# Mining Novel Multivariate Relationships in Time Series Data Using Correlation Networks

Saurabh Agrawal, Michael Steinbach, Daniel Boley, Snigdhansu Chatterjee, Gowtham Atluri, Anh The Dang, Stefan Liess, and Vipin Kumar

Abstract—In many domains, there is significant interest in capturing novel relationships between time series that represent activities recorded at different nodes of a highly complex system. In this paper, we introduce multipoles, a novel class of linear relationships between more than two time series. A multipole is a set of time series that have strong linear dependence among themselves, with the requirement that each time series makes a significant contribution to the linear dependence. We demonstrate that most interesting multipoles can be identified as cliques of negative correlations in a correlation network. Such cliques are typically rare in a real-world correlation network, which allows us to find almost all multipoles in discovering new physical phenomena in two scientific domains: climate science and neuroscience. In particular, we discovered several multipole relationships that are reproducible in multiple other independent datasets and lead to novel domain insights.

Index Terms—multivariate linear patterns; correlation mining; spatio-temporal; climate teleconnections; fMRI

## **1** INTRODUCTION

I N many domains, understanding the relationships between time series is essential for obtaining actionable insights. For instance, in climate science, pressure dipoles, which are pairs of locations with strong negative correlations in their Sea Level Pressure time series, have been extensively studied, and have been linked with anomalous weather events all over the globe such as forest fires, hurricanes etc. [1], [2], [3]. Similarly, in neuroscience, researchers have discovered pairs of brain regions that exhibit positive correlations between their activity time series. These correlated time series represent brain regions exhibiting synergistic activity [4].

In this work, we define a novel class of linear relationships, called **multipoles**, that involve more than two time series, also referred to as variables. We say that a set of variables is a multipole if i) the variables show strong linear dependence, and ii) each variable makes a significant contribution to the linear dependence, i.e., excluding any of the variables from the set significantly weakens the strength of the linear dependence among the remaining variables. We define linear dependence in terms of the variance of a linear combination of the standardized (zero mean, unit variance) vectors. We measure the strength of the linear dependence of a set of time series by the variance of their least variant linear combination, i.e., a (unit-length) linear combination that has the minimum variance. The smaller the variance of this linear combination, the more constant (less variable) the linear combination, and thus higher the linear dependence, and vice-versa. We define the contribution of each variable

to the linear dependence as the reduction in the strength of the linear dependence of the smaller set when this variable is removed. We define linear gain as the minimum contribution of included variables to the linear dependence.

We next illustrate multipoles with a real-world example. Consider a set *S* of three time series  $T_1$ ,  $T_2$ , and  $T_3$ , shown in Figure 2(a), which capture the traffic volume on three different roads in Minnesota as shown in Figure 1. The three bottommost plots of Figure 2(b) show the least variant linear combination for each possible pair, while the top plot of Figure 2(b) shows the least variant linear combination for all three. Note that the linear combination of all three time series has a variance of only 0.08. Thus, we say that there exists a strong, although not perfect, linear dependence among the above three time series.

Further, note that each of the three time series forms a crucial component of the relationship, as excluding any one of them will significantly weaken the strength of the linear dependence among the remaining two time series. For example, if  $T_3$  is excluded from the set, the variance of the least variant linear combination  $Z_{12}$  of  $T_1$  and  $T_2$  turns out to be 0.33, which is much higher than the variance 0.08 of Z, thus indicating a significant contribution from  $T_3$ . Similarly, if we exclude  $T_1$  or  $T_2$  instead of  $T_3$ , the variance of least variant combinations become 0.58 and 0.74. (See the three bottommost plots of Figure 2(b).) The linear gain for  $\{T_1, T_2, T_3\}$  is min $\{0.33, 0.58, 0.74\} - 0.08 = 0.25$ .

An explanation for the relationships illustrated above can be provided by the notion of conservation of flow. Two of the time series,  $T_1$  and  $T_2$ , were observed on the roads that act as major tributaries to the highway where  $T_3$  is being observed. All the southbound traffic coming from the tributaries is likely to merge at the highway, thus leading to a strong linear dependence among the three time series. Omitting any single time series leads to a weaker linear dependence because, unlike the highway, traffic flow

Saurabh Agrawal, Michael Steinbach, Daniel Boley, Snigdhanshu Chatterjee, Stefan Liess, and Vipin Kumar, are with University of Minnesota, Minneapolis, MN.

E-mail: [agraw066,stei0062,boley,chatt019,liess,kumar001]@umn.edu

Gowtham Atluri and Anh The Dang are with University of Cincinnati, Cincinnati, OH. E-mail: atlurigm@ucmail.uc.edu

in the tributaries significantly changes during weekends (see peaks and drops every Sunday), which weakens the strength of their pairwise linear relationships with highway. However, the simultaneous rise and fall in traffic on both tributaries complements each other, resulting in a stronger linear dependence among the three time series.

Multipoles can also be seen in other domains. For instance, we used the techniques presented in this paper to find several previously unknown multipole relationships between climate variables observed at more than two distant locations. Similarly, in neuroscience, we found novel multipole relationships between different brain regions that are triggered by specific visual and auditory stimuli. Some of these relationships were found to have insightful domain interpretations. (See Section 6.)

The relationships and conditional independence of two or more variables have been widely studied in different contexts using a variety of techniques including regression models [5], [6], [7], [8], PCA-based approaches [9], [10], structure learning methods [11], [12], correlation network analysis [13], [14], [15], [16], etc. However, as discussed in Section 3, multipole relationships as defined here can be viewed as novel in the sense that none of these previous approaches can be used to find the type of relationship represented by multipoles in the data.

A naïve approach to find all multipoles in a time series dataset would be to enumerate all possible combinations of time series and measure the strength of their linear dependence and linear gain. However, this is computationally infeasible due to the combinatorial nature of the search space. In this paper, we present an efficient approach to find multipoles.

This approach formulates the multipole-search problem as a clique-enumeration problem in a correlation network, where each node represents a time series, and the weight of an edge between two nodes represents the strength of the linear correlation between the corresponding time series. Our proposed problem formulation is motivated by the following two key empirical observations: the upper limit on the linear gain of a multipole is dependent on i) the size of multipole, and ii) the maximum correlation strength among two variables in a multipole.

Leveraging these observations, we propose a novel Clique Based Multipole Search (CoMEt) approach to find most interesting multipoles in a time series dataset. The central idea of the approach is to identify and restrict the search for multipoles to family of subsets, which we refer to as 'promising candidates' that are more likely to exhibit multipole relationships with stronger linear gain between their members. Using the above empirical observations, we show that most promising candidates of multipoles with high linear gain appear in a correlation network either as negative cliques, i.e. sets with all pairwise correlations being negative, or negative-equivalent cliques, that can be transformed into negative cliques by flipping the signs of one or more of its member variables (see Definitions 11 and 12 for further details). The number of promising candidates typically turns out to be much smaller in scenarios where multipoles with high linear gain are desired, thereby contributing to the remarkably high computational efficiency of our approach, although with some loss of completeness in



Fig. 1: A multipole among daily traffic time series  $T_1$ ,  $T_2$ , and  $T_3$  observed for three roads in Minneapolis



Fig. 2: (a) Standardized (zero mean and unit variance) daily traffic volume for the three roads shown in Figure 1. (b) Time series of their linear combinations

the final output.

Furthermore, we propose CoMEtExtended, a more generalized version of CoMEt approach, where we redefine what constitutes a promising candidate. In particular, CoMEtExtended involves an additional parameter that can be tuned to expand or prune the scope of promising candidates beyond negative cliques and negative-equivalent cliques, which allows one to achieve a better trade-off between computational efficiency and completeness at different thresholds of linear gain (see Section 4.5 for further details).

Our paper makes several key contributions: 1) We formally define multipole, a novel relationship in time series data and devise measures to quantify its interestingness. 2) We formulate a novel and computationally efficient patternmining approach to find most interesting multipoles in a time series dataset. 3) Further, we propose an empirical framework to evaluate discovered multipoles that includes an empirical procedure to assess the statistical significance of multipoles. 4) Using our proposed framework, we demonstrate the relevance of multipoles to two scientific domains: climate science and neuroscience.

## 2 DEFINITIONS AND NOTATIONS

Let  $S = \{X_1, X_2, ..., X_k\}$  denote a set of k standardized (zero mean, unit variance) time series observed over Tconsecutive timestamps. (We will use the terms 'time series' and 'variables' interchangeably in what follows.) Also, let **X** be the corresponding  $T \times k$  data matrix and let  $\Sigma = \mathbf{X}^T \mathbf{X}$ be the  $k \times k$  covariance matrix. Since all the variables (time series) are standardized, the correlation matrix and covariance matrices will be exactly the same for X. We next define a few measures on a set of variables which will be eventually used to formally define a multipole.

- **Definition 1.** A Normalized Linear Combination (NLC) refers to a linear combination with normalized weights. Specifically, for a given vector  $l \in \mathbb{R}^k$  with  $||l||_2 = 1$ , a normalized linear combination of variables in the set, S, is given as  $Z_S = \mathbf{X}l$ .
- **Definition 2.** Given a set of variables  $S = \{X_1, X_2, ..., X_k\}$ , let  $Z_S^*$  be the **Least Variant Normalized Linear Combination** (LVNLC) of variables, i.e., the NLC that has the least variance across the *T* observations. Formally,

$$Z_S^* = \mathbf{X}l^*$$
 where  $l^* = \underset{l \in \mathbb{R}^k, ||l||_2=1}{\operatorname{argmin}} var(\mathbf{X}l).$ 

In linear algebra, a set  $S = \{X_1, X_2, ..., X_k\}$  of variables is said to be linearly dependent if there exists a linear combination of the variables that equals the 0 vector. Since we have subtracted the mean from our vectors (time series), any linear combination will also have a mean of zero, and thus the variance of a linear combination of S is equivalent to finding the L2 norm of the linear combination. Hence, the LVNLC of S finds the linear combination of S that is the closest approximation to the 0 vector in terms of the L2 norm. If the linear combination of our time series (vectors) is perfectly constant, i.e., the time series vectors are linearly dependent, the variance of LVNLC will be exactly zero. In the other extreme case, when all the variables are mutually orthogonal to each other, the variance of LVNLC will be equal to 1. (This is because we have standardized the times series in S to have unit variance.) Thus, variance of LVNLC can be used as an inverse indicator of the strength of linear dependence. Based on this observation, we define the linear dependence of a set as follows:

## **Definition 3 (Linear Dependence:).** The linear dependence, $\sigma_S$ , of a set of vectors S is given by $1 - var(Z_S^*)$ .

Relation between  $\sigma_S$  and least eigenvalue of correlation matrix: Notice that on performing the eigenvalue decomposition of correlation matrix  $\Sigma$  of the variables in the set, the eigenvalues so obtained are equal to the variances of the projections of the data along their corresponding eigenvectors. Since  $Z_S^*$  corresponds to the direction of least variance,  $Z_S^*$  is nothing but the eigenvector corresponding to the least eigenvalue  $\lambda_{min}$  of  $\Sigma$ . Thus, the variance of  $Z_S^*$ is exactly equal to the least eigenvalue of  $\Sigma$  and therefore,

$$\sigma_S = 1 - \lambda_{min} \tag{1}$$

Before proceeding further, the following two properties of linear dependence are noteworthy:

*Lemma* 1. For any set S of standardized variables,  $\sigma_S \in [0, 1]$ .

*Lemma* **2**. The linear dependence of a set *S* is always less than or equal to that of its supersets.

**Proof**: See the Supplemental material. 
$$\Box$$

Although linear dependence indicates a strong relationship among the variables, it does not exclude the presence of irrelevant variables in the set. For instance, let S = $\{X_1, X_2, X_3, X_4\}$  be a set of linearly dependent variables with the linear relation being  $X_1 + X_2 + X_3 + 0X_4 = 0$ . Although the four variables are linearly dependent,  $X_4$  is an irrelevant variable and can be pruned from S without weakening the linear dependence among remaining variables. Hence, to avoid irrelevant variables in the pattern, we next propose a measure called *linear gain* that checks the minimum contribution from all member variables to the linear dependence of the set.

**Definition 4 (Linear Gain:).** The linear gain of a set S with |S| > 2 is measured as the gain in the linear dependence of S with respect to one of its proper subsets S' that has strongest linear dependence. Mathematically, we can write linear gain of S as

$$\Delta \sigma_S = \sigma_S - \max_{S' \subset S} \sigma_{S'} \tag{2}$$

From Lemma 2, we get that the linear dependence of a set is always greater than that of its subsets, which implies that the linear gain of a set will always be positive. Furthermore, we can say that the subset with strongest linear dependence will be of size |S| - 1. Thus, the linear gain can be more precisely written as

$$\Delta \sigma_S = \sigma_S - \max_{X_i \in S} \sigma_{S-X_i} \tag{3}$$

Higher values of linear gain imply that a more significant drop in linear dependence would be observed if any one of the variables are excluded from the set, thereby ensuring that no irrelevant variables are included. Furthermore, a high threshold on linear gain will avoid redundancies in the set. For instance, in the example from traffic data described in previous section, suppose we insert  $T_4$  into the set  $\{T_1, T_2, T_3\}$ , where  $T_4$  comes from a sensor close to the sensor for  $T_3$ , on the the same road. Since  $T_4$  is almost a duplicate of  $T_3$ , then the linear dependence of the resultant set  $S' = \{T_1, T_2, T_3, T_4\}$  would be almost 1, since there would exist a linear combination  $T_4 - T_3 \approx 0$  with near perfect linear dependence. However, that would also imply that many of its subsets, e.g.  $\{T_1, T_3, T_4\}$ , would have near perfect linear dependence. Hence by definition, the linear gain of S' will be very close to zero. More generally, a high threshold on linear gain will also avoid multicollinearity in the set. For instance, consider a set  $S = S_1 \cup S_2$  that consists of two independent subsets  $S_1$  and  $S_2$  of perfectly linearly dependent variables. By definition, the linear gain of such a set *S* will be 0, and hence will be discarded.

Using the above definitions, we next present the formal definition of a multipole.

**Definition 5.** A **multipole** refers to the set *S* of variables with  $|S| \ge 2$  such that  $\sigma_S \ge \sigma$  and  $\Delta \sigma_S \ge \delta$ , where  $\sigma$  and  $\delta$  are user-specified thresholds.

IEEE TRANSCATIONS ON KNOWLEDGE AND DATA ENGINEERING

We next define the notion of maximality in a multipole.

**Definition 6** (Maximal Multipole). In a set Q of multipoles, a multipole S is considered to be maximal if none of its supersets are in Q.

A maximal multipole is likely to capture the underlying signal more comprehensively compared to its subsets. Hence, all non-maximal multipoles could potentially be pruned in the final output of the search.

Using the above definitions, we formulate the multipolediscovery problem as the following:

**Definition 7 (Problem Formulation:).** Given  $\delta$  and  $\sigma$ , find the set *P* of all maximal multipoles in a given time series dataset.

## **3 RELATED WORK**

In this section, we present an overview of different techniques that have been applied to study relationships between two or more variables and discuss their similarities and differences with multipoles.

Eigenanalysis-based Approaches: Eigenanalysis basedapproaches such as Principal Component Analysis (PCA) [9], [10] and Independent Component Analysis [17] commonly focus on finding linear combinations of variables that capture the dominant global signals of variability in the data. Thus, they are interested in the largest eigenvalues and often treat the smaller eigenvalues as noise and discard them. If we consider a dataset of the traffic time series at all roads in the city and apply PCA on it, it is expected to capture dominant global patterns of variability such as work-home traffic patterns, patterns influenced by social events, etc. In contrast, multipoles capture the least-variant local linear dependencies among small subsets of vectors. (Roads in this example.) This corresponds to subsets of vectors whose least eigenvalue (of the correlation matrix) is small, i.e., close to 0.

**Regression Models:** Regression models, such as ordinary least squares (OLS) and its regularized variants , e.g. LASSO [5], are used to find a linear combination of independent variables that predicts the given dependent variable with high accuracy [6], [7]. The independent variables can therefore be considered as showing a strong linear dependence with the dependent variable. However, such techniques do not have a notion of gain and thus are not designed to find multipole relationships. To illustrate, consider once again the highway time series example. If we use LASSO to find a set of predictors for the highway time series  $T_1$ , LASSO would always include the time series that is most strongly correlated with  $T_1$  as the predictor, say  $T'_1$ , collected at another sensor on the same highway. Consequently, it will miss all multipole patterns that include  $T_1$  but not  $T'_1$ , many of which otherwise could be capturing interesting and non-trivial relationships of  $T_1$  with distant road stations.

**Error-In-Variables Models:** Error-in-Variables (EIV) models are a special class of regression models that account for uncertainties in the measurements of both dependent and independent variables (unlike standard regression

models, which assume that independent variables are measured accurately). Like EIV models, the definition of multipoles does not create any distinctions among the participating variables. One of the multivariable linear EIV models, named Total Least Squares (TLS) has striking similarities with the proposed definition of multipoles [8]. In particular, TLS focuses on learning a linear combination of a given set of dependent and independent variables to minimize the joint residual error in all the variables, which is exactly same as finding a linear combination of variables with highest linear dependence. Like multipoles, the solution to TLS is obtained by computing the eigenvector corresponding to the least eigenvalue of the covariance matrix of the given set of variables. TLS does not have a notion of linear gain, although it can be shown that a given set of variables obtains a high linear gain only if 1) TLS obtains a unique solution, and 2) all the regression coefficients obtained in the solution of TLS are significantly higher than zero in magnitude. Thus, applying TLS on a given set of variables could be an alternative approch to evaluate the goodness of a multipole relationship formed between the variables of the set. However, TLS does not provide any approach for searching through a large set of variables, e.g., the time series that capture temperature on the Earth's surface. Thus,

variables forming multipole patterns from a larger dataset. Structure Learning: Another stream of related work in machine learning literature is that of structure learning methods that learn the structure of stochastic dependencies among variables in a dataset in the form of a graphical model called Markov network [11], [12], which is a graph where each node represents a variable and follows the pairwise Markovian property, according to which it is independent of any non-neighboring node in the network conditioned on all of its neighboring nodes. Markov networks are typically studied to infer the conditional independence between different subsets of variables using various statistical inference techniques [18]. In contrast the multipole patterns are defined to capture direct or indirect dependencies between different subsets of variables. An experimental demonstration on the limitations of structure learning methods in finding multipoles is provided in supplemental material.

it cannot not be used as a tool to find all the subsets of

**Correlation Networks:** Linear relationships in time series data have also been studied in past using correlation networks, where each time series represent a node, and the weight of an edge between any two nodes represents the strength of the linear correlation between the corresponding time series. Correlation networks have been used in past for studying a variety of patterns, the most popular being 'community', which refers to a group of nodes (time series) with strong mutual positive correlations [13], [14]. If considered as a potential multipole, a community would have very low linear gain since its time series are highly similar, i.e., they show considerable collinearity. In contrast, time series in a multipole with high linear gain cannot be highly similar.

Some works, including our own [15], have further studied pairs of negatively correlated communities, which form *dipoles*. Multipoles often have negative correlations among vectors—see Section 4.2. Nonetheless, those links can be relatively weak, i.e., not meaningful dipoles. Further, there is no guarantee that a dipole will show up as part of a multipole pattern. We also recently defined *tripoles* [16], but a multipole is not a generalization of a tripole. A tripole consists of a root and a pair of leaf time series, such that the sum of the leaf time series shows much stronger correlation with the root compared to either of their individual correlations with the root. Thus, as with regression, one of the variables (root or dependent variable) has a special role, which is not the case for multipoles. Further, tripoles are restricted to only one linear combination (i.e. sum of leaves), whereas multipoles allow arbitrary normalized linear combination to attain linear dependence. More importantly, there does not seem to be a way to generalize the tripole concept beyond three time series in a way that would facilitate efficient search for such patterns in a large data sets.

In summary, the problem of finding multipoles in the data is a novel and unique problem and to the best of our knowledge, there doesn't exist any method in the relevant literature that is directly suitable to solve this problem.

## 4 FINDING MULTIPOLES

The combinatorial aspect of the problem makes it extremely challenging to come up with an approach that is both computationally efficient and guarantees completeness of the search. A brute-force approach would examine all subsets of size k, varying k from 3 to N, the total number of time series in the dataset. While such an approach guarantees completeness of the search, it will easily become computationally infeasible even for very small datasets due to its exponential time complexity. To give an estimate, one of our real-world datasets is quite small in size and has only 171 time series. However, performing a brute-force search on a regular desktop over all subsets of i) size 4 takes about 4 hours, ii) size 5 takes more than 5 days, and so on for subsets of size beyond 5.

In this paper, we propose a correlation graph-based approach to capture most interesting multipoles in a computationally efficient manner. Our approach is primarily motivated by some empirical observations that indicate a direct relationship between the linear gain of a set and the strengths of pairwise correlations between the members of the set. Leveraging these observations, our approach identifies and restricts the search for multipoles to a family of subsets, which we refer to as 'promising candidates' that are more likely to exhibit multipole relationships with stronger linear gain between their members. Such a family of subsets are usually much rarer in the data, thereby contributing to the remarkably high computational efficiency of our approach, although with some loss of completeness in the final output of multipoles, especially at lower thresholds on linear gain.

In the remainder of this section, we first define a canonical form for a set of variables called *self-canceling* form that has linear dependence and linear gain identical to that of the original set. We then present our empirical observations that indicate an inverse relationship between linear gain of a set and the largest pairwise correlation in its canonical form. Finally, we describe how the empirical observations can be leveraged to identify promising candidates for multipoles and describe our proposed approach in detail.

#### 4.1 Self-Cancellation

*Definition 8.* A set is said to be **self-canceling** if all the weights in its LVNLC are non-negative.

Any non self-canceling set *S* can be converted into a selfcanceling set  $\hat{S}$ , by flipping the signs of all the members in S that have negative weights in LVNLC of S. For example, consider the set  $S \equiv T_1, T_2, T_3$  of three traffic time series from transportation example discussed in Section 1. The LVNLC of S was observed to be  $0.6T_1 + 0.65T_2 - 0.47T_3$ (see top panel of Figure 2). By flipping the sign of  $T_3$ , we get a new set  $\widetilde{S} \equiv \{T_1, T_2, -T_3\}$ , whose LVNLC is given by  $0.6T_1 + 0.65T_2 + 0.47(-T_3)$ . Since all the variables in  $\tilde{S}$  have non-negative weights in its LVNLC,  $\tilde{S}$  is a self-canceling set. It is important to note that flipping the signs of one or more variables does not affect the eigenvalues of their correlation matrix. Therefore, S and  $\overline{S}$  will have identical linear dependence and linear gain. Based on the above ideas, we can define a canonical form of the set called the selfcanceling form as the following.

**Definition 9.** A self-canceling form of a set S is a canonical form  $\tilde{S}$  that is obtained by flipping the signs of all the variables that have negative weights in LVNLC of S.

#### 4.2 Empirical Observations

We are now in a position to describe our empirical observations that relate linear gain of a set S with largest pairwise correlation in its self-canonical form  $\tilde{S}$ . Specifically, let  $\rho_S$  denote the highest pairwise correlation observed in  $\tilde{S}$ . Figure 3 then shows scatter plots between  $\Delta \sigma_S$  (Xaxis) and  $\rho_S$  (Y-axis) for more than a million correlation matrices of sizes  $k \times k$  for k = 3, 4, 5. Each of these matrices were generated by sampling all  $\binom{k}{2}$  pairwise correlations between [-1, 1] using a uniform distribution. Among the generated matrices, only the ones that satisfy positive semidefiniteness were considered to be valid correlation matrices. The implementation of the procedure can be accessed at this URL <sup>1</sup>. From Figure 3, we make two key observations:

First key observation: All the plots in Figure 3 show that the sets with largest linear gain (the rightmost end of the distribution) always lie below horizontal green line that corresponds to equation  $\rho_S = 0$ . This implies that the linear gain of a set S tends to be higher when the largest pairwise correlation in  $\tilde{S}$  is strongly negative. More generally, for all sets that form multipoles with linear gain at least  $\delta$ , there exists an upper bound on the largest pairwise correlation in their self-canceling forms.

Second key observation: The maximum possible linear gain of a multipole of size k is  $\frac{1}{k-1}$ . (In Figure 3, we can see that for k equal to 3,4, and 5, the maximum linear gain is 0.5, 0.33, and 0.25, respectively.) The maximum linear gain corresponds to the set where all the pairwise correlations in its self-canceling version equals  $\frac{-1}{k-1}$ . This implies that the linear gain is smaller for larger multipoles. Hence, if we are only interested in finding multipoles with linear gain at least  $\delta$ , we can safely ignore all sets beyond of size  $\lfloor \frac{1+\delta}{\delta} \rfloor$ . Both of the above two observations were found to empirically hold true also for sets of sizes beyond 5.



Fig. 3: Empirical relationship between linear gain  $\Delta \sigma_S$  of a set S (X-axis) and largest pairwise correlation in its self-canceling form  $\tilde{S}$  (Y-axis) on more than  $10^6$  generated  $k \times k$  correlation matrices for  $k \in [3, 5]$ . Vertical red-dash line in each plot indicates largest value of linear gain observed. All promising candidates for CoMEt approach lie below the  $\rho_S = 0$  (green solid line). All multipoles with linear gain larger than a given threshold  $\delta$  lie to the right of the vertical solid black line.



Equivalent clique

Fig. 4: Illustrating equivalence between a negativeequivalent clique and a negative clique.

Based on the above empirical observations, we next define a promising candidate for a multipole in the following subsection.

### 4.3 Promising Candidates

**Definition 10.** A set *S* is said to be a **promising candidate** if i)  $|S| \leq \lfloor \frac{1+\delta}{\delta} \rfloor$ , where  $\delta$  is the user-specified threshold on linear gain, and ii) maximum pairwise correlation in its self-canceling form  $\tilde{S}$  is negative.

Promising candidates can be classified into two types: i) Negative Cliques, and ii) Negative Equivalent Cliques.

*Definition 11.* A **negative clique** refers to a set where all the pairwise correlations between its members are negative.

The terminology is motivated from the appearance of such sets in a correlation graph, where each vertex represents a variable and the weight of an edge  $e(X_i, X_j)$  is equal to  $corr(X_i, X_j)$ . Such sets would appear as a clique of negative edges in the correlation graph. It can be shown that the self-canceling form of any negative clique is itself, i.e.  $S = \tilde{S}$ , and by definition satisfies the requirement of all pairwise correlations in  $\tilde{S}$  to be negative. Therefore, a negative clique is a promising candidate.

## **Definition 12.** A negative-equivalent clique refers to a set S whose self-canceling form $\tilde{S}$ is a negative clique.

All negative-equivalent cliques can be identified using the following lemma.

*Lemma* 3. A set S is a **negative-equivalent clique** *iff* it can be partitioned into two negative cliques  $S_1$  and  $S_2$  such that the all the cross correlations between members of  $S_1$ and  $S_2$  are non-negative.

**Proof**: See supplemental material.

## 4.4 Proposed Approach: CoMEt

Leveraging the above empirical observations and the concept of promising candidates discussed in previous section, we propose our Clique Based Multipole SEarch (CoMEt) to find multipoles. The central idea of CoMEt is to find all promising candidates for multipoles (negative cliques and negative-equivalent cliques) and then check each of them to obtain true multipoles. To find all promising candidates, we first construct a graph such that every promising candidate forms a clique in it. We then obtain all maximal promising candidates by enumerating all the maximal cliques of the constructed graph. Note that maximal clique-enumeration problem is NP-complete in general. However, the cliques of our interest tend to be rare in the graphs generated from real-world datasets, which allows us to recover most of the promising candidates in much less computing time. For each maximal clique obtained, we examine its subcliques and select all those that form true multipoles. Finally we eliminate all duplicate and non-maximal multipoles to obtain the final set of maximal multipoles.

Algorithm 1 summarizes the CoMEt approach. We begin by finding all maximal promising candidates in line 3, viz. maximal negative cliques and maximal negative-equivalent cliques using Algorithm 2. Among the obtained maximal promising candidates, we then obtain all the multipoles in line 3 using Algorithm 3. Finally, in line 5, we eliminate duplicate and non-maximal multipoles using Algorithm 4. We next describe each of the modules used in the different steps of CoMEt.

## 4.4.1 FIND MAXIMAL PROMISING CANDIDATES

This module is used to find all maximal promising candidates in the data that include all negative and negativeequivalent cliques. Algorithm 2 summarizes this module. The key idea is to construct a correlation graph (network) G

#### IEEE TRANSCATIONS ON KNOWLEDGE AND DATA ENGINEERING

#### Algorithm 1 CoMEt (Clique Based Multipole Search)

(lines 1-4) such that a set of nodes of size  $\geq 3$  in G would form a clique iff it is a negative clique or negative-equivalent clique. To accomplish that, we first construct two identical graphs  $G_1$  and  $G_2$ , such that for both graphs, a set of nodes will form a clique iff it is a negative clique. To obtain such graphs, we begin by creating a set  $V_1 = \{v_{11}, v_{12}, ..., v_{1n}\}$  of n nodes such that each node  $v_{1i}$  corresponds to a time series  $X_i$  in the given dataset D. For any pair of nodes  $(v_i, v_j)$ , an edge is drawn iff  $corr(X_i, X_j) \leq 0$ . Let  $E_1$  denote the set of all such edges. Then  $G_1 = (V_1, E_1)$  is the desired correlation graph where a clique will be formed among a set of nodes iff their corresponding variables in D form a negative clique. Similarly, an identical correlation graph  $G_2 = (V_2, E_2)$  can be constructed on a set of nodes  $V_2 = \{v_{21}, v_{22}, ..., v_{2n}\},\$ where each node  $v_{2i}$  corresponds to a time series  $X_i$  in the given dataset D.

Next, to include edges for non-negative correlations in all negative-equivalent cliques, we construct a set E of crossedges between nodes of  $G_1$  and  $G_2$ . Specifically, we connect each node  $v_{1i} \in V_1$  to all the nodes  $v_{2j} \in V_2$  such that  $corr(X_i, X_j) \geq 0$ . As a result of this operation, for any negative-equivalent clique  $S = S_1 \cup S_2$ , where  $S_1$  and  $S_2$  are its two negative subcliques, all the non-negative correlations across  $S_1$  and  $S_2$  are now included. Hence, in the resultant graph  $G = (V_1 \cup V_2, E_1 \cup E_2 \cup E)$ , every negative-equivalent clique will also appear as a clique. Also note that every clique S of size  $\geq 3$  in G would be a promising candidate; it would either be a negative clique (if  $S \subset G_1$  or  $S \subset G_2$ ), or a negative-equivalent clique (if it includes nodes from both  $G_1$  or  $G_2$ ).

From the resultant graph  $G_{i}$  all the maximal promising candidates could be obtained by enumerating all maximal cliques in graph G using any of the standard cliqueenumeration algorithms. In this work, we used an efficient algorithm proposed in [19] (line 11) to enumerate maximal cliques in sparse graphs (implementation provided by authors in [20]).

Note that every maximal promising candidate will result in formation of two maximal cliques in G. For instance, every negative clique will form two cliques: one in  $G_1$  and  $G_2$  each. Similarly, a negative-equivalent clique  $S = (S_1 = \{X_i, X_i\}, S_2 = \{X_k\})$ , where  $S_1$  and  $S_2$  are two negative sub-cliques, will result in formation of two cliques:  $(v_{1i}, v_{1j}, v_{2k})$  and  $(v_{2i}, v_{2j}, v_{1k})$  in G. For every maximal promising candidate, exactly one of the maximal cliques is retained (line 12). Finally, the set of retained maximal cliques is returned as the set of all maximal promising candidates.

## 4.4.2 GET MULTIPOLES

This procedure is applied to each of the obtained maximal promising candidates to extract all multipole relationships.

#### Algorithm 2 FIND MAXIMAL PROMISING CANDIDATES

 $\triangleright \rho = 0$  for CoMEt **Input** Dataset: $\mathcal{D}, \rho$ **Output** set *C* of all maximal promising candidates **Correlation Graph Construction:** 1:  $V_1 \leftarrow$  set of *n* vertices  $v_{11}, v_{12}, ..., v_{1n}$  corresponding to *n* variables

- $X_1, X_2, ..., X_n$  in D
- 2:  $E_1 \leftarrow \text{All pairs } (v_{1i}, v_{1j}) \text{ s.t. } corr(X_i, X_j) \leq 0$ 3:  $G_1 = (V_1, E_1)$
- 4:  $G_2 = (V_2, E_2)$  be the exact duplicate of  $G_1$ 5:  $E \leftarrow \phi$  $\triangleright$  An empty set of edges
- 6: for each  $v_{1i} \in V_1$  do
- 7:  $E' \leftarrow \text{all pairs } (v_{1i}, v_{2j}) \text{ s.t. } corr(X_i, X_j) \ge 0$
- $E = E \cup E'$ 8:
- 9: end for
- 10:  $G = (V_1 \cup V_2, E_1 \cup E_2 \cup E)$
- 11:  $C \leftarrow All \text{ maximal cliques of } G$ ▷ algorithm proposed in [19]
- 12: Remove all duplicate maximal cliques from C
- 13: return C

As summarized in Algorithm 3, we begin by checking if the given candidate S forms a multipole by comparing its linear dependence and linear gain with user-specified thresholds  $\sigma$  and  $\delta$ , respectively. If so, then S is added to the set of discovered multipoles and we move on to the next promising candidate. Otherwise, it could be possible that one or more of the subsets of S might form a multipole that satisfies the thresholds. However, if  $\sigma_S$  turns out to be lower than the threshold  $\sigma$ , then by Lemma 2, all the subsets of S would also have weaker linear dependence than  $\sigma$  and thus, could be safely ignored. Therefore, only if  $\sigma_S \geq \sigma$ , do we perform an exhaustive search on all subsets of Sof sizes  $[3, \lfloor \frac{1+\delta}{\delta} \rfloor]$  and select all the ones that satisfy the thresholds. The range of sizes of subsets is derived based on the observation made in section 4.2, which states that for a given threshold  $\delta$  on linear gain, all sets of sizes beyond  $\left|\frac{1+\delta}{\delta}\right|$  can be safely discarded.

Applying the above procedure to all maximal promising candidates might result in inclusion of several non-maximal multipoles in the output. Furthermore, a multipole could be generated multiple times from different maximal promising candidates. Hence, in the final step of CoMEt, we eliminate all non-maximal and duplicate multipoles using the module described below.

Algorithm	3 GET	MULTIPOLES	5 FROM	CANDIDATE(S)	
					1

Input: set  $S, \sigma, \delta$ **Output** All multipoles S' with  $\sigma_{S'} \geq \sigma$ , and  $\Delta \sigma_{S'} \geq \delta$ 1: if  $\Delta \sigma_S \geq \delta$  and  $\sigma_{S'} \geq \sigma$  then 2:  $M \leftarrow S$ 3: else if  $\sigma_{S'} \geq \sigma$  then  $M \leftarrow all$  subsets S' of S s.t.  $|S'| \in [3, |1 + \frac{1}{\delta}|], \sigma_{S'} \geq \sigma$ , and 4:  $\Delta \sigma_{S'} \geq \delta$ 5: end if 6: return M

#### 4.4.3 REMOVE NON-MAXIMALS & DUPLICATES

This module is called in line 5 of Algorithm 1 to eliminate all non-maximal and duplicate multipoles that are generated in the previous step. As described in Algorithm 4, we first initialize two empty sets (lines 1-2): i) U' that collects the final set of non-redundant multipoles, and ii) IsIncluded, that maintains a collection of all the subsets of multipoles that are included in the final set U' at any point of time. We then scan through all the multipoles in input set U. For



Fig. 5: **Motivation of CoMEtExtended:** (a) shows the multipoles missed by CoMEt for a given threshold  $\delta$  on linear gain, which could be recovered by CoMEtExtended by setting parameter  $\rho$  to a higher value as shown in (b).

each multipole S in U, we first check if it is present in the IsIncluded set, and if not, insert it into the final output set U' as well as insert all the subsets of S in IsIncluded that are not previously present in it (lines 5-8). We then remove S from U and repeat the entire procedure in lines (4-9) until U is empty. All the multipoles in the resultant set U' are consequently distinct and maximal.

Algorithm 4 REMOVE DUPLICATES & NON-MAXIMALS

	<b>Input:</b> a set $U$ of multipoles <b>Output:</b> a set $U'$ of non-maximal and distinct multipoles
1.	
1.	$C \leftarrow \varphi$
2:	$IsIncluded \leftarrow \phi$
3:	while $U \neq \phi$ do
4:	$S \leftarrow A$ multipole in U
5:	if S not in IsIncluded then
6:	Insert S' to $IsIncluded \ \forall S' \subseteq S$
7:	$U' \leftarrow U' \cup S$
8:	end if
9:	$U \leftarrow U - S$
10:	end while
11:	return $U'$

#### 4.5 CoMEtExtended

In certain cases, the completeness and computational efficiency of CoMEt might be unsatisfactory. For instance, at lower thresholds of linear gain, the obtained set of promising candidates by CoMEt approach could potentially miss some of the interesting multipoles (see Figure 5(a)). On the other hand, at high thresholds of linear gain, there could potentially be many false positives among promising candidates, which would compromise the computational efficiency. Hence, to overcome these limitations of CoMEt, we further propose CoMEtExtended, a generalized version of CoMEt, where we redefine what constitutes a promising candidate and allow user to expand or prune the scope of promising candidates in different scenarios. The ability to adjust the search space helps CoMEtExtended in achieving a better trade-off between computational efficiency and completeness at different thresholds of linear gain (demonstrated further in Section 5.3).

We first begin with redefining the notion of a promising candidate. Specifically, instead of enforcing that all the pairwise correlations in  $\tilde{S}$  be negative, we require all of them to be below a threshold  $\rho$ , where  $\rho \in [-1, 1]$ . Likewise, we redefine the notion of negative and negative equiva-

lent cliques as *pseudo-negative* and *pseudo negative equivalent* cliques respectively as the following:

- **Definition 13.** For a given  $\rho \in [-1, 1]$ , a **pseudo negative** clique refers to a set where all the pairwise correlations between its members are less than or equal to  $\rho$ .
- **Definition 14.** A pseudo negative-equivalent clique refers to a set S whose self-canceling form  $\tilde{S}$  is a pseudonegative clique.

A pseudo negative-equivalent clique can be identified using the following lemma which is a general version of Lemma 3.

*Lemma 4.* A set *S* is a **pseudo negative-equivalent clique** *iff* it can be partitioned into two pseudo negative cliques  $S_1$  and  $S_2$  such that the all the cross correlations between members of  $S_1$  and  $S_2$  are pseudo non-negative.

Proof: See the Supplemental material.

Using the above definitions and results, we now propose CoMEtExtended which is the more general version of CoMEt. The steps of CoMEtExtended are exactly the same as that of CoMEt, the only difference being in the lines 2 and 7 of Algorithm 2 where the threshold of 0 on pairwise correlations is replaced by  $\rho$  and  $-\rho$ , respectively. By incorporating  $\rho$  in thresholds, we ensure that every clique in the resultant graph is either a pseudo negative clique or a pseudo negative-equivalent clique. Note that setting  $\rho$  to zero gives us the original definition of promising candidates, whereas setting  $\rho$  to a higher positive value would include more sets as promising candidates and hence lead to recovery of more of the missed multipoles (see Figure 5(b)). However, it should be noted that setting  $\rho$  to positive values would increase the number of candidate cliques and thus increase the computational cost. In fact, setting  $\rho$  to 1 would turn CoMEtExtended into an exhaustive bruteforce search, which guarantees completeness, but would be computationally infeasible. On the other hand, setting  $\rho$  to a high negative value would reduce the search space and hence leads to much faster recovery of the multipoles with a very high linear gain in larger datasets.

## 5 DATA AND EXPERIMENTAL EVALUATION

In this section, we discuss results and computational evaluation of proposed approach CoMEtExtended. Specifically, we evaluate the completeness of CoMEtExtended against a regularized linear regression-based baseline, analyze tradeoff between completeness and computational efficiency of CoMEtExtended at multiple parameter settings, study the scalability of CoMEtExtended, analyze the statistical significance of the multipoles, and evaluate their utility using real-world datasets from climate science and neuroscience domains. All experiments were run on a computer with 20 processors, each processor being Intel(R) Xeon(R) CPU E5-2470 0 running at 2.30GHz with a total shared RAM of 100 GB, running Linux version 2.6.32-696. We begin by describing all the datasets along with the data pre-processing steps that were applied to each of them.

#### 5.1 Data and Preprocessing

#### 5.1.1 Sea Level Pressure (SLP) data:

We used monthly Sea Level Pressure (SLP) dataset provided by NCEP/National Center for Atmospheric Research (NCAR) Reanalysis Project [21], which is available from 1979-2014 (36 years) at a spatial resolution of  $2.5 \times 2.5$ degree (10512 grid points, also referred to as locations). In this paper, we constructed SLP time series for each location using only the months of winter season (December, January, and February) from each year, thereby resulting in 108 observations in every time series. For each of the time series, we followed the standard pre-processing steps followed in climate science to remove the annual seasonality and linear trends [1].

Relationships in climate datasets are preferably studied between regions (sets of spatially contiguous locations) as opposed to individual locations because of spatial autocorrelation, due to which locations in a spatial neighborhood have highly similar time series that will lead to discovery of redundant relationships. Therefore, we next converted the given location-based time series dataset into a set of 171 region-based time series dataset using a simple clustering procedure that is described in supplemental material. In addition, for purposes of validation of obtained multipoles, we used monthly Hadley Center SLP (HadSLP2) observational data available for years prior to 1979 to obtain the time series of these regions.

#### 5.1.2 Brain fMRI data:

We used neuroimaging data collected at the University of Utah as part of a reproducibility study [22]. In this study, a set of 50 functional-Magnetic Resonance Imaging (fMRI) scans of one subject were acquired while the subject was involved in an audio-visual task (watching cartoons). Another set of 50 fMRI scans were collected from the same subject while the subject was resting. The spatial resolution and the temporal resolution of every scan was  $3mm \times 3mm \times 3mm$ and 2 secs, respectively. A number of fMRI pre-processing steps-described in [22]-were performed including motion correction, unwarping, and filtering. In addition, we used an Automated Anatomical Labeling Atlas [23], which maps grey matter locations to 90 anatomical regions, to compute a mean time series of each brain region from each scan. As a result, we obtained a set of 90 time series for each of the 100 fMRI scans. We applied our approach to find multipoles in one of the 50 audio-visual fMRI scans, while the other 49 scans were used for evaluation purposes that we will describe later in this section.

## 5.2 Parameter settings

Two user-specified thresholds, minimum linear gain ( $\delta$ ) and minimum linear dependence ( $\sigma$ ), are needed for discovering multipoles. The choice of values for these parameters needs to be determined based on domain knowledge and the availability of computational resources to find multipoles. In particular, a relaxed linear gain threshold  $\delta$  will increase the search space of multipoles and thus will require more computational time for search and evaluation. Similarly, a lower threshold of  $\sigma$  will result in a larger number of In this work, we performed computational evaluation and scalability analysis at different combinations of values of  $\sigma \in \{0.4, 0.5, 0.6\}$  and  $\delta \in \{0.1, 0.15, 0.2\}$  respectively. Statistical significance analysis was performed on multipoles obtained at  $\sigma = 0.50$  and  $\delta = 0.15$  for both SLP and brain fMRI datasets.

	Total Multipoles in	Completeness		
$(\sigma, \delta)$	Pseudo-complete set	LAB	CoMEtExtended	
(0.4,0.1)	70150	0.09%	81%, <i>ρ</i> =0.01	
(0.4,0.15)	6255	0.33%	96%, <i>ρ</i> =0.01	
(0.4,0.2)	1264	0.39 %	99%, <i>ρ</i> =0.01	
(0.5,0.1)	41126	0.15%	76%, <i>ρ</i> =0.01	
(0.5,0.15)	3348	0.62%	92%, <i>ρ</i> =0.01	
(0.5,0.2)	930	0.54%	99%, <i>ρ</i> =0.01	
(0.6,0.1)	13743	0.47%	75.5%, <i>ρ</i> =0.03	
(0.6,0.15)	1525	1.38%	85%, <i>ρ</i> =0.01	
(0.6,0.2)	488	1.0%	98%, <i>ρ</i> =0.01	

TABLE 1: Completeness evaluation of CoMEtExtended against LASSO-based baseline (LAB) at different combinations of  $\sigma$  and  $\delta$  in the SLP dataset. The parameter  $\rho$  in CoMEtExtended was set so as to keep the computational time under 90 minutes.

	Total Multipoles in		Completeness	
(σ,δ)	Pseudo-complete set	LAB	CoMEtExtended	
(0.4,0.1)	15855	0.006%	71%, <i>ρ</i> =0.2	
(0.4,0.15)	3019	0%	98%, <i>ρ</i> =0.2	
(0.4,0.2)	716	0%	100%, <i>ρ</i> =0.2	
(0.5,0.1)	15258	0.006%	70%, <i>ρ</i> =0.2	
(0.5,0.15)	2805	0%	98%, ρ=0.2	
(0.5,0.2)	697	0%	100%, <i>ρ</i> =0.2	
(0.6,0.1)	13721	0.007%	82%, <i>ρ</i> =0.25	
(0.6,0.15)	2172	0%	97%, <i>ρ</i> =0.2	
(0.6,0.2)	547	0%	100%, <i>ρ</i> =0.2	

TABLE 2: Same as Table 1, but in the fMRI dataset.

## 5.3 Evaluation

Comparison of CoMEtExtended to another approach is difficult, since it defines local patterns in terms of linear algebra concepts and finds such patterns by searching for negative cliques in a similarity graph. None of the related works mentioned in Section 3 do the same thing. Nonetheless, to provide some comparison, we evaluate the completeness of the search of CoMEtExtended with respect to a LASSObased baseline approach (LAB) that is described as follows:

**LASSO-based baseline approach (LAB):** LASSO is a variant of regularized linear regression that obtains a subset of variables from a larger set that could be linearly combined to predict the given predictand with high accuracy. Specifically, given a predictand *Y* and a set of predictors  $\mathbf{X} = [X_1, X_2, ..., X_n]$ , it learns a sparse set of regression coefficients  $\beta = [\beta_1, \beta_2, ..., \beta_k]^T$  by minimizing the following objective function:

$$\min_{\beta} ||Y - \mathbf{X}\beta||_2 + \lambda ||\beta||_1$$

where  $\lambda$  is the hyperparameter that can be tuned to control the number of non-zero regression coefficients in the solution. In this baseline, we use LASSO to find potential candidates that could form multipoles. Specifically, in a dataset of *N* variables  $\{X_1, X_2, ..., X_n\}$ , for any variable  $X_i$ , we consider the remaining set of variables as predictors and apply LASSO find a subset of predictors from the remaining N-1 variables that could be linearly combined to best predict  $X_i$ . The combined set of  $X_i$  and the predictors with non-zero regression coefficients is then added to the set of potential candidates of multipoles. For each variable  $X_i$ , we obtain up to N-1 different solutions of LASSO by varying  $\lambda$  such that every solution yields a different number of predictors with non-zero regression coefficients. By applying the above procedure to all possible  $X_i$ , we obtain up o  $N \times (N-1)$  potential candidates of multipoles in total. For each the obtained candidates, we next compute linear dependence and linear gain and discard all those that do not satisfy the given thresholds  $\sigma$  and  $\delta$ .

We next evaluate the completeness of CoMEtExtended with that of LAB. The completeness of an algorithm is measured as the fraction of multipoles of a *complete set* (a set that includes all the multipoles present in the data) that it finds from the data. An ideal approach to generate a complete set of multipoles will be to perform an exhaustive brute-force search over all possible subsets of variables and select those subsets that show linear dependence  $\geq \sigma$  and linear gain  $\geq \delta$ . However, in practice, such an approach is computationally infeasible even for small-size datasets. For instance, in case of SLP data that has 171 time series, the estimated time on our computer to examine all subsets of size 5 is more than 5 days, size 6 is more than 144 days, and so on. Hence, for the purposes of evaluation, we generated a pseudo-complete set of multipoles that includes unique and non-redundant multipoles generated by the following approaches: i) an exhaustive brute-force search over all subsets of variables upto size 4 for SLP, and 5 for fMRI dataset, ii) A random-search approach that is run for 24 hours to examine random subsets of variables of size k, where  $k \ge 5$  for SLP, and  $k \ge 6$  for fMRI datasets, and select the ones that show linear dependence  $\geq \sigma$  and linear gain  $\geq \delta$ , iii) LAB approach, and iv) CoMEtExtended approach, where the choice of parameters is set so as to keep its computational time within 90 minutes.

Tables 1 and 2 summarize the results of completeness evaluation on SLP and fMRI datasets respectively. The first column in both the tables indicate the total size of the pseudo-complete set of multipoles at different combinations of parameters  $\sigma$  and  $\delta$ . Second and third columns indicate the completeness of LAB and CoMEtExtended approach. As can be seen in both tables, LAB recovers a negligible or very tiny fraction of multipoles in all the parameter settings. For instance, in the SLP dataset, LAB recovers less than 1.3% of total multipoles in pseudo-complete set in all the parameter settings. The completeness is negligible (< 0.1%) for cases where  $\delta = 0.1$ . Further, in the fMRI dataset, LAB is able to recover only one multipole for cases where  $\delta = 0.1$ , whereas for the remaining other cases where  $\delta \ge 0.15$ , LAB is unable to find any multipole. This strongly indicates that linear regression-based approaches like LASSO are not suitable for finding multipoles in the data.

In contrast, the completeness of CoMEtExtended is more than 80 % for all the parameter settings. Moreover, the completeness of CoMEtExtended increases at higher thresholds of linear gain and is close to 100 % for  $\delta = 0.2$ , which shows that our approach is less likely to miss multipoles with high linear gain. This is in concordance with our empirical observations made in Figure 3, according to which, whenever linear gain is high, the largest pairwise correlation in the self-canceling version of all of the multipoles tends to be much lower and approaches strong negative values. As a result, they are more likely to be captured as promising candidates by our approach. In summary, CoMEtExtended outperforms the LASSO-based baseline approach in completeness at different thresholds on linear dependence and linear gain. Further, it is much more efficient and relatively complete in finding multipoles at higher thresholds of linear gain, compared to brute-force and the LASSO-based baseline approach.

**Trade-off between completeness and efficiency:** We also evaluated the performance of CoMEtExtended based on the trade-off made between completeness and efficiency by varying parameter  $\rho$ . Table 3 shows completeness and computational time (in seconds) on SLP and fMRI datasets respectively at different values of  $\rho$  for different combinations of  $\sigma$  and  $\delta$ . As  $\rho$  increases, the completeness improves but also adds to the total computational cost. This is expected since  $\rho$  signifies the upper bound on the largest pairwise correlation in the self-canceling version of a promising candidate. Therefore, as  $\rho$  increases, more sets qualify as promising candidates, which further expands the search space, leading to higher computational cost.

	Completeness, Computing Time (in minutes)				
$(\sigma, \delta)$	$\rho = -0.2$	$\rho = -0.1$	$\rho = 0$	$\rho = 0.01$	
(0.4,0.1)	1%,0.1	12%,0.5	77%,13	81%,28	
(0.4,0.15)	21%,0.1	76%,0.3	95%,9	96%,21	
(0.4,0.2)	73%,0.1	92%,0.2	99%,9	99%,20	
(0.5,0.1)	2%,0.1	12%,0.3	71%,11	76%,23	
(0.5,0.15)	27%,0.1	63%,0.2	91%,8	72%,18	
(0.5,0.2)	70%,0.1	90%,0.2	98%,7	99%,17	
(0.6,0.1)	6%,0.1	14%,0.2	63%,7	68%,14	
(0.6,0.15)	31%,0.1	55%,0.2	82%,6	85%,12	
(0.6,0.2)	67%,0.1	83%,0.2	97%,5	98%,11	

TABLE 3: Performance of CoMEtExtended at different values of  $\rho$  for different combinations of  $\sigma$  and  $\delta$  in SLP dataset. Each cell in the table contains two values: i)Completeness of search, and ii) Computational time (in minutes) taken by CoMEtExtended.

Also note that at higher  $\delta$ , the completeness of CoMEtExtended reaches close to 100 % at much smaller values of  $\rho$  and requires much less computing time. This indicates that our approach is much more efficient in finding multipoles with high linear gain. This is again consistent with the empirical observations made in Figure 3 where we observed that for a multipole with high linear gain, all the pairwise correlations in its self-canceling version tend to have stronger negative values. Therefore, they get included among the promising candidates at much lower values of  $\rho$ .

#### 5.4 Scalability Analysis

As is common with many pattern finding techniques, such as frequent pattern mining, the CoMEt algorithm is inherently exponential. (For a detailed analysis of the the time complexity of the the three parts of the algorithm, see the Supplemental material.) However, as with association analysis, adjusting the parameter settings—in this case,  $\sigma_S$  and  $\Delta \sigma_S$ —can make the pattern search quite tractable in many situations. To illustrate, we next discuss the scalability of our approach using SLP data, i.e., how does the computational time vary with the size of datasets (number of time series). To obtain datasets of different sizes, we first generated 10 additional seasonal datasets for different seasons, each season being a set of three consecutive months: ((Jan.,Feb.,Mar.), (Feb.,Mar.,Apr.) ,...(Oct. Nov. Dec.)). Each dataset consists of time series from the same 171 regions that we chose for our original SLP dataset. We then generated 10 datasets of sizes 171\*k, where  $k \in [1, 10]$  by merging k of the above 10 seasonal datasets.

Figure 6(a) shows the total computational time of CoMEtExtended on all of the above datasets, at  $\sigma = 0.5$  and  $\delta = 0.15$  for different values of  $\rho$  in range [-0.15, -0.06]. As can be seen in figure, the computing time increases with the increase in the size of the datasets at different rates depending on  $\rho$ . For stronger negative values of  $\rho$ , the computing time increases almost linearly with the increase in size of datasets, which highlights its scalability<sup>2</sup>. However, as  $\rho$ approaches zero, the scalability is weak, and the computing time increases dramatically for bigger datasets. Similar observations are also made at other parameter settings, e.g. at  $(\sigma = 0.4, \delta = 0.15)$  and  $(\sigma = 0.4, \delta = 0.2)$ , as shown in Figures 6(b) and 6(c) respectively. The observed loss in scalability could be attributed to the typical distribution of pairwise correlations in the correlation graph of any time series dataset, as shown in Figure 6(d) for one of the SLP datasets. The distribution is bell-shaped with major fraction of edges having strengths close to zero in magnitude. Consequently, as  $\rho$  approaches zero, the number of cliques found in Step 1 of the algorithm increase exponentially. Note that many of the cliques that have all or most of the edges being weak are expected to have weak linear dependence among their variables, and hence are unlikely to form a multipole. By setting  $\rho$  to stronger negative values, we avoid such cliques and save a lot of computing time. However, that also leads to missing some of the interesting cliques where only one or two edges were close to zero.<sup>3</sup> Such cliques could potentially be recovered by heuristic approaches, which could be an interesting direction to pursue for future work.

#### 5.5 Evaluation of Multipoles

One of the key challenges of this work is distinguishing between reliable and spurious multipoles, i.e., those multipole patterns that arise due to random variation, from the large number of discovered multipoles. Domain validation is an ideal approach to evaluate multipoles, but most of the multipole relationships discovered in this work are currently unknown to domain scientists. Thus, in our work, we used an empirical evaluation framework that consists of two steps. The first step involves a procedure for estimating



Fig. 6: **Scalability analysis:** Figures 6(a), 6(b), and 6(c) plot the computing time (Y-axis) of COMETExt on different sizes of SLP datasets (X-axis) at different values of  $\rho$  for three parameter settings (indicated in subcaptions). Figure 6(d) shows the distribution of correlation strengths of edges in a correlation network in the SLP dataset that has 1710 time series. See section 5.4 for further details.

the statistical significance of a multipole. This procedure is then used in the second step to assess the reproducibility of a discovered multipole in multiple datasets, where each dataset is collected during a time period different from that of original dataset used for finding multipoles. Intuitively, spurious multipoles are less likely to reproduce in time periods that were not used for finding them. In contrast, multipoles that do reproduce with high statistical significance are more likely to be patterns that are outcome of a real phenomenon, and thus they would be ideal candidates for further investigation by domain experts.

### 5.5.1 Step 1: Statistical Significance Evaluation:

To filter spurious multipoles, it is important to answer the following two questions: i) how likely is it that the observed level of linear dependence,  $\sigma_S$ , of a multipole *S* is due to chance? and ii) does every member in *S* contribute significantly to the linear dependence of *S*?

To address the first question, we generate a null distribution of linear dependence by randomly generating 100,000 sets of time series and evaluating each set S for its level of linear dependence. Each of these randomly generated sets is created by sampling time series from different time periods. For instance, for our SLP investigations, a random set of size |S| is constructed by sampling a time series from any |S| of the nine time windows of HadSLP2 data. Similarly, for brain fMRI data, a random set of |S| time series is constructed by sampling a time series from any |S| of the 50 scans. Generating a random set in this manner is an approximation to independently generating |S| time series while ensuring that the general underlying nature of domain time series (e.g. autocorrelation, periodicity etc.) is retained in the randomly generated data. Using the resultant

<sup>2.</sup> We have also demonstrated the scalability of our approach on larger synthetic datasets containing up to 100k time series (see supplemental)

<sup>3.</sup> The exact number of missing cliques could not be computed due to the absence of ground truth and computational intractability of bruteforce approach.

IEEE TRANSCATIONS ON KNOWLEDGE AND DATA ENGINEERING



Fig. 7: **Reproducibility analysis:** Figures show the number of multipoles that are found to be reproducible (Y-axis) in different number of independent datasets (X-axis) for both SLP and fMRI data for  $\sigma = 0.5$ ,  $\delta = 0.15$ , and  $\rho = 0$ . See section 5.5.2 for further details

null distribution, we then determine the statistical significance of the multipoles we originally found. We evaluate  $\sigma_S$  at a 0.01 level of significance.

We next describe our approach to assess the second point, i.e, the significance of the contribution of each of the k variables in a given set  $S = \{X_1, X_2, ..., X_k\}$  to its linear dependence. Specifically, to assess the significance of the contribution of time series  $X_i$ , we replace it with a random time series  $X_R$  that is sampled from an independent dataset in a manner similar to the procedure above. We then compute the linear dependence of the resultant set, which we call S'. If the contribution of  $X_i$  is not spurious, it would be unlikely for a randomly chosen time series  $X_R$  to replicate it in which case,  $\sigma_{S'} \leq \sigma_S$ . We repeat the above process 1000 times and compute the fraction of the population for which  $\sigma_{S'} \leq \sigma_S$  holds true. This fraction is our significance level. We again use a significance level of 0.01. The above procedure is repeated for each of the kmembers.

#### 5.5.2 Step 2: Reproducibility in Independent Datasets

In this step, we estimate the reproducibility of a given multipole in multiple time periods. A multipole is considered to be reproducible in a dataset *D*, if, at a 0.01 level of significance, it is found to have statistically significant linear dependence, as well as a statistically significant contribution from each of its members. Specifically, for each multipole discovered in 1979-2014 SLP dataset, we computed its linear gain and linear dependence in HadSLP2 data in 9 time windows 1901-1936, 1906-1941,...,1941-1976. Likewise, for multipoles discovered in one of the 50 brain fMRI scans, we studied their reproducibility in the remaining 49 scans.

Figure 7(a) shows number of multipoles reproduced in different numbers of HadSLP2 time windows during 1901-1976. More than 40% of multipoles reproduced in all 9 time windows. Similarly, Figure 7(b) shows number of multipoles from fMRI data reproduced in different number of fMRI scans taken while the subject was watching a video. At least 25% of multipoles reproduced in more than 10 other scans. Higher reproducibility of multipoles suggests that they are more robust to noise in the data and unlikely to be spurious.

## 6 CASE STUDIES

Results discussed in the previous section indicate the existence of several multipole relationships with high reproducibility in multiple independent time periods, which makes a compelling case for their connection to underlying physical phenomena that might be currently unknown to domain scientists, but could potentially be discovered by domain experts upon further analysis. In this section, we present case studies on the physical interpretation of two of the discovered multipoles in SLP and brain fMRI data.

#### 6.1 Discovering Climate Phenomena

One of the multipoles in SLP data was found between the four regions shown in Figure 8. The time series of the four regions show negative correlations with each other, resulting in a multipole relationship with linear dependence of 0.7 and a linear gain of 0.15 during 1979-2014. Further, as indicated in Figure 9(a), the multipole was found to be reproducible in 7 out of 9 time windows during the period of 1901-1976, showing strong linear dependence (red curve) and linear gain (indicated by gap between red and black curves). This multipole appears to be strongly related to the well-known climate phenomenon known as the El-Nino Southern Oscillation (ENSO) as indicated by regions  $R_2$ and  $R_3$ . This phenomenon appears not only in the tropical Pacific Ocean, but also has large-scale impacts on regional climate outside the tropics [24]. More recently, a connection between the West Siberian Plain and ENSO was discovered as a tripole relationship (as defined in [16]) between three regions, that are re-captured as regions  $R_1$ ,  $R_2$ , and  $R_3$ in above multipole. This phenomenon was attributed to a wave train that originates from the sub-tropical Atlantic and propagates north-eastwards towards the north of the West Siberian Plain, where it is deflected southeastwards and reaches all the way to the central Pacific Ocean, where the two centers of action of ENSO are located. Notably, region  $R_4$  in the northern Atlantic Ocean is located near the proposed location of origin of the wave train. This finding can be further used to study the detailed path of the wave train and to attribute weather and climate characteristics over wave-affected regions to their potential ENSO origins.

Evaluation of Climate Models: Multipoles, being potential representatives of physical processes, could serve as useful benchmarks to evaluate various climate models that are often used to study climate change under different greenhouse gas emission scenarios. In particular, climate models can be evaluated based on their ability to reproduce the physical processes represented by these multipoles. For instance, Figure 9 compares the linear dependence and linear gain of the above multipole across multiple time windows obtained in observations data (HadSLP2) and a couple of climate models used in the IPCC (Intergovernmental Panel on Climate Change) CMIP5 (Couple Model Intercomparison Project Phase 5) evaluation. It can be seen that climate models differ in their ability to simulate the multipole effectively. Specifically, climate model MPI-ESM-MR is able to reproduce multipole with statistically significance in at least 4 time windows, whereas the other model BNU-ESM could not reproduce the given multipole at all.



Fig. 8: A 4-pole in 1979-2014 that was reproducible in 7 out of 9 time windows during 1901-1976

#### 6.2 Studying Complex Dynamics in Brain

The notion of multipoles proposed in this paper is highly suited for capturing complex signaling relationships in the brain. For example, one of the multipoles we discovered in one of the brain fMRI scans captures a relationship between three brain regions: Right Frontal Inferior  $(T_2)$ , Right Parietal Inferior( $T_1$ ), and Right Temporal Pole Superior( $T_3$ ). This multipole was found to be reproducible in 23 out of 50 'task' scans (collected while the subject was watching a cartoon video) while only in 5 of the 50 other 'resting state' scans (collected while subject was resting). This relationship is interesting for multiple reasons. First, the parietal and temporal regions are known to be the first recipients of the visual and auditory stimuli [25], respectively, that are presented to the subject in the form of cartoon videos. Second, the fact that the dependence of the above multipole relationship in task scans is greater than that of resting-state scans provides support for the argument that visual and auditory stimuli would have triggered signaling from parietal and temporal regions to the frontal region. While some existing studies [26] have observed activity in the frontal region due to visual and auditory stimuli, existence of a direct signaling pathway between visual and auditory cortices with the frontal region is yet to be fully investigated. In summary, the multipole framework could serve as a promising tool for discovering new signaling pathways that exist in the brain.

## 7 CONCLUSION AND FUTURE WORK

In this paper, we introduced and formally studied a novel class of multivariate linear relationships called *multipoles* in time series data. A multipole corresponds to a set of time series that show much stronger linear dependence compared to any of its subsets. We presented a series of empirical observations to show that most interesting multipoles could be found as cliques of negative correlations in a correlation network, and proposed a novel and computationally efficient correlation network-based approach to find multipoles in the data. We demonstrated the utility of our proposed approach to find multipoles on real-world datasets from climate and neuroscience domain. Furthermore, we presented case studies from both domains to highlight the potential of multipoles in discovering novel physical processes. While the approach proposed in this paper is based on a series of empirical observations, it is noteworthy that all of our observations are universal in nature as opposed to being

specific to a particular time series dataset. Moreover, there are certain scenarios in which our approach could be empirically shown to guarantee completeness of the search (see supplemental material for further details). Derivation of theoretical proofs of these observations are subject of future research. Other useful extensions of this work could be to extend the notion of multipoles to non-linear relationships, and generalization to time-lagged multipole-relationships.

## ACKNOWLEDGMENTS

We would like to thank Siddhant Agrawal, Department of Mathematics, University of Michigan for invaluable discussions. This work was supported by NSF grants IIS-1029771 and IIS-1319749 and NASA grant 14-CMAC14-0010. Access to the computing facilities was provided by the University of Minnesota Supercomputing Institute.

## REFERENCES

- [1] J. Kawale, S. Liess, A. Kumar, M. Steinbach, P. Snyder, V. Kumar, A. R. Ganguly, N. F. Samatova, and F. Semazzi, "A graph-based approach to find teleconnections in climate data," *Statistical Analysis and Data Mining: The ASA Data Science Journal*, vol. 6, no. 3, pp. 158–179, 2013.
- [2] A. S. Taschetto, A. S. Gupta, N. C. Jourdain, A. Santoso, C. C. Ummenhofer, and M. H. England, "Cold tongue and warm pool ENSO events in CMIP5: mean state and future projections," *Journal* of Climate, vol. 27, no. 8, pp. 2861–2885, 2014.
- [3] J. M. Wallace and D. S. Gutzler, "Teleconnections in the geopotential height field during the northern hemisphere winter," *Monthly Weather Review*, vol. 109, no. 4, pp. 784–812, 1981.
- [4] G. Atluri, A. MacDonald III, K. O. Lim, and V. Kumar, "The brainnetwork paradigm: Using functional imaging data to study how the brain works," *Computer*, vol. 49, no. 10, pp. 65–71, 2016.
  [5] R. Tibshirani, "Regression shrinkage and selection via the lasso,"
- [5] R. Tibshirani, "Regression shrinkage and selection via the lasso," *Journal of the Royal Statistical Society. Series B (Methodological)*, pp. 267–288, 1996.
- [6] A. C. Lozano *et al.*, "Spatial-temporal causal modeling for climate change attribution," in *Proceedings of the 15th ACM SIGKDD international conference on data mining*. ACM, 2009, pp. 587–596.
- [7] S. Chatterjee, K. Steinhaeuser, A. Banerjee, S. Chatterjee, and A. R. Ganguly, "Sparse group lasso: Consistency and climate applications." in SDM. SIAM, 2012, pp. 47–58.
- [8] R. J. Carroll, D. Ruppert, C. M. Crainiceanu, and L. A. Stefanski, Measurement error in nonlinear models: a modern perspective. Chapman and Hall/CRC, 2006.
- [9] A. G. Barnston and R. E. Livezey, "Classification, seasonality and persistence of low-frequency atmospheric circulation patterns," *Monthly weather review*, vol. 115, no. 6, pp. 1083–1126, 1987.
- [10] Q. Ding and B. Wang, "Circumglobal teleconnection in the northern hemisphere summer," *Journal of Climate*, vol. 18, no. 17, pp. 3483–3505, 2005.
- [11] J. Friedman, T. Hastie, and R. Tibshirani, "Sparse inverse covariance estimation with the graphical lasso," *Biostatistics*, vol. 9, no. 3, pp. 432–441, 2008.
- [12] N. Meinshausen and P. Bühlmann, "High-dimensional graphs and variable selection with the lasso," *The annals of statistics*, pp. 1436– 1462, 2006.
- [13] A. A. Tsonis, G. Wang et al., "Community structure and dynamics in climate networks," *Climate dynamics*, vol. 37, no. 5-6, pp. 933– 940, 2011.
- [14] J. F. Donges, Y. Zou, N. Marwan, and J. Kurths, "The backbone of the climate network," *EPL (Europhysics Letters)*, vol. 87, no. 4, p. 48007, 2009.
- [15] J. Kawale, M. Steinbach, and V. Kumar, "Discovering dynamic dipoles in climate data." in SDM. SIAM, 2011, pp. 107–118.
- [16] S. Agrawal, G. Atluri et al., "Tripoles: A new class of relationships in time series data," in Proceedings of the 23rd ACM SIGKDD International Conference on Data Mining. ACM, 2017, pp. 697–706.
- [17] V. G. van de Ven, E. Formisano, D. Prvulovic, C. H. Roeder, and D. E. Linden, "Functional connectivity as revealed by spatial independent component analysis of fmri measurements during rest," *Human brain mapping*, vol. 22, no. 3, pp. 165–178, 2004.



Fig. 9: **Climate Models Inter-comparison based on multipole simulations:** Summary statistics of multipole shown in Figure 8 in different time windows of HadSLP2 data (Figure 9(a)) and CMIP5 climate model datasets (Figures 9(b) and 9(c)). In each plot, the red curve indicates the strength of the linear dependence of the multipole, while the blue curve indicates the highest strength of linear dependence obtained for one of its subsets. The difference in red and blue curves thus indicate the linear gain of the multipole.

- [18] D. Koller and N. Friedman, Probabilistic graphical models: principles and techniques. MIT press, 2009.
- [19] D. Eppstein and D. Strash, "Listing all maximal cliques in large sparse real-world graphs," in *International Symposium on Experimental Algorithms*. Springer, 2011, pp. 364–375.
- [20] https://github.com/darrenstrash/quick-cliques.
- [21] R. Kistler, W. Collins, S. Saha, G. White, J. Woollen, E. Kalnay et al., "The NCEP-NCAR 50-year reanalysis: Monthly means CD-ROM and documentation," Bulletin of the American Meteorological society, vol. 82, no. 2, pp. 247–267, 2001.
  [22] J. S. Anderson, M. A. Ferguson, M. Lopez-Larson, and
- [22] J. S. Anderson, M. A. Ferguson, M. Lopez-Larson, and D. Yurgelun-Todd, "Reproducibility of single-subject functional connectivity measurements," *American journal of neuroradiology*, vol. 32, no. 3, pp. 548–555, 2011.
- [23] N. Tzourio-Mazoyer *et al.*, "Automated anatomical labeling of activations in spm using a macroscopic anatomical parcellation of the mni mri single-subject brain," *Neuroimage*, vol. 15, no. 1, pp. 273–289, 2002.
- [24] G. A. Vecchi and A. T. Wittenberg, "El niño and our future climate: where do we stand?" Wiley Interdisciplinary Reviews: Climate Change, vol. 1, no. 2, pp. 260–270, 2010.
- [25] A. Tachibana, J. A. Noah *et al.*, "Parietal and temporal activity during a multimodal dance video game: an fnirs study," *Neuroscience letters*, vol. 503, no. 2, pp. 125–130, 2011.
  [26] A. P. Saygin and M. I. Sereno, "Retinotopy and attention in human
- [26] A. P. Saygin and M. I. Sereno, "Retinotopy and attention in human occipital, temporal, parietal, and frontal cortex," *Cerebral Cortex*, vol. 18, no. 9, pp. 2158–2168, 2008.



Saurabh Agrawal is a PhD. candidate at University of Minnesota in Computer Science. His dissertation focuses on finding novel multivariate relationships in time series and/or spatiotemporal data. He received the Bachelors of Technology degree from the Indian Institute of technology in Bombay, India in 2012.



**Michael Steinbach** is a researcher in the Department of CSE at UMN. Steinbach's research interests include data mining, healthcare, bio-informatics, and statistics. Steinbach received his PhD in CS, M.S in CS and Statistics, and B.S in Math from UMN.



**Daniel Boley** is a Full Professor in Computer Science and Engineering at University of Minnesota. His research interests include numerical linear algebra methods, parallel algorithms, iterative methods for matrix eigenproblems, and error correction for floating point computations. He received the B.A. summa cum laude (Hons.) degree in mathematics from Cornell University (1974), and the M.S. (1976) and Ph.D. (1981) degrees in CS from Stanford University.







Snigdhansu Chatterjee is a Professor in the School of Statistics, and the Director of the Institute for Research in Statistics and its Applications (IRSA, http://irsa.stat.umn.edu/) at the University of Minnesota. His research interests include Data Science theory and methods and Big Data applications to climate sciences, networks and graphs, neuroimaging, and social sciences.

**Gowtham Atluri** is an Assistant Professor in the Department of Electrical Engineering and Computer Science (EECS) at University of Cincinnati. Atluris research interests include data mining, neuroimaging, and climate science. Atluri received his PhD in Computer Science (CS) from UMN and M.Tech in CS from IIT Roorkee.

Anh The Dang is a PhD student in Computer Science at the University of Cincinnati. Anh received his BS in Mechatronics and Computer Science from the Ho Chi Minh City University, Vietnam (2006) and the University of Cincinnati (2018), respectively. He received his MBA from the Newcastle University, England in 2011.



Stefan Liess is currently a researcher in the Department of Soil, Water, and Climate at the University of Minnesota. He analyzes observational data and performs dynamical model simulations to study global and regional climate variability and climate change. He received his M.Sc. (1997) and Ph.D. (2002) degrees in meteorology from the University of Hamburg in collaboration with the Max Planck Institute for Meteorology.



**Vipin Kumar** is a Regents Professor and holds William Norris Chair in Computer Science and Engineering at the University of Minnesota. His research interests include data mining, high performance computing, and their applications in climate/ecosystems and biomedical domains. He received his PhD in CS from University of Maryland College Park, ME in EE from Philips International Institute Eindhoven, and B.E in Electronics & Communication Engineering from IIT Roorkee.