, in the graph-theoretical approach, a network is represented by a graph  $G = (V, E)$ , where  $V$  is the set of nodes of  $G$ ,  $n = \#V$  is the number of nodes and  $E \subset V^2$  is the set of the edges (or links) of  $G$ . In weighted graphs, each edge connecting nodes  $i$  and  $j$  can be assigned a weight  $a_{i,j}$  representing the strength of the link. Thus, the pairwise causality matrix  $\mathcal{T}$  with entries  $\mathcal{T}_{i,j} = \mathcal{T}_{X_i \rightarrow X_j}$  can be understood as a weighted graph. Commonly, the graph is transformed into an unweighted matrix by suitable thresholding, keeping only links with weights higher than some threshold (and setting their weights to 1) while removing the weaker links (setting their weights to zero).

There are three principal strategies to choose the threshold. One can choose a fixed value based on expert judgment of what constitutes a strong link, or adaptively to enforce a required density of the graph (relative number of links with respect to the maximum number possible, *i.e.*, in a full graph of given size). The third option is to use statistical testing to detect statistically significant links. In the current paper, we start with the original unthresholded graphs, but provide also example results for thresholded graphs using the above approaches.

## 2.8. Reliability Assessment

In line with the terminology of psychometrics or classical test theory, by reliability we mean the overall consistency of a measure (consistency here does not mean the statistical sense of asymptotic behavior). In the context of network construction, we considered a method reliable if the networks constructed from different samples of the same dynamical process are similar to each other. Note that this does not necessarily imply validity or accuracy of the method—under some circumstances, a method could consistently arrive at wrong results. In a way, reliability/consistency can be considered a first step to validity. In practical terms, even if the validity was undoubted, reliability can give the researcher an estimate on the confidence he/she can have in the reproducibility of the results.

To assess the similarity of two matrices, many methods are available, including the (entry-wise) Pearson's linear correlation coefficient. Inspection of the causality matrices suggests a heavily non-normal distribution of the values with many outliers. Therefore, the correlation of ranks, using Spearman's correlation coefficient, may be more suitable.

Apart from reliability of the full weighted causality graphs, we also study the unweighted graphs derived by thresholding. Based on inspection of the causality matrices, a density of 0.01 (keeping 1 percent of strongest links) was chosen for the analysis. To assess the similarity of two binary matrices, we use the Jaccard similarity coefficient. This is the relative number of links that are shared by the matrices with respect to the total number of links that appear at least in one of the matrices. Such a ratio is a natural measure of matrix overlap, ranging from 0 for matrices with no common links to 1 for identical matrices.

### 2.8.1. Model

A convenient method for assessing the reliability of a method on time series is to compare the results obtained from different temporal windows. However, dissimilarity among the results can be theoretically attributed to both lack of reliability of the method and hypothetical true changes in the underlying system over time (nonstationarity).

Therefore, to isolate the effect of method properties, we test the methods on a realistic but stationary model of the data. To provide such a stationary model of the potentially non-stationary data, surrogate time series were constructed.

Technically, surrogate time series are conveniently constructed as multivariate Fourier transform (FT) surrogates [29,30]. They are computed by first performing a Fourier transform of the series and keeping unchanged the magnitudes of the Fourier coefficients (the amplitude spectrum). Then the same random number is added to the phases of coefficients of the same frequency bin. The surrogate is then the inverse FT into the time domain.

The surrogate data represent a realization of a linear stationary process conserving the linear structure (covariance and autocovariance) of the original data, and hence also the linear component of causality. Note that any nonlinear component of causality should be removed, and the nonlinear methods should therefore converge to the linear methods (as discussed in Section 2.1).

After testing the reliability on the stationary linear model, we assess the stability of the methods also on real data. The variability here should reflect a mixture of method non-reliability and true climate

changes. Note that also the nonlinear methods may potentially diverge from the linear ones due to nonlinear causalities in the data.

### 2.8.2. Implementation Details

For estimation, both the stationary model and the real data time series were split into 6 windows (one for each decade, *i.e.*, with approximately 3650 time points). For each of the windows, the causality matrix has been estimated with several causality methods.

In particular, we have used pairwise Granger causality as a representative linear method and computed transfer entropy by two standard algorithms using a range of critical parameter values. The first is an algorithm based on the discretization of studied variables into  $Q$  equiquantal bins (EQQ [24],  $Q \in \{2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14\}$ ) and the second is a  $k$ -nearest neighbor algorithm (kNN[15,25],  $k \in \{2, 4, 8, 16, 32, 64, 128, 256, 512\}$ ).

Each of these algorithms provides a matrix of causality estimates among the 67 climate components within the respective decade. We further assess the similarity of these matrices across both time and methods, first in stationary data (where temporal variability is attributable to method instability only) and then in real data. Apart from direct visualization, the similarity of the constructed causality matrices is quantified by the Spearman's rank correlation coefficient of off-diagonal entries. The reliability is then estimated as the average Spearman's rank correlation coefficient across all  $(6 \times 5)/2 = 15$  pairs of temporal windows.

To inspect the robustness of the results, the analysis was repeated with several possible alterations to the approach. Firstly, the similarity among the thresholded rather than unweighted graphs was assessed by means of the Jaccard similarity coefficient instead of Spearman's rank correlation coefficient. Secondly, we repeated the analysis using linear multivariate AR(1) process for the generation of the stationary model instead of Fourier surrogates. Thirdly, the analysis was repeated on sub-sampled data (by averaging each 6 days to give one data point). This way, the same methods should provide causality on a longer time scale. To keep the same (and sufficiently high) number of time points, the sub-sampled data were not split into windows, but 6 realizations were generated from a fitted multivariate AR(1) process.

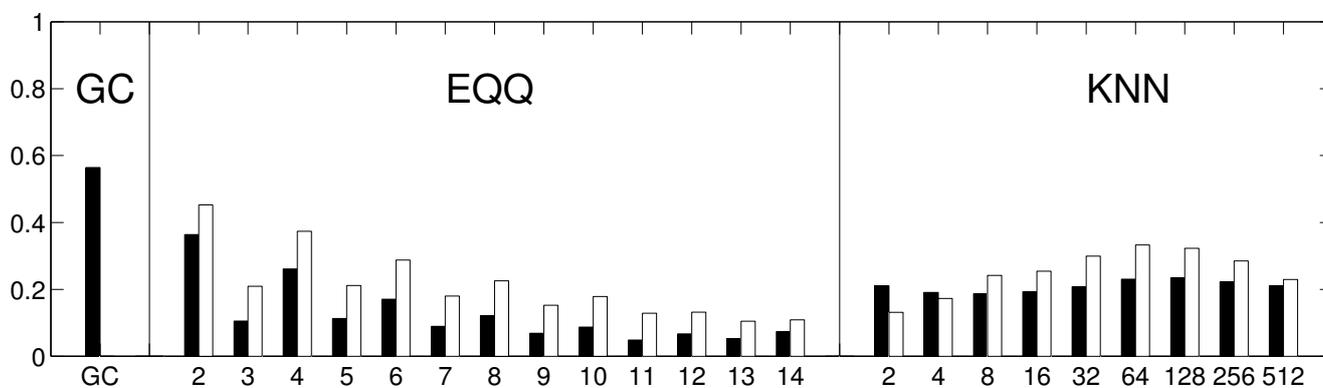
Finally, to assess the role that the reliability of different methods may have for statistical inference, we tested the number of links that the method was able to statistically distinguish compared with an empirical distribution corresponding to independent linear processes. This null hypothesis was realized by computing the causalities on a set of  $N = 19$  univariate Fourier surrogates. Under the hypothesis of no dependence between the processes, the probability of the data's causality value for a given pair of variables being the highest among the total 20 values available is  $p = 0.05$ , providing a convenient nonparametric test of causality.

### 3. Results

#### 3.1. Weighted Causality Networks

The reliability of weighted causality networks computed from a decade of stationary model data is shown in Figure 2 (for all methods and parameter values), along with the average similarity of the nonlinear network estimate by each method with the one obtained for the linear Granger causality method. The linear Granger causality shows the highest reliability, with the average Spearman's rank coefficient being  $\sim 0.6$ . The equiquantal binning method provided the most reliable network estimates for  $Q = 2$  ( $\bar{r} \sim 0.36$ ), with reliability generally decreasing for larger  $Q$ . The k-nearest neighbors algorithm provided even less reliable network estimates, with only weak dependence on the value of  $k$  and optimum reliability of  $\bar{r} \sim 0.33$  for  $k = 64$ .

**Figure 2.** Reliability of causality network detection using different causality estimators, and the similarity to linear causality network estimates using the Fourier surrogate model. For each estimator, six causality networks are estimated, one for each decade-long section of model stationary data (a Fourier surrogate realization of the original data). Black: the height of the bar corresponds to the average Spearman's correlation across all 15 pairs of decades. White: the height of the bar corresponds to the average Spearman's correlation of nonlinear causality network and linear causality network across 6 decades.



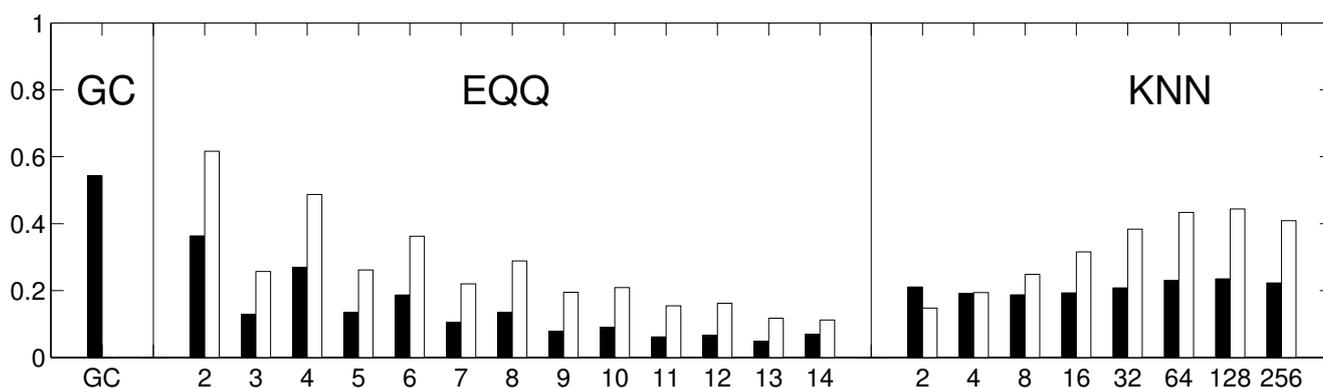
The causality networks constructed by each nonlinear method have been compared with the causality network obtained using the linear Granger causality analysis, see white bars in Figure 2. In general, the nonlinear causality networks have shown higher similarity to linear estimates than to nonlinear estimates for different sections of the stationary model time series. Interestingly, the parameter settings that optimized the reliability also provided the (almost) closest results to the linear methods. We have also observed generally lower reliability of the EQQ method for odd  $Q$ -values, an effect that will be investigated in detail elsewhere.

Figure 3 shows the results of an analogous analysis on original data rather than the stationary model. Note that here the computed causality network similarities reflect a combination of the (lack of) reliability of the methods and the real variability in the dynamical properties of the time series across time (*i.e.*, true changes in the causality pattern). The results are both qualitatively and quantitatively

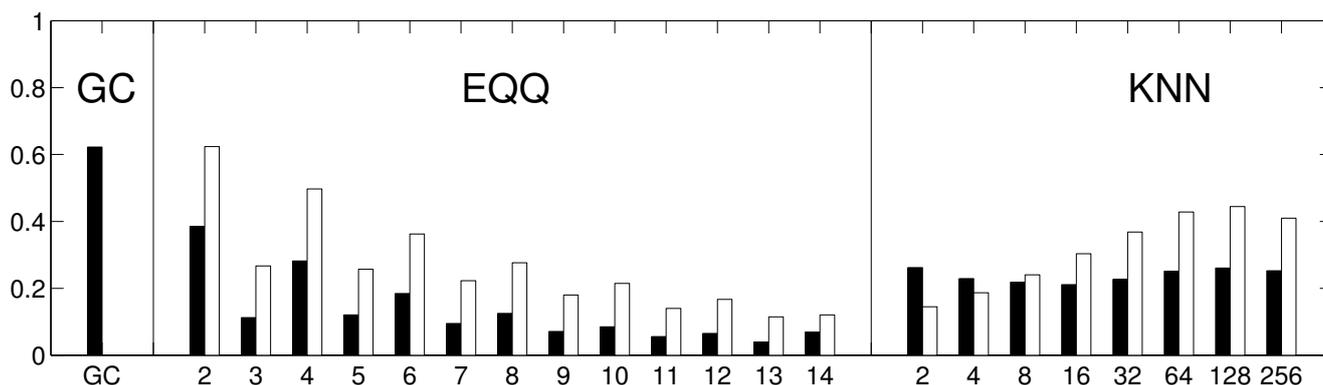
similar to those shown in Figure 2, suggesting that the true variability of the causal networks on this time scale is likely rather small compared with the coarseness of the causality assessment methods.

The results for other settings are shown in Figure 4 (use of multivariate AR(1) as the stationary model) and Figure 5 (6-day averages), generally confirming the main observations. However, some differences were observable. For instance, in the 6-days-averaged data, the reliability dependence of the kNN method on the  $k$ -parameter was more pronounced and peaked for a higher value of  $k = 256$ . The increase of reliability of the EQQ method for high  $Q$  was found to be spurious and is discussed in Section 4.

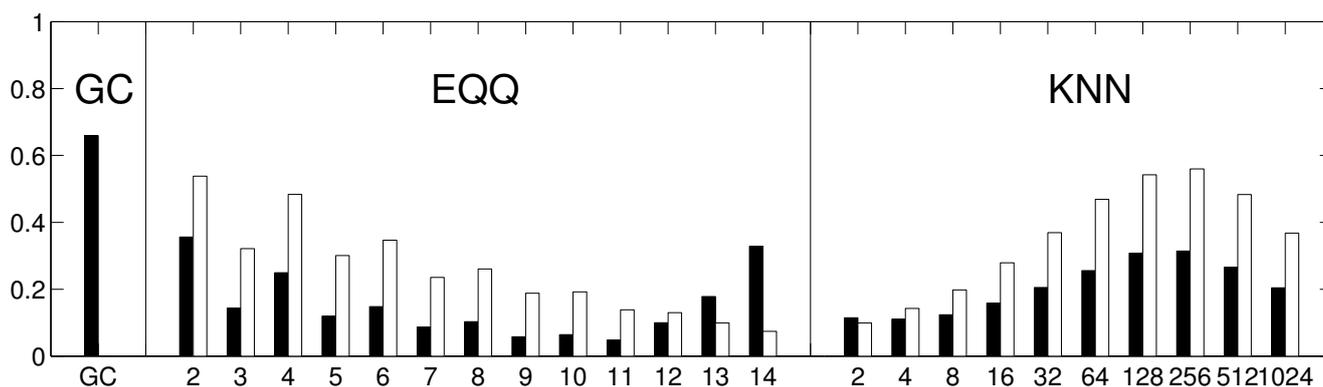
**Figure 3.** The variability of causality network detection using different causality estimators and the similarity to linear causality network estimates for the original data. For each estimator, six causality networks are estimated, one for each decade of the data. Black: the height of the bar corresponds to the average Spearman’s correlation across all 15 pairs of decades. White: the height of the bar corresponds to the average Spearman’s correlation of nonlinear causality network and linear causality network across 6 decades.



**Figure 4.** The reliability of causality network detection using different causality estimators and the similarity to linear causality network estimates for the stationary model constructed as multivariate AR(1) surrogate of the original data. For each estimator, six causality networks are estimated, one for each decade of modeled stationary data. Black: the height of the bar corresponds to the average Spearman’s correlation across all 15 pairs of decades. White: the height of the bar corresponds to the average Spearman’s correlation of nonlinear causality network and linear causality network across 6 decades.



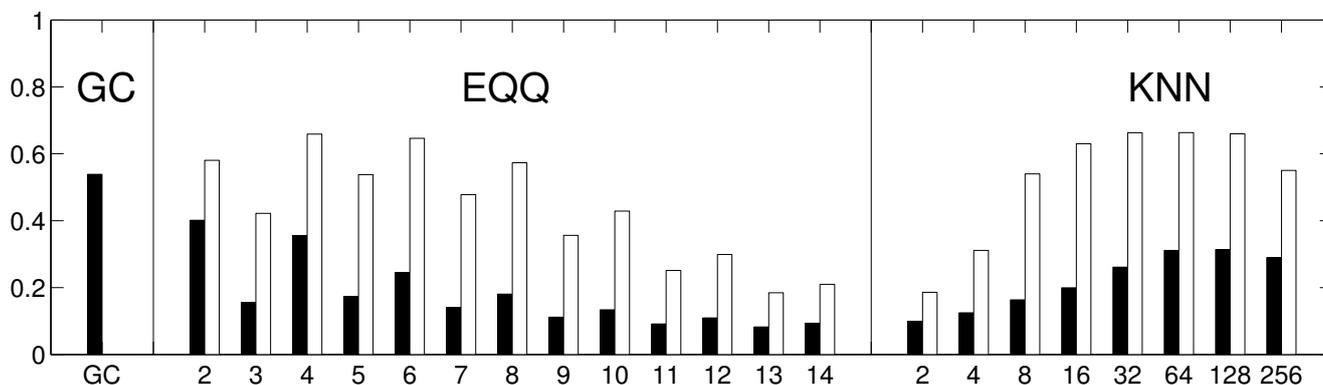
**Figure 5.** The reliability of causality network detection using different causality estimators and the similarity to linear causality network estimates for the stationary model constructed as multivariate AR(1) surrogate of the original data. For each estimator, six causality networks are estimated, each for a separate realization of the multivariate AR(1) process fitted to the original data. Black: the height of the bar corresponds to the average Spearman’s correlation across all 15 pairs of decades. White: the height of the bar corresponds to the average Spearman’s correlation of nonlinear causality network and linear causality network across 6 decades.



### 3.2. Unweighted Causality Networks

For unweighted causality networks, after thresholding to keep 1% of the strongest links, the network similarity was assessed by the Jaccard correlation coefficient. The results are plotted as in the previous figures, see Figure 6.

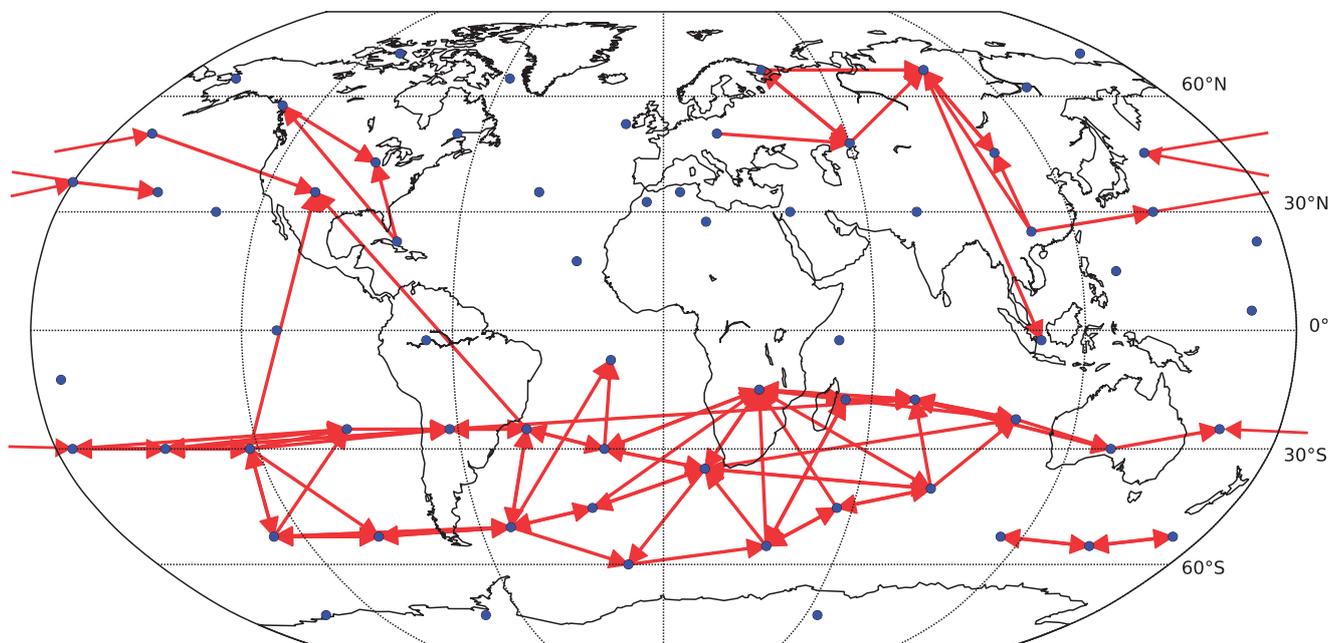
**Figure 6.** The reliability of causality network detection using different causality estimators and the similarity to linear causality network estimates. For each estimator, six causality networks are estimated, one for each decade of modeled stationary data. Black: the height of the bar corresponds to the average Jaccard similarity coefficient across all 15 pairs of decades. White: the height of the bar corresponds to the average Jaccard similarity coefficient of nonlinear causality network and linear causality network across 6 decades.



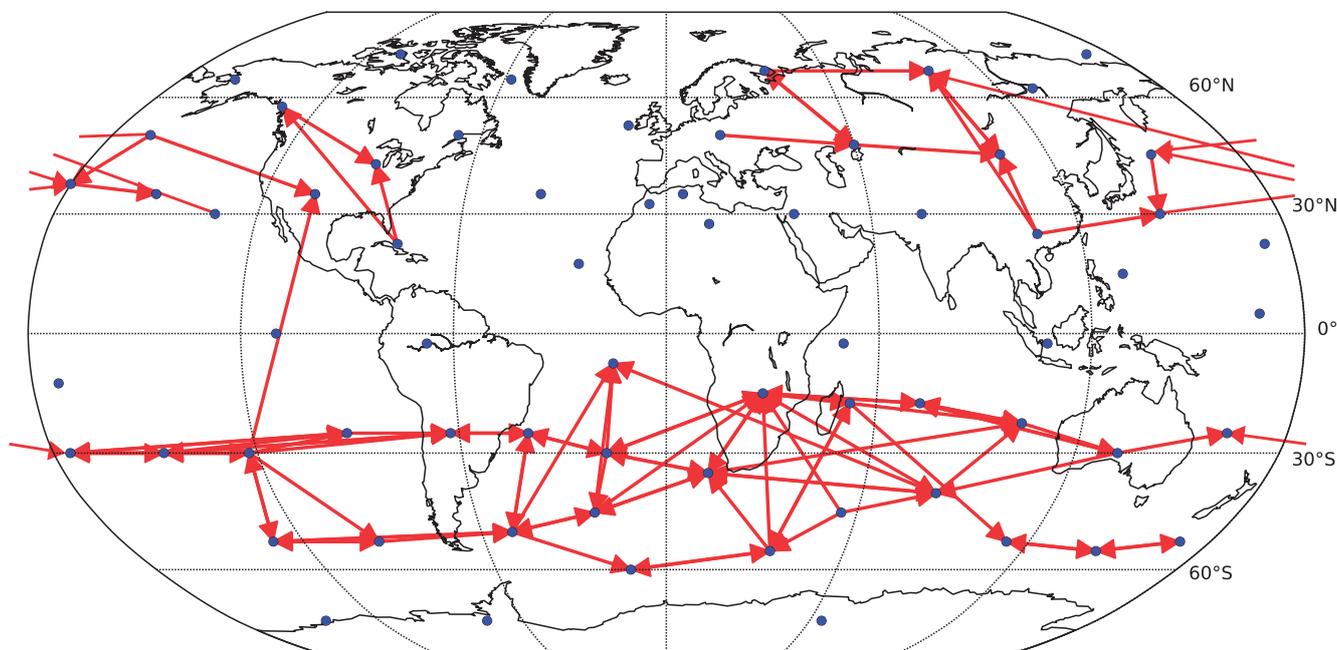
### 3.3. Components and Resulting Networks

To visualize the climatic networks, we first provide an overview of the localization of the networks in Figure 1 showing a parcellation of the globe with the components. As the variability among the decades has not been proven to be substantially higher in the data than in the stationary linear surrogate model (compare Figure 2 and Figure 3) and seems to mainly correspond to random fluctuations, the causality matrices are well represented by their average, shown below. For detailed results for each temporal window, see the Supplementary Material. In Figure 7, a graph of the 100 strongest links in the causality network computed using linear Granger causality and averaged over the six decades is provided. For comparison, the causality network obtained using the EQQ with  $Q = 2$  is shown in Figure 8. The overlap between the 100 strongest links in the average linear and nonlinear causality graph is relatively high—the graphs share 91 out of 100 links. In both graphs, the detected links are located predominantly in the extra-tropical regions. The prevailing eastward direction in the oceanic areas in these latitudes is in line with the expected circulation direction; this is manifested mostly in the north Pacific and southern Atlantic. However, many of the detected interactions are bidirectional. The long links crossing the equator are likely spurious.

**Figure 7.** Causality network obtained by averaging the results for the six decades (total time span 1948–2007) for decomposed data (67 components represented by center of mass). Only the 100 strongest links are shown. For each decade, the network was estimated by linear Granger causality.



**Figure 8.** Causality network obtained by averaging the results for the six decades (total time span 1948–2007) for decomposed data (67 components represented by center of mass). Only the 100 strongest links are shown. For each decade, the network was estimated by (nonlinear) transfer entropy using the equiqantal binning method with  $Q = 2$ .



#### 4. Discussion

The series of examinations provided evidence that both nonlinear and linear methods may be used to construct directed climate networks in a reliable way under a range of settings, with the basic linear Granger causality outperforming the studied nonlinear methods. On the side of nonlinear causality methods, we focused on the prominent family of methods based on the estimation of conditional mutual information in the form of transfer entropy. Two algorithms were used that represent key approaches of estimating conditional mutual information and have been extensively used and proven efficient on real-world data. Alternative approaches also exist, including the use of recurrence plots [31].

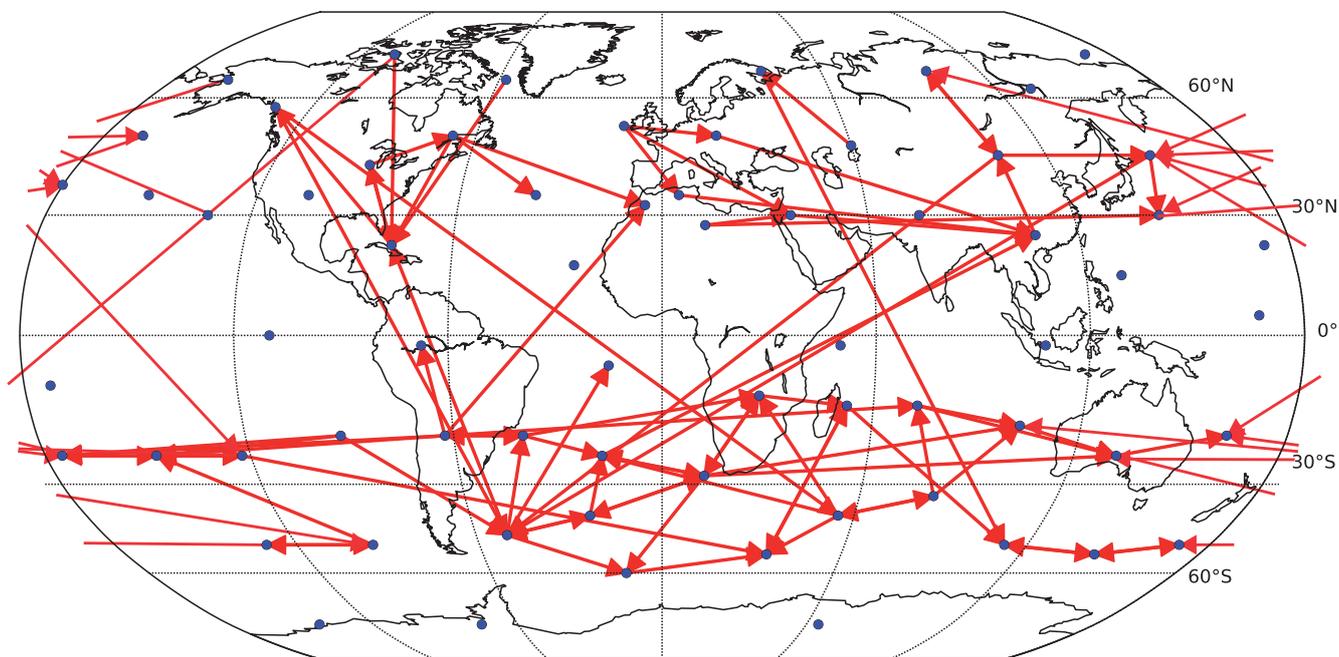
Indeed, for the sake of tractability, we have limited the investigation in several ways. As the motivation of the current study was to explore the rationale for use of linear/nonlinear methods in causal climate network construction and to provide a basis for more detailed analyses of climate networks, some simple pragmatic methodological choices have been made mainly for testing the method's differences, rather than studying specific phenomena.

Most prominently, linear Granger causality was chosen for its theoretical equivalence to transfer entropy (under the assumption of linearity), as this provides a fair comparison. However, the use of the strictly pairwise causality estimators suffers from severe inherent limitations. To give an example, a system consisting of three processes  $X, Y, Z$ , where  $Z$  drives both  $X$  and  $Y$ , but with different temporal lags, may erroneously show causal influence between  $X$  and  $Y$  even though these were not directly coupled. To deal with such situations, the concepts can be generalized to the multivariate case [32]. In

the case of linear Granger causality this leads to fitting a multivariate linear process including, apart from  $X$  and  $Y$ , also a (possibly multivariate)  $Z$ .

On the other side, there is a price for including too many variables in the model, as for short time series the theoretical advantage of including possible common causes may be outweighed by the numerical instability of fitting a higher order model. In Figure 9 we include an example result for applying the fully multivariate linear Granger causality, with the pairwise causality for each pair taking into account variability explained by all the other 65 component time series. Note that many of the links recovered by the simplified bivariate linear Granger causality are recovered (see Figure 7), however, the multivariate model leads to the detection of many long-ranging links crossing the equator, which are likely spurious and a result of lower reliability of the algorithm. Quantitatively, the average similarity (Spearman's rank correlation coefficient) among the multivariate linear Granger causality matrices obtained for six decades of linear surrogate data was  $\hat{r}_{sp} = 0.45$ , compared with  $\hat{r}_{sp} = 0.57$  obtained for the simplified bivariate linear Granger causality. We do not provide analogous results for nonlinear methods, as the estimation of information-theoretical functionals in the fully multivariate setting is computationally prohibitively difficult.

**Figure 9.** Causality network obtained by averaging the results for the six decades (total time span 1948–2007) for decomposed data (67 components represented by center of mass). Only the 100 strongest links are shown. For each decade, the network was detected by the fully multivariate linear Granger causality.



Similarly, the assumption of a single possible lag (1 time step of 1 or 6 days respectively in our investigation) may not be suitable in real context, although at least the relative reliability of different methods may not be strongly affected by this within reasonable range of parameters. For precise geophysical interpretation, fitting an appropriate model order and allowing multiple lags would be warranted.

In general, estimation of these generalized causality patterns from relatively short time series is technically challenging, particularly in the context of nonlinear, information-theory based causality measures, due to the exponentially increasing dimension of probability distributions to be estimated. However, recent work has provided promising approaches to tackle this curse of dimensionality by decomposing TE into low-dimensional contributions [32]. For theoretical and numerical considerations on how a causal coupling strength can be defined in the multivariate context, see [33]. Assessment of the performance of these novel methods in the detection of causal climate networks is a subject of future work. For completeness we mention that apart from the time-domain treatment of causality, the whole problem can also be reformulated in the spectral domain, leading to frequency-resolved causality indices such as partial directed coherence (PDC [34]) or Directed Transfer Function (DTC [35]).

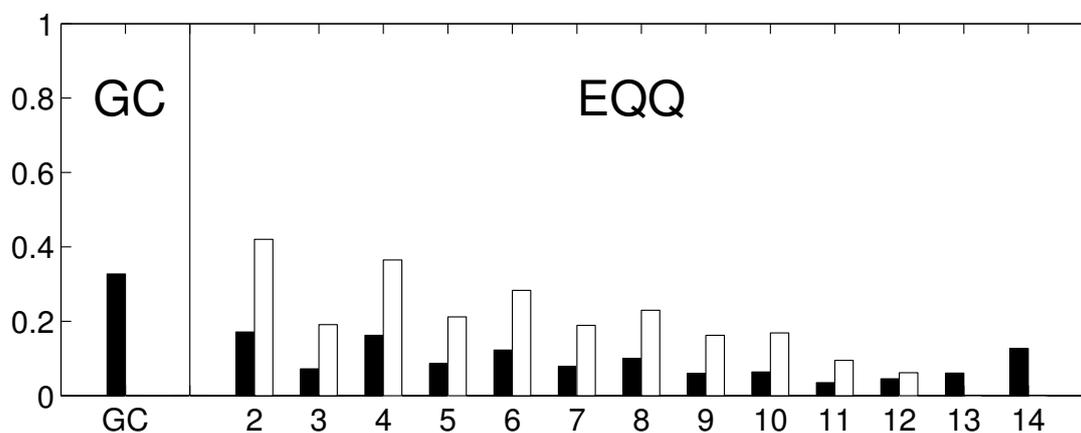
Our study also shows some key properties of the conditional mutual information estimates. Note that for instance for the 6-days networks (Figure 5), the EQQ method reliability increases again for  $Q \geq 11$ , however, the similarity to the linear estimate further decreases. A direct check of the causality networks shows that they tend to a trivial column-wise structure, with the intensities for each column highly correlated to the autoregressive coefficient of the given time series. This corresponds to a manifestation of a dominant autocorrelation-dependent bias in the EQQ estimator for too high  $Q$  values (note that  $Q = 14$  corresponds to less than one time point per average in a 4D bin, an unsuitable sampling of the space for effective probability distribution function approximation).

Reconstruction of networks directly from gridded climatic field data is challenging and perhaps not always the best approach, for reasons including efficient computation and visualization of the results. Instead, we apply here a decomposition of the data in order to get the most important components, *i.e.*, by using a varimax-rotated principle component analysis. This provides a useful dimension-reduction of the studied problem. In particular, the gridded data does not reflect the real climate subsystems, which may be better approximated by the decomposition modes. However, as the decomposition is an implicit (weighted) coarse-graining, the detected difference between the linear and nonlinear methods may be different from that in the original time series data. In particular, the spatial averaging may increase the reliability of both approaches by suppressing noise, but also suppress any highly spatially localized heterogeneous patterns of both nonlinear and linear character. This might be reflected in the specific results of the paper, in particular the obtained quantitative reliability estimates.

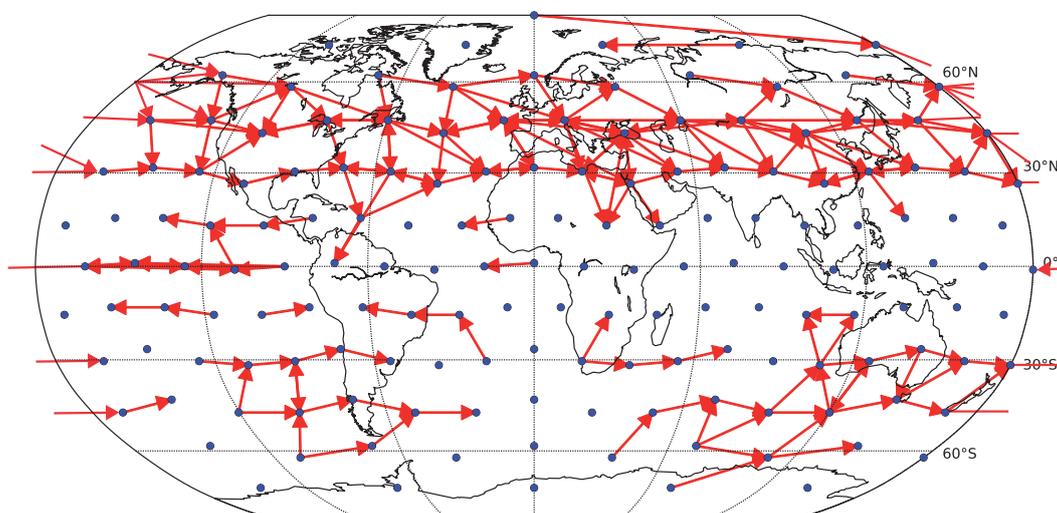
To assess how the results generalize to gridded data, we carried out a partial analysis using an approximately equidistant grid consisting of 162 spatial locations. The time series were preprocessed by removing anomalies in mean as well as in variance. The reliability assessment for the gridded data shows qualitatively similar results, see Figure 10. As for the component data, the linear Granger causality has proven to be most reliable, and also relatively reliable results are obtained for nonlinear conditional mutual information with a small number of bins ( $Q = 2$ ). The overall decrease in reliability is probably due to a combination of larger matrix size and less sophisticated dimensionality reduction strategy. The kNN algorithm was not assessed on this data to limit additional computational time demands. The average graphs for linear and nonlinear causality are shown in Figures 11 and 12 respectively. Interestingly, the spatial sparsity may be more suited to the 1-day lag causality estimation, as the detected average causality graphs hold structure that is easier to interpret. In particular, the direction of detected information flow in the extra-tropical areas approximately between  $30^\circ$  and  $60^\circ$  on both hemispheres

almost without exception resembles the prevailing eastward direction of air circulation in the Ferrell cells (mid-latitude westerlies). On the other side, in particular the graph in Figure 11, the direction of detected information flow corresponding to the strongest 200 links for the more robust linear Granger causality method also includes many arrows aiming westwards in the tropical regions (e.g., in the Pacific), which may correspond to the easterly trade winds in the regions corresponding to Hadley cells.

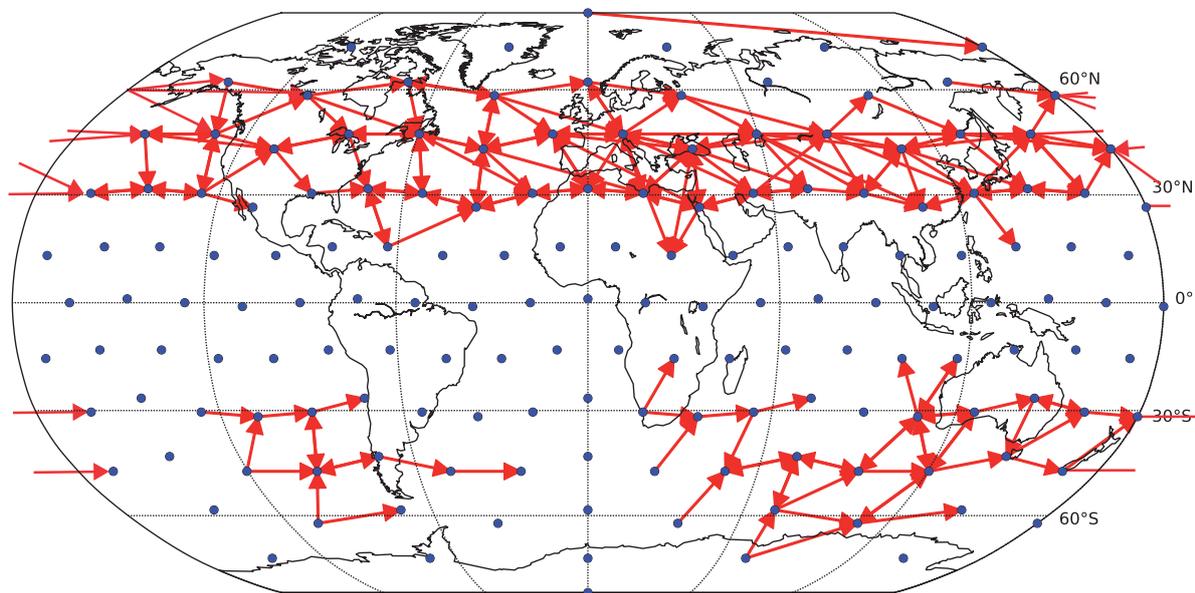
**Figure 10.** The reliability of causality network detection using different causality estimators and the similarity to linear causality network estimates for the Fourier surrogates model. For each estimator, six causality networks are estimated, one for each decade-long section of the model stationary data (a Fourier surrogate realization of the original data). Black: the height of the bar corresponds to the average Spearman’s correlation across all 15 pairs of decades. White: the height of the bar corresponds to the average Spearman’s correlation of nonlinear causality network and linear causality network across 6 decades.



**Figure 11.** Causality network obtained by averaging the results for the six decades (total time span 1948–2007) for gridded data (162 spatial locations). Only the 200 strongest links are shown. For each decade, the network was estimated by linear Granger causality.



**Figure 12.** Causality network obtained by averaging the results for the six decades (total time span 1948–2007) for gridded data (162 spatial locations). Only the 200 strongest links are shown. For each decade, the network was estimated by (nonlinear) conditional mutual information, using the equiangular binning method with  $Q = 2$ .



In general, the differences in reliability may have important consequences for the detectability of causal links as well as their changes. Of course, theoretically this is not necessary since a nonlinear system may have links that have negligible linear causality but strong enough nonlinear causality fingerprint. Recent works have proposed approaches for the explicit quantification of the nonlinear contribution of equal-time dependence [36], applied to neuroscientific [36] as well as climate data [11]. However, the generalization to higher dimensional information-theoretic functionals is not straightforward and is the subject of ongoing work. Thus, we give here at least an illustrative example of the trade-off between the generality of transfer entropy and the higher reliability of linear Granger causality: using the original 67 components data (divided into 6 decadal sections, see Section 2), the basic statistical test at the 5% significant level (described in Section 2) marked on average 1592 links as statistically significant (out of 4422 possible), while the EQQ method showed in general a lower number of significant links, depending on the parameter value in a way similar to the reliability estimate, with a maximum of 983 significant links for  $Q = 2$  and a minimum of 287 significant links for  $Q = 13$ . Note that given the significance level of the test, on average  $221 \approx 4422 \times 0.05$  significant links would be expected to appear by chance in a collection of completely unrelated processes.

## 5. Conclusions

A meaningful interpretation of climate networks and their observed temporal variability requires knowledge and minimization of the methodological limitations. In the present work, we discussed the problem of the reliability of network construction from time series of finite length, quantitatively assessing the reliability for a selection of standard bivariate causality methods. These included two major algorithms for estimating transfer entropy with a wide range of parameter choices, as well as linear

Granger causality analysis, which can be understood as linear approximation of transfer entropy. Overall, the causality methods provided reproducible estimates of climate causality networks, with the linear approximations outperforming the studied nonlinear methods in reliability. Interestingly, optimizing the nonlinear methods with respect to reliability has led to improved similarity of the detected networks to those discovered by linear methods, in line with the hypothesis of near-linearity of the investigated climate reanalysis data, in particular the surface air temperature time series.

The latter hypothesis regarding the surface air temperature has been supported by the study in [11] which extended the older results of [37] who tested for possible nonlinearity in the dynamics of the station (Prague-Klementinum) SAT time series and found that the dependence between the SAT time series  $x(t)$  and its lagged twin  $x(t + \tau)$  was well-explained by a linear stochastic process. This result about a linear character of the temporal evolution of SAT time series, as well as the results of the present study, cannot be understood as arguments for a linear character of atmospheric dynamics *per se*. Rather, these results characterize the properties of measurement or reanalysis data at a particularly coarse level of resolution. The data, thus, reflect a spatially and temporally averaged mixture of dynamical processes. For instance, a closer look on the dynamics of the specific temporal scales in temperature and other meteorological data has led to the identification of oscillatory phenomena with nonlinear behavior, exhibiting phase synchronization [38–42]. Also the leading modes of atmospheric variability exhibit nonlinear behavior [43,44] and can be inferred the directionality of coupling with conditional mutual information, in a nonlinear way [45].

Further work is needed to assess the usability and advantages of more sophisticated, recently proposed causality estimation methods. The current work also provides an important step towards the reliable characterization of climate networks and the detection of potential changes over time.

## Acknowledgements

This study is supported by the Czech Science Foundation, Project No. P103/11/J068 and by the DFG grant No. KU34-1. We thank our anonymous reviewers for their thorough evaluation and constructive recommendations for improving the manuscript.

## References

1. Newman, M.E.J. The structure and function of complex networks. *SIAM Rev.* **2003**, *45*, 167–256.
2. Boccaletti, S.; Latora, V.; Moreno, Y.; Chavez, M.; Hwang, D.U. Complex networks: Structure and dynamics. *Phys. Rep.* **2006**, *424*, 175–308.
3. Tsonis, A.; Roebber, P. The architecture of the climate network. *Physica A* **2004**, *333*, 497–504.
4. Tsonis, A.A.; Swanson, K.L.; Roebber, P.J. What do networks have to do with climate? *Bull. Am. Meteorol. Soc.* **2006**, doi:10.1175/BAMS-87-5-585.
5. Yamasaki, K.; Gozolchiani, A.; Havlin, S. Climate networks around the globe are significantly affected by El Nino. *Phys. Rev. Lett.* **2008**, doi: 10.1103/PhysRevLett.100.228501.
6. Malik, N.; Bookhagen, B.; Marwan, N.; Kurths, J. Analysis of spatial and temporal extreme monsoonal rainfall over South Asia using complex networks. *Clim. Dyn.* **2012**, *39*, 971–987.

7. Steinhäuser, K.; Ganguly, A.R.; Chawla, N.V. Multivariate and multiscale dependence in the global climate system revealed through complex networks. *Clim. Dyn.* **2012**, *39*, 889–895.
8. Steinhäuser, K.; Chawla, N.V.; Ganguly, A.R. Complex Networks in Climate Science: Progress, Opportunities and Challenges. In Proceedings of the 2010 Conference on Intelligent Data Understanding, CIDU 2010, Mountain View, CA, USA, 5–6 October 2013; pp. 16–26.
9. Donges, J.F.; Zou, Y.; Marwan, N.; Kurths, J. The backbone of the climate network. *EPL* **2009**, doi:10.1209/0295-5075/87/48007.
10. Donges, J.F.; Zou, Y.; Marwan, N.; Kurths, J. Complex networks in climate dynamics. *Eur. Phys. J.* **2009**, *174*, 157–179.
11. Hlinka, J.; Hartman, D.; Vejmelka, M.; Novotna, D.; Paluš, M. Non-linear dependence and teleconnections in climate data: Sources, relevance, nonstationarity. *Clim. Dyn.* **2012**, doi:10.1007/s00382-013-1780-2.
12. Ebert-Uphoff, I.; Deng, Y. Causal discovery for climate research using graphical models. *J. Clim.* **2012**, *25*, 5648–5665.
13. Ebert-Uphoff, I.; Deng, Y. A new type of climate network based on probabilistic graphical models: Results of boreal winter versus summer. *Geophys. Res. Lett.* **2012**, *39*, 5648–5665.
14. Granger, C.W. Investigating causal relations by econometric model and cross-spectral methods. *Econometrica* **1969**, *37*, 424–438.
15. Vejmelka, M.; Palus, M. Inferring the directionality of coupling with conditional mutual information. *Phys. Rev. E* **2008**, doi:10.1103/PhysRevE.77.026214.
16. Schreiber, T. Measuring information transfer. *Phys. Rev. Lett.* **2000**, *85*, 461–464.
17. Kistler, R.; Kalnay, E.; Collins, W.; Saha, S.; White, G.; Woollen, J.; Chelliah, M.; Ebisuzaki, W.; Kanamitsu, M.; Kousky, V.; *et al.* The NCEP-NCAR 50-year reanalysis: Monthly means CD-ROM and documentation. *Bull. Am. Meteorol. Soc.* **2001**, *82*, 247–267.
18. Kalnay, E.; Kanamitsu, M.; Kistler, R.; Collins, W.; Deaven, D.; Gandin, L.; Iredell, M.; Saha, S.; White, G.; Woollen, J.; *et al.* The NCEP/NCAR 40-year reanalysis project. *Bull. Am. Meteorol. Soc.* **1996**, *77*, 437–471.
19. Wiener, N. *Modern Mathematics for Engineers*; McGraw-Hill: New York, NY, USA, 1956; chapter The theory of prediction, pp. 165–190.
20. Ding, M.; Chen, Y.; Bressler, S.L. Granger Causality: Basic Theory and Application to Neuroscience. In *Handbook of Time Series Analysis*; Wiley-VCH Verlag GmbH & Co. KGaA: Weinheim, Germany, 2006; pp. 437–460.
21. Geweke, J.F. Measurement of linear dependence and feedback between multiple time series. *J. Am. Stat. Assoc.* **1982**, *79*, 907–915.
22. Geweke, J.F. Measures of conditional linear dependence and feedback between time series. *J. Am. Stat. Assoc.* **1984**, *79*, 907–915.
23. Barnett, L.; Barrett, A.B.; Seth, A.K. Granger causality and transfer entropy are equivalent for gaussian variables. *Phys. Rev. Lett.* **2009**, *103*, doi:10.1103/PhysRevLett.103.238701.
24. Palus, M.; Albrecht, V.; Dvorak, I. Information theoretic test for nonlinearity in time series. *Phys. Lett.* **1993**, *175*, 203–209.

25. Frenzel, S.; Pompe, B. Partial mutual information for coupling analysis of multivariate time series *Phys. Rev. Lett.* **2007**, *99*, doi:10.1103/PhysRevLett.99.204101
26. Kaiser, H. The varimax criterion for analytic rotation in factor analysis. *Psychometrika* **1958**, *23*, 187–200.
27. Schwarz, G. Estimating the dimension of a model. *Ann. Stat.* **1978**, *5*, 461–464.
28. Benjamini, Y.; Hochberg, Y. Controlling the false discovery rate: A practical and powerful approach to multiple testing. *J. R. Stat. Soc. Series B-Methodol.* **1995**, *57*, 289–300.
29. Prichard, D.; Theiler, J. Generating surrogate data for time series with several simultaneously measured variables. *Phys. Rev. Lett.* **1994**, *73*, 951–954.
30. Palus, M. Detecting phase synchronization in noisy systems. *Phys. Lett.* **1997**, *235*, 341–351.
31. Zou, Y.; Romano, M.C.; Thiel, M.; Marwan, N.; Kurths, J. Inferring indirect coupling by means of recurrences. *Int. J. Bifurc. Chaos* **2011**, *21*, 1099–1111.
32. Runge, J.; Heitzig, J.; Petoukhov, V.; Kurths, J. Escaping the curse of dimensionality in estimating multivariate transfer entropy. *Phys. Rev. Lett.* **2012**, doi:10.1103/PhysRevLett.108.258701.
33. Runge, J.; Heitzig, J.; Marwan, N.; Kurths, J. Quantifying causal coupling strength: A lag-specific measure for multivariate time series related to transfer entropy. *Phys. Rev.* **2012**, doi:10.1103/PhysRevE.86.061121.
34. Baccala, L.A.; Sameshima, K. Partial directed coherence: A new concept in neural structure determination. *Biol. Cybern.* **2001**, *84*, 463–474.
35. Kaminski, M.; Ding, M.; Truccolo, W.A.; Bressler, S.L. Evaluating causal relations in neural systems: Granger causality, directed transfer function and statistical assessment of significance. *Biol. Cybern.* **2001**, *85*, 145–157.
36. Hlinka, J.; Palus, M.; Vejmelka, M.; Mantini, D.; Corbetta, M. Functional connectivity in resting-state fMRI: Is linear correlation sufficient? *NeuroImage* **2011**, *54*, 2218–2225.
37. Palus, M.; Novotna, D. Testing for nonlinearity in weather records. *Phys. Lett.* **1994**, *193*, 67–74.
38. Palus, M.; Novotna, D. Enhanced Monte Carlo Singular System Analysis and detection of period 7.8 years oscillatory modes in the monthly NAO index and temperature records. *Nonlinear Process. Geophys.* **2004**, *11*, 721–729.
39. Palus, M.; Novotna, D. Quasi-biennial oscillations extracted from the monthly NAO index and temperature records are phase-synchronized. *Nonlinear Process. Geophys.* **2006**, *13*, 287–296.
40. Palus, M.; Novotna, D. Phase-coherent oscillatory modes in solar and geomagnetic activity and climate variability. *J. Atmos. Solar-terr. Phys.* **2009**, *71*, 923–930.
41. Palus, M.; Hartman, D.; Hlinka, J.; Vejmelka, M. Discerning connectivity from dynamics in climate networks. *Nonlinear Process. Geophys.* **2011**, *18*, 751–763.
42. Feliks, Y.; Ghil, M.; Robertson, A.W. Oscillatory climate modes in the eastern mediterranean and their synchronization with the north atlantic oscillation. *J. Clim.* **2010**, *23*, 4060–4079.
43. Boucharel, J.; Dewitte, B.; du Penhoat, Y.; Garel, B.; Yeh, S.W.; Kug, J.S. ENSO nonlinearity in a warming climate. *Clim. Dyn.* **2011**, *37*, 2045–2065.
44. Osprey, S.M.; Ambaum, M.H.P. Evidence for the chaotic origin of northern annular mode variability. *Geophys. Res. Lett.* **2011**, doi:10.1029/2011GL048181.

45. Mokhov, I.I.; Smirnov, D.A.; Nakonechny, P.I.; Kozlenko, S.S.; Seleznev, E.P.; Kurths, J. Alternating mutual influence of El-Nino/Southern Oscillation and Indian monsoon. *Geophys. Res. Lett.* **2011**, doi:10.1029/2010GL045932.

© 2013 by the authors; licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution license (<http://creativecommons.org/licenses/by/3.0/>).