

Technical Report

Department of Computer Science
and Engineering
University of Minnesota
4-192 Keller Hall
200 Union Street SE
Minneapolis, MN 55455-0159 USA

TR 14-008

Algorithms for Mining the Coevolving Relational Motifs in Dynamic
Networks

Rezwan Ahmed and George Karypis

March 24, 2014

Algorithms for Mining the Coevolving Relational Motifs in Dynamic Networks

Rezwan Ahmed and George Karypis

Department of Computer Science & Engineering,
University of Minnesota, Minneapolis, MN 55455
ahmed@cs.umn.edu, karypis@cs.umn.edu

Abstract

Computational methods and tools that can efficiently and effectively analyze the temporal changes in dynamic complex relational networks enable us to gain significant insights regarding the entity relations and their evolution. This paper introduces a new class of dynamic graph patterns, referred to as coevolving relational motifs (CRMs), which are designed to identify recurring sets of entities whose relations change in a consistent way over time. Coevolving relational motifs can provide evidence to the existence of, possibly unknown, coordination mechanisms by identifying the relational motifs that evolve in a similar and highly conserved fashion. We developed an algorithm to efficiently analyze the frequent relational changes between the entities of the dynamic networks and capture all frequent coevolutions as CRMs. Our algorithm follows a depth-first exploration of the frequent CRM lattice and incorporates canonical labeling for redundancy elimination. Experimental results based on multiple real world dynamic networks show that the method is able to efficiently identify CRMs. In addition, a qualitative analysis of the results shows that the discovered patterns can be used as features to characterize the dynamic network.

1 Introduction

Networks are generic models used in various applications from different domains to model the relations between various entities. Examples of some widely studied networks include the friends-networks of popular social networking sites like Facebook and Myspace [6], the Enron email network [5, 11], co-authorship and citation networks [32, 28], web-page linking networks [7, 31], protein-protein interaction networks [17, 37], and networks derived from the IMDB movie database[23]. Many of these networks are dynamic in nature as the entities and relations that need to be modeled change over time. With the emergence of new application areas, there has been a great need to analyze the temporal changes underlying many of these systems, and to develop tools capable of capturing the dynamic aspects of the data, as well as analyze and process this data. Such analysis can identify relational patterns that provides strong observational evidence of the existence of mechanisms that control, coordinate, and trigger the evolutionary changes, which can be used to focus further studies and analysis.

Dynamic networks have recently emerged as an important research area. The focus of this growing research has been mainly defining important recurrent structural patterns and developing algorithms for their identification. In recent years, studies have been done for finding frequent patterns [5], clustering [10], characterizing network evolution rules [4], detecting related cliques [9, 36], identifying conserved relational states [1], finding subgraph subsequences [20], and identifying co-evolution patterns capturing attribute trends [22, 12] in dynamic networks. Although the existing techniques can detect various patterns in a dynamic network, they have not focused on identifying frequent conserved changes of the relational patterns over time. Understanding the evolution of relational patterns that occur frequently can provide evidence of, possibly unknown, coordination mechanisms and to the existence of external factors that are responsible for changing the stable relational patterns in the dynamic network.

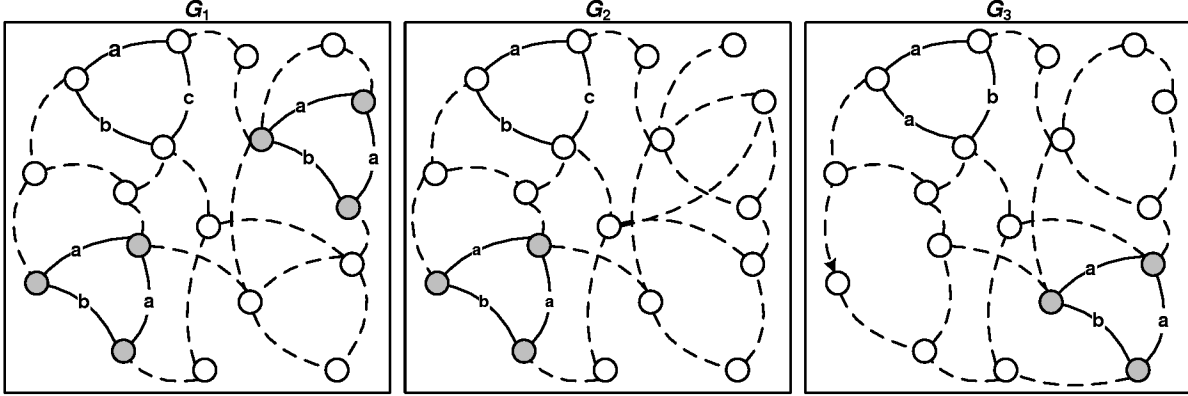


Figure 1: Examples of relational motifs.

Our contribution in this paper is two fold: Firstly, we introduce a new class of patterns, referred to as coevolving relational motifs (CRM), designed to capture patterns that change in a consistent way in a dynamic networks. CRMs identify consistent patterns of relational motif evolution that can provide valuable insights on the processes of the underlying networks. Secondly, we present an algorithm, referred to as CRMminer, to efficiently mine a subclass of these patterns by identifying all frequent coevolving relational motifs such that the motifs that make up the CRM share at least one relation (i.e., an edge) that changes over time. Our algorithm follows a depth-first exploration of the frequent CRM lattice and incorporates canonical labeling for redundancy elimination. We impose both frequency based and node overlap based constraints for pruning the search space to increase efficiency and an approximate pruning to reduce the complexity of our algorithm.

We provide a comprehensive evaluation of the performance of CRMminer and the usefulness of the discovered patterns. We experimentally evaluate the performance and scalability of CRMminer on a large co-authorship network, on a cell culture bioprocess network (multivariate time series data), and on a market sales network (multivariate time series data) by varying different input parameters. Furthermore, our experiments show that the approximate version of CRMminer is able to identify the majority of the CRMs while reducing the amount of time that is required. In addition, we investigate some discovered coevolving relational motifs from all three datasets and provide a qualitative analysis of the information captured in them. We show that some of the discovered CRMs capture relational changes between the nodes that are thematically different (i.e., edge labels transition between two clusters of topics that have very low similarity). In addition, we investigate the extent to which these discovered CRMs can lead to high quality features to build a predictive model. Our results, in the context of a bioprocess dataset, show that the discovered CRMs are present mostly as part of the high yield production runs.

The rest of the paper is organized as follows. Section 2 reviews some graph-related definitions and introduces the notation that is used in the paper. Section 3 captures a brief survey of the related research in this area. Section 4 introduces the concept of coevolving relational motifs and Section 5 describes in detail the algorithm to detect such patterns. Section 6 describes the datasets and metrics used in the experimental evaluation of our algorithm. Section 7 presents a detailed analysis of the experimental results and Section 8 provides some concluding remarks and future directions.

2 Definitions and Notation

A relational network is represented via labeled graphs. A *labeled graph* $G = (V, E, L[V], L[E])$ is composed of a set of nodes V modeling the entities of the network, a set of edges E modeling the relations between these entities, a set of node labels $L[V]$ modeling the type of the entities ($|V| = |L[V]|$), and a set of edge labels $L[E]$ modeling the type of the relations ($|E| = |L[E]|$). The labels assigned to the vertices (edges) are typically not unique and multiple vertices (edges) can have the same label. A *subgraph* $G' = (V', E', L[V'], L[E'])$ of G is a graph such that $V' \subseteq V$, $E' \subseteq E \cap (V' \times V')$. An *induced subgraph* $G'' = (V'', E'', L[V''], L[E''])$ of G

is a graph such that $V'' \subseteq V$, $E'' \subseteq E$ and $\forall(u, v) \in E$ such that $v \in V''$ and $u \in V''$, $(u, v) \in E''$.

Given a connected graph $G = (V, E)$ and a depth-first search (DFS) traversal of G , its *depth-first search tree* T is the tree formed by the forward edges of G . All nodes of G are encoded with subscripts to order them according to their discovery time. Given a DFS tree T of graph G containing n nodes, the root node is labeled as (v_0) and the last discovered node is labeled as (v_{n-1}) . The *rightmost path* of the DFS tree T is the path from vertex v_0 to vertex v_{n-1} .

A *dynamic network* $\mathcal{N} = \langle G_1, G_2, \dots, G_T \rangle$ is modeled as a finite sequence of graphs, where each $G_t = (V, E_t, L_t[V], L_t[E_t])$ is a labeled graph describing the state of the system at a discrete time interval t . The term *snapshot* will be used to refer to each of the graphs in the sequence. Snapshots are assumed to contain the same set of nodes, which will also be referred to as the nodes of \mathcal{N} , denoted by $V_{\mathcal{N}}$, but potentially different sets of edges and node/edge labels. When nodes appear or disappear over time, the set of nodes of each snapshot is the union of all the nodes over all snapshots. Also, the nodes across the different snapshots are numbered consistently, so that the i th node of G_k ($1 \leq k \leq T$) will always correspond to the same i th node of \mathcal{N} .

We define the *span sequence* of an edge as the sequence of maximal-length time intervals in which an edge is present in a consistent state. An edge (u, v) is in a *consistent state* over a maximal time interval $s:e$ if it is present in all snapshots G_s, \dots, G_e with the same label l and it is different in both G_{s-1} and G_{e+1} (assuming $s > 1$ and $e < T$). The span sequence of an edge will be described by a sequence of vertex labels, edge labels and time intervals of the form $\langle (l_{u_1}, l_{e_1}, l_{v_1}, s_1:e_1), \dots, (l_{u_n}, l_{e_n}, l_{v_n}, s_n:e_n) \rangle$, where $l_{u_i}, l_{v_i} \in L[V]$, $l_{e_i} \in L[E]$, $s_i \leq e_i$, and $e_i \leq s_{i+1}$.

An *injection* is defined as a function $f : A \rightarrow B$ such that $\forall a, b \in A$, if $f(a) = f(b)$, then $a = b$. A function $f : A \rightarrow B$ is a bijective function or *bijection*, iff $\forall b \in B$, there is a unique $a \in A$ such that $f(a) = b$. A function composition $g \circ f$ implies that for function $f : X \rightarrow Y$ and $g : Y \rightarrow Z$, then *composite* function $g \circ f : X \rightarrow Z$.

A *relational motif* is a subgraph that occurs frequently in a single snapshot or a collection of snapshots. In Figure 1, the three-vertex subgraph consisting of the shaded nodes that are connected via the labeled edges corresponds to a relational motif that occurs a total of four times (twice in G_1 and once in each of G_2 and G_3). The set of nodes that support the multiple occurrences of a relational motif do not need to be the same. In order to determine if a snapshot supports a relational motif (and how many times), we need to perform subgraph isomorphism operations (i.e., identify the embeddings of the relational motif’s graph pattern). We will use M to denote a relational motif and the underlying subgraphs will be denoted by the tuple $(N, A, L[N], L[A])$, where N is the set of nodes, A is the set of edges (arcs), $L[N]$ is the set of node labels, and $L[A]$ is the set of edge labels.

3 Related Work

Unlike the problem of finding patterns in static networks for which there exist a considerable body of research [25, 34, 43, 2, 21, 18, 26, 16, 35, 41], the problem of mining dynamic networks is comparatively less studied.

Borgwardt et. al. [5] introduced the notion of the *dynamic subgraph*, which extends the traditional notion of the subgraph to include the sequence of subgraphs that exist in a consecutive sequence of snapshots and developed algorithms to identify the set of dynamic subgraphs that occur frequently in a dynamic network. Jin et. al. [22] introduced the notion of the *trend motif*, which is a connected subgraph with nodes containing time series and each node’s time series exhibits a consistent trend over a time interval. Their work focused only on developing algorithms for finding frequently occurring trend motifs that show either an increasing or a decreasing trend. With a similar objective, Desmier et. al. [12] defined the problem of mining cohesive co-evolving patterns which capture the local co-evolution of similar vertices at several timestamps based on their attributes trends. Lahiri et. al. [27] proposed a mining problem of finding periodic or near periodic subgraphs in dynamic networks.

Berlingerio et. al. [4] also provided an algorithm to detect frequent subgraphs in time-evolving graphs for deriving graph-evolution rules that satisfy a minimum confidence constraint. Robardet [36] represented the frequent patterns of a graph as *pseudo-cliques* and proposed an algorithm that first mines each graph snapshot of a dynamic graph for local patterns and then combines these with patterns from previous snapshot

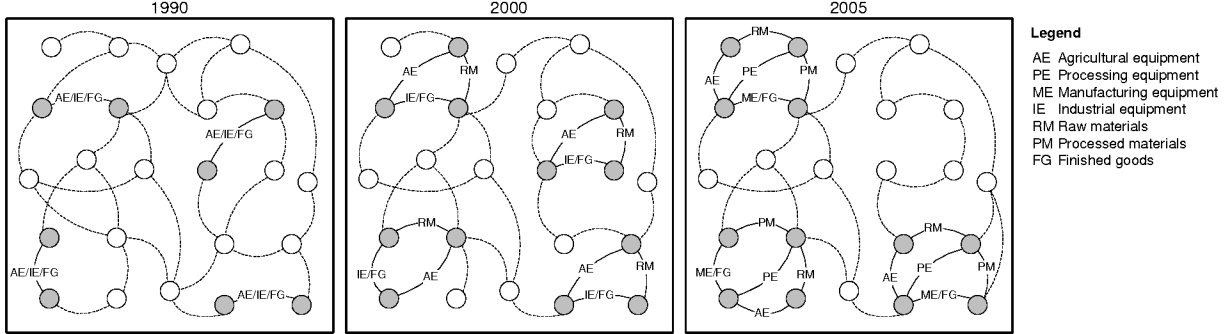


Figure 2: An example of a coevolving relational motif.

based on some constraints to form evolving patterns. Inokuchi et. al. [19, 20] solved a similar problem of finding frequent induced subgraph subsequences from graph sequence data and capturing the changes of a subgraph over the subsequence. Ahmed et. al. [1] introduced a new class of patterns that captures the time-persistent relations among the nodes, called *relational states* and presented an algorithm to identify all maximal non-redundant evolution paths of the stable relational states of a dynamic network.

A related body of research has investigated the task of identifying and tracking patterns in biological networks [42, 24, 39], evolving communities in social networks [3, 13], and evolutionary clustering [10, 38]. Even though these methods provide valuable insights on the evolution of dynamic networks, the nature of co-evolving relational patterns focused in this paper is different from the various types of evolution addressed in other papers.

4 Coevolving Relational Motifs

Coevolving relational motifs are designed to identify the relational patterns that change in a consistent way over time. An example of this type of conservation is illustrated in Figure 2, in the context of a hypothetical country-to-country trading network where labels represent the commodities been traded. The network for 1990 shows a simple relational motif (M_1) between pairs of nodes that occurs four times (shaded nodes and solid labeled edges). This relational motif has evolved in the network for 2000 in such a way so that in all four cases, a new motif (M_2) that includes an additional node has emerged. Finally, in the network for 2005 we see that three out of these four occurrences have evolved to a new motif (M_3) that now involves four nodes. This example shows that the initial relational motif among the four sets of nodes has changed in a fairly consistent fashion over time (i.e., it *coevolved*) and such a sequence of motifs $M_1 \rightsquigarrow M_2 \rightsquigarrow M_3$ represents an instance of a CRM.

CRMs identify consistent patterns of relational motif evolution that can provide valuable insights on the processes of the underlying networks. For example, the CRM of Figure 2 captures the well-known phenomenon of production specialization due to economic globalization, in which the production of goods have been broken down into different components performed by different countries [14]. Similarly, CRMs in health-care networks can capture how the set of medical specialties required to treat certain medical conditions have changed over the years, in communication networks CRMs can capture the evolution of themes being discussed among groups of individuals as their lifestyles change, whereas CRMs in corporate email networks can capture how the discussion related to certain topics moves through the companies' hierarchies.

The formal definition of a CRM that is used in this paper is as follows:

Definition 1 A CRM of length m is a tuple $\{N, \langle M_1, \dots, M_m \rangle\}$, where N is a set of vertices and each $M_j = (N_j, A_j)$ is a relational motif defined over a subset of the vertices of N that satisfies the following constraints:

- i) it occurs at least ϕ times,
- ii) each occurrence uses a non-identical set of nodes,
- iii) $M_j \neq M_{j+1}$, and

iv) $|N_j| \geq \beta|N|$ where $0 < \beta \leq 1$.

A relational motif M_j is *defined* over a subset of vertices N if there is a injection ξ_j from N_j to N . An m -length CRM *occurs* in a dynamic network whose node set is V if there is a sequence of m snapshots $\langle G_{i_1}, G_{i_2}, \dots, G_{i_m} \rangle$ and a subset of vertices B of V (i.e., $B \subseteq V$) such that:

- i) there is a bijection ξ from N to B
- ii) the injection $\xi \circ \xi_j$ is an embedding of M_j in G_{i_j}
- iii) there is no embedding of M_j via the injection $\xi \circ \xi_j$ in $G_{i_{j+1}}$ or no embedding of M_{j+1} via the injection $\xi \circ \xi_{j+1}$ in G_{i_j} .

For example, the number of occurrences of the CRM in Figure 2 is 3. Note that the third condition in the above definition is designed to ensure that for each pair of successive motifs at least one of them is not supported by the snapshot-nodes pair that supported the other motif. This is done to ensure that there is a relational change between the nodes associated with those embeddings in each others snapshot.

The purpose of the parameters ϕ and β in Definition 1 are as follows: The parameter ϕ is used to eliminate sequences of evolving motifs that are not frequent enough to indicate the existence of an underlying process driving these changes. The parameter β is used to control the degree of change between the sets of nodes involved in each motif of a CRM and enforces a minimum node overlap among all motifs of CRM. Finally, the third constraint of Definition 1 limits the discovered CRMs to only an evolving sequence of motifs and not a sequence that remains the same. Note that the frequent dynamic subgraphs introduced by Borgwardt et. al. [5] (Section 3) correspond to CRMs in which the snapshots supporting each set of nodes are restricted to be consecutive and $\beta = 1$.

In this paper we focus on developing an efficient algorithm to mine a subclass of the CRMs, such that in addition to the conditions mentioned in Definition 1, the motifs that make up the CRM, also share at least one edge that itself is a CRM. Formally, we focus on identifying the CRMs $c = \{N_c, \langle M_1, \dots, M_m \rangle\}$ that contain at least a pair of vertices $\{u, v\} \in N_C$ such that each induced subgraph M'_i of M_i on $\{u, v\}$ is connected and $x = \{\{u, v\}, \langle M'_1, \dots, M'_m \rangle\}$ is a CRM. A CRM like x that contains only one edge and two vertices will be called an *anchor*. We focus our CRM enumeration problem around the anchors, since they ensure that the CRM's motifs contain at least a pair of nodes in common irrespective of the specified overlap constraint that evolves in a conserved way. It also characterizes how the network around these core set of entities coevolved with them. We will refer to the class of CRMs that contain an anchor as *anchored CRMs*. For the rest of the discussion, any references to a CRM will assume it is an anchored CRM.

Given the above definition, the work in this paper is designed to develop efficient algorithm for solving the following problem:

Problem 1 *Given a dynamic network \mathcal{N} containing T snapshots, a user defined minimum support ϕ ($1 \leq \phi$), a minimum number of edges k_{\min} per CRM, and a minimum number of motifs m_{\min} per CRM, find all CRMs such that the motifs that make up the CRM, also share an anchor CRM.*

A CRM that meets the requirements specified in Problem 1 is referred as a *frequent CRM* and it is *valid* if it also satisfies the minimum node overlap constraint (Definition 1(iv)).

5 Finding Coevolving Relational Motifs

A consequent of the way anchored CRMs are defined is that the number of motifs that they contain is exactly the same as the number of motifs that exist in the anchor(s) that they contain. As a result, the CRMs can be identified by starting from all the available anchors and try to grow them by repeatedly adding edges as long as the newly derived CRMs satisfy the constraints and have exactly the same number of motifs as the anchor. Since a CRM can contain more than one anchor, this approach may identify redundant CRMs by generating the same CRM from multiple anchors. Therefore, the challenge is to design a strategy that is complete and non-redundant. To achieve this, we develop an approach that generates each CRM from a unique anchor. Given a CRM, the anchor from which it will be generated is referred to as its *seed anchor*.

The algorithm that we developed, named CRMminer, for finding all non-redundant valid CRMs (Problem 1) initially identifies all frequent anchors and then performs a depth-first exploration of each anchor pattern space along with a canonical labeling that is derived by extending the ideas of the minimum DFS

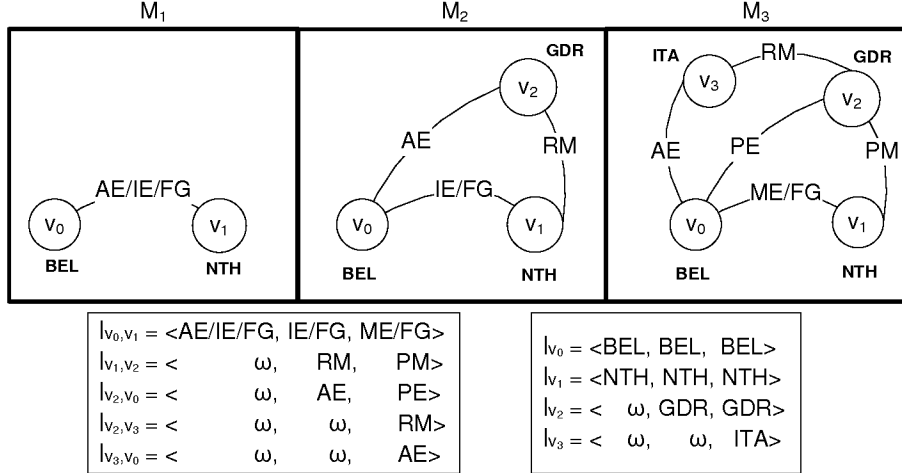


Figure 3: A CRM Representation.

code [40] to the case of CRMs for redundancy elimination. We impose frequency-based constraints by stopping any further exploration of a CRM when the pattern does not occur at least ϕ times in the dynamic network \mathcal{N} .

5.1 CRM Representation

A CRM $c = \{N, \langle M_1, \dots, M_m \rangle\}$ is represented as a graph $G_c = (N, E_c)$, such that an edge $(u, v) \in E_c$ is a 5-item tuple $(u, v, l_u, l_v, l_{u,v})$, where $u, v \in N$, the vectors l_u and l_v contain the vertex labels and $l_{u,v}$ contains the edge labels of all motifs. If the CRM consists of m motifs, then $l_u = \langle l_{u_1}, \dots, l_{u_m} \rangle$, $l_v = \langle l_{v_1}, \dots, l_{v_m} \rangle$ and $l_{u,v} = \langle l_{u_1,v_1}, \dots, l_{u_m,v_m} \rangle$. The k th entry in each vector l_u , $l_{u,v}$, and l_v records the connectivity information among the vertices of the k th motif (M_k). If an edge (u, v) is part of motif M_k , then the k th entry of l_u , l_v , and $l_{u,v}$ are set to the labels of u and v vertices, and the label of the (u, v) edge respectively. If both vertices u and v are part of motif M_k , but the (u, v) edge is not, or at least one of the vertices u or v is not part of the M_k motif (i.e., no (u, v) edge is possible), then ω is inserted at the k th entry of the $l_{u,v}$ to capture the disconnected state. Similarly, if u or v does not have any incident edges in the M_k motif (i.e., the vertices are not present in that motif), then ω is added as the vertex label at the k th entry of l_u or l_v . Note that the value of ω is lexicographically greater than the maximum edge and vertex label. This representation is illustrated in Figure 3.

5.2 Mining Anchors

The search for CRMs is initiated by locating the frequent anchors that satisfy the CRM definition and the restrictions defined in Problem 1. This is done as following: Given a dynamic network \mathcal{N} , we sort all the vertices and edges by their label frequency and remove all infrequent vertices and edges. Remaining vertices and edges are relabeled in decreasing frequency. We determine the span sequences of each edge and list every edge's span sequence if that sequence contains at least a span with an edge label that is different from the rest of the spans. At this point, we use the sequential pattern mining technique prefixSpan [33] to determine all frequent span sequences. Since the frequent sequences can be partial sequences of the original input span sequences, it is not guaranteed that they all include at least one span with a different edge label. Thus, the frequent sequences that contain at least one different span in terms of edge label are considered as the anchors. The number of spans in a frequent span sequence corresponds to the total number of motifs in the anchor.

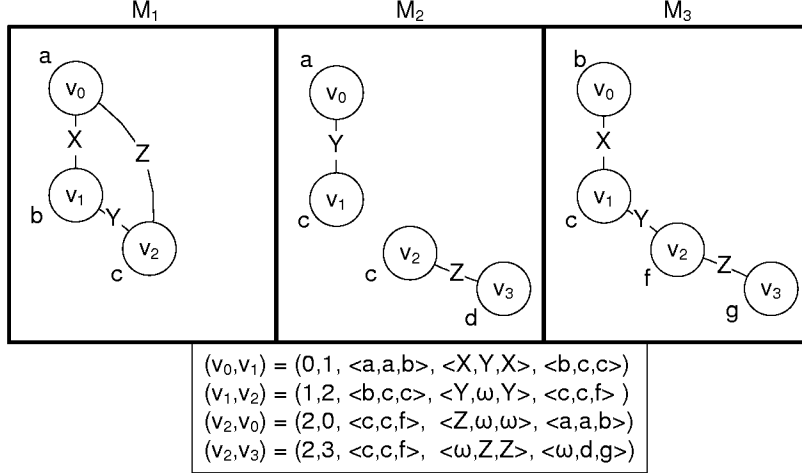


Figure 4: DFS tree for a CRM.

5.3 CRM Enumeration

Given an anchor c , we generate the set of desired CRMs by growing the size of the current CRM one edge at a time following a depth-first approach. To ensure that each CRM is generated only once in the depth-first exploration, we use an approach similar to the gSpan algorithm [40], which we have extended for the problem of CRM mining.

gSpan explores the frequent pattern lattice in a depth-first fashion. The pattern lattice is represented as a hierarchical search space where each node corresponds to a connected frequent pattern, the highest node being an empty pattern (i.e., a single vertex), the next level nodes represent 1-edge patterns, and so on. The n th level nodes, which represent n -edge patterns, contain one more edge than the corresponding $(n-1)$ level nodes. To ensure that each frequent pattern in this lattice is visited exactly once, gSpan's exploration amounts to visiting the nodes of the lattice (i.e., frequent patterns) by traversing a set of edges that form a spanning tree of the lattice. This spanning tree is defined by assigning to each node of the lattice a canonical label, called the minimum DFS code. A DFS code of a graph, is a unique label that is formed based on the sequence of edges added to that node during a depth-first exploration. The minimum DFS code, is the DFS code that is lexicographically the smallest. Given this canonical labeling, the set of lattice edges that are used to form the spanning tree correspond to the edges between two successive nodes of the lattice (parent and child) such that the minimum DFS code of the child can be obtained by simply appending the extra edge to the minimum DFS code of the parent. For example, given a DFS code $\alpha = \langle a_0, a_1, \dots, a_m \rangle$, a valid child DFS code is $\gamma = \langle a_0, a_1, \dots, a_m, b \rangle$, where b is the new edge. This spanning tree guarantees that each node has a unique parent and all nodes are connected [40]. To efficiently grow a node, gSpan generates the child nodes by only adding those edges that originate from the vertices on the rightmost path of the DFS-tree representation of the parent node. It then checks whether the resulting DFS code of the child node corresponds to the minimum DFS code. Construction of the child nodes generated by adding other edges (i.e., not from the rightmost path) are skipped, since such child nodes will never contain a DFS code that corresponds to the minimum DFS code.

In order to apply the ideas introduced by gSpan to the problem of efficiently mining CRMs, we need to develop approaches for (i) representing the DFS code a CRM, (ii) ordering the DFS codes of a CRM using the DFS lexicographic ordering, (iii) representing the minimum DFS code of a CRM to use as the canonical label, and (iv) extending a DFS code of a CRM by adding an edge. Once properly defined, the correctness and completeness of frequent CRM enumeration follows directly from the corresponding proofs of gSpan.

5.3.1 DFS code of a CRM

In order to derive a DFS code of a CRM, we need to develop a way of ordering the edges. Given a CRM c , represented as a graph $G_c = (N, E_c)$, we perform a depth-first search in G_c to build a DFS tree T_c . The

vertices (N) are assigned subscripts from 0 to $n-1$ for $|N| = n$ according to their discovery time. The edges (E_c) are grouped into two sets: the forward edge set $E_{T_c, fw} = \{(v_i, v_j) \in E_c \mid i < j\}$ and the backward edge set $E_{T_c, bw} = \{(v_i, v_j) \in E_c \mid i > j\}$. Let us denote a partial order on $E_{T_c, fw}$ as $\prec_{T_c, fw}$, a partial order on $E_{T_c, bw}$ as $\prec_{T_c, bw}$, and a partial order on E_c as $\prec_{T_c, bw+fw}$. Given two edges $e_1 = (v_{i_1}, v_{j_1})$ and $e_2 = (v_{i_2}, v_{j_2})$, the partial order relations are defined as:

- a) $\forall e_1, e_2 \in E_{T_c, fw}$, if $j_1 < j_2$, then $e_1 \prec_{T_c, fw} e_2$.
- b) $\forall e_1, e_2 \in E_{T_c, bw}$, if $i_1 < i_2$ or $(i_1 = i_2 \text{ and } j_1 < j_2)$, then $e_1 \prec_{T_c, bw} e_2$.
- c) $\forall e_1 \in E_{T_c, bw}$ and $\forall e_2 \in E_{T_c, fw}$, if $i_1 < j_2$, then $e_1 \prec_{T_c, bw+fw} e_2$.
- d) $\forall e_1 \in E_{T_c, fw}$ and $\forall e_2 \in E_{T_c, bw}$, if $j_1 \leq i_2$, then $e_1 \prec_{T_c, bw+fw} e_2$.

The combination of the three partial orders defined above enforces a linear order \prec_{T_c, E_c} on E_c .

Given this linear order \prec_{T_c, E_c} , we can order all edges in G_c and construct an edge sequence to form a DFS code of a CRM, denoted as $code(c, T_c)$. An edge of the DFS code of a CRM is represented similar to the CRM edge definition and uses a 5-tuple representation $(i, j, \mathbf{l}_i, \mathbf{l}_{i,j}, \mathbf{l}_j)$, where i and j are the DFS subscripts (i.e., the discovery time) of the vertices, and \mathbf{l}_i , \mathbf{l}_j , and $\mathbf{l}_{i,j}$ are the label vectors of the vertices and edge, respectively. The k th entry in each vector \mathbf{l}_i , \mathbf{l}_j , and $\mathbf{l}_{i,j}$ contains the labels of vertices and edges of motif M_k .

For example, a DFS tree of a CRM illustrated in Figure 4 contains the DFS code (v_0, v_1) , (v_1, v_2) , (v_2, v_0) , and (v_2, v_3) . The DFS code for the sequence of edges are: $(0, 1, \langle a, a, b \rangle, \langle X, Y, X \rangle, \langle b, c, c \rangle)$, $(1, 2, \langle b, c, c \rangle, \langle Y, \omega, Y \rangle, \langle c, c, f \rangle)$, $(2, 0, \langle c, c, f \rangle, \langle Z, \omega, \omega \rangle, \langle a, a, b \rangle)$, and $(2, 3, \langle c, c, f \rangle, \langle \omega, Z, Z \rangle, \langle \omega, d, g \rangle)$. Note that the k th entry of a vertex and edge label vector of a DFS code is filled with ω if the corresponding vertex or edge is not present in motif M_k . DFS code's Neighborhood restriction property defined in [40] still holds for DFS code of a CRM.

5.3.2 DFS Lexicographic Ordering

To establish a canonical labeling system for a CRM, CRMminer defines the DFS lexicographical ordering based on the CRM's DFS code definition. The linear ordering is defined as follows. Let two DFS codes of a CRM consisting of m motifs be $\alpha = code(c_\alpha, T_\alpha) = \langle e_\alpha^0, e_\alpha^1, \dots, e_\alpha^p \rangle$ and $\delta = code(c_\delta, T_\delta) = \langle e_\delta^0, e_\delta^1, \dots, e_\delta^q \rangle$. Let the forward and backward edge set for T_α and T_δ be $E_{\alpha, fw}$, $E_{\alpha, bw}$, $E_{\delta, fw}$, and $E_{\delta, bw}$, respectively. Also, let $e_\alpha^x = (i_\alpha^x, j_\alpha^x, \mathbf{l}_{i_\alpha^x}, \mathbf{l}_{i_\alpha^x, j_\alpha^x}, \mathbf{l}_{j_\alpha^x})$ and $e_\delta^y = (i_\delta^y, j_\delta^y, \mathbf{l}_{i_\delta^y}, \mathbf{l}_{i_\delta^y, j_\delta^y}, \mathbf{l}_{j_\delta^y})$ be two edges, one in each of the α and δ DFS codes. Now, $e_\alpha^x < e_\delta^y$, if any of the following is true:

- i) $e_\alpha^x \in E_{\alpha, bw}$ and $e_\delta^y \in E_{\delta, fw}$, or
- ii) $e_\alpha^x \in E_{\alpha, bw}$, $e_\delta^y \in E_{\delta, bw}$, and $j_\alpha^x < j_\delta^y$, or
- iii) $e_\alpha^x \in E_{\alpha, bw}$, $e_\delta^y \in E_{\delta, bw}$, $j_\alpha^x = j_\delta^y$ and $\mathbf{l}_{i_\alpha^x, j_\alpha^x} < \mathbf{l}_{i_\delta^y, j_\delta^y}$, or
- iv) $e_\alpha^x \in E_{\alpha, fw}$, $e_\delta^y \in E_{\delta, fw}$, and $i_\delta^y < i_\alpha^x$, or
- v) $e_\alpha^x \in E_{\alpha, fw}$, $e_\delta^y \in E_{\delta, fw}$, $i_\alpha^x = i_\delta^y$ and $\mathbf{l}_{i_\alpha^x} < \mathbf{l}_{i_\delta^y}$, or
- vi) $e_\alpha^x \in E_{\alpha, fw}$, $e_\delta^y \in E_{\delta, fw}$, $i_\alpha^x = i_\delta^y$, $\mathbf{l}_{i_\alpha^x} = \mathbf{l}_{i_\delta^y}$ and $\mathbf{l}_{i_\alpha^x, j_\alpha^x} < \mathbf{l}_{i_\delta^y, j_\delta^y}$, or
- vii) $e_\alpha^x \in E_{\alpha, fw}$, $e_\delta^y \in E_{\delta, fw}$, $i_\alpha^x = i_\delta^y$, $\mathbf{l}_{i_\alpha^x} = \mathbf{l}_{i_\delta^y}$, $\mathbf{l}_{i_\alpha^x, j_\alpha^x} = \mathbf{l}_{i_\delta^y, j_\delta^y}$, and $\mathbf{l}_{j_\alpha^x} < \mathbf{l}_{j_\delta^y}$.

Based on the above definitions, we derive the following conditions to compare two DFS codes of a CRM. We define $\alpha \leq \delta$, iff any of the following conditions is true:

- a) $F_\omega(\alpha) \leq F_\omega(\delta)$, where $F_\omega(x)$ is the number of edges $(\mathbf{l}_{i_\alpha^t, j_\alpha^t})$ that are set to ω , or
- b) $\exists t, 0 \leq t \leq \min(p, q)$, $e_\alpha^k = e_\delta^k$ for $k < t$, and $e_\alpha^t < e_\delta^t$, or
- c) $e_\alpha^k = e_\delta^k$ for $0 \leq k \leq p$ and $p \leq q$.

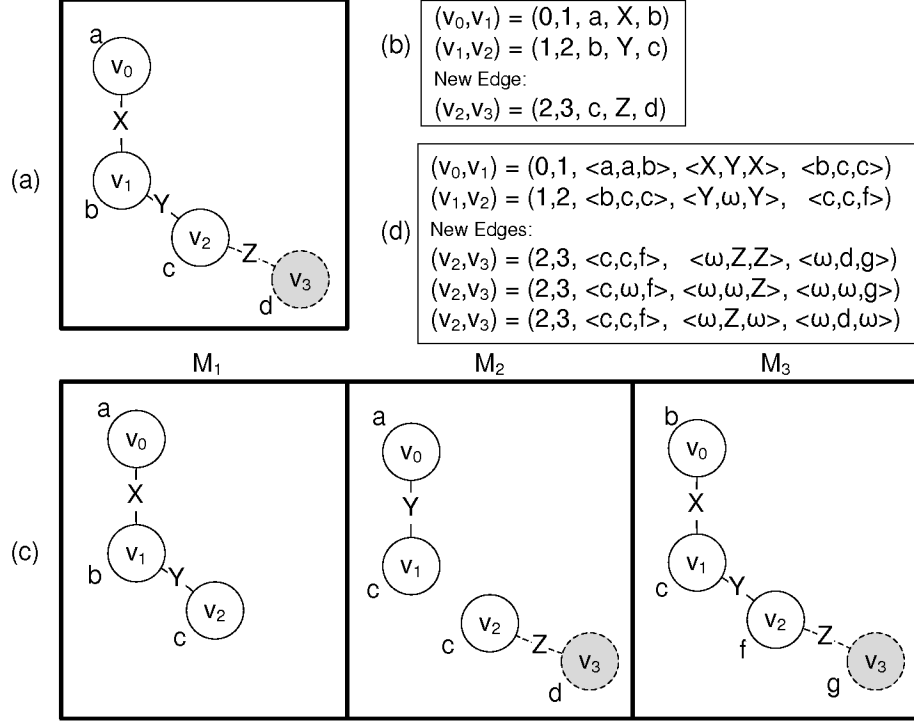


Figure 5: Adding an edge according to rightmost extension rules. (a) Extending a simple graph, (b) DFS code of the simple graph, (c) extending a CRM, and (d) DFS code of the CRM.

Note that the DFS lexicographical ordering ranks edges with label vectors containing no ω labels higher than the edges which does. To define the relation between ω and valid vertex/edge label, the value of ω is set to a lexicographically higher value than the maximum edge and vertex label. This is important as we show later in Section 5.3.5.

5.3.3 Minimum DFS Code

A CRM can be represented by different DFS trees resulting in different DFS codes. To define a canonical label for a CRM, we select the minimum DFS code according to the DFS lexicographic order, represented by $min_code(c)$. Similar to simple graph isomorphism, given two CRMs c and c' , c is isomorphic to c' if and only if $min_code(c) = min_code(c')$. Thus, by searching frequent minimum DFS codes, we can identify the corresponding frequent CRMs.

5.3.4 Pattern lattice growth

The difference between a simple graph pattern and a CRM is the vertex/edge label representation. One contains a single label and the other contains a label vector (i.e., sequence of labels including an empty label ω). The growth of a simple graph by one edge results in a number of simple graphs based on the number of unique labels of the frequent edges on the rightmost path. However, a CRM extended by an edge can generate large number of CRMs, since these are formed based on the combination of the label vectors of the frequent edges on the rightmost path.

In Figure 5, we illustrated one edge growth of a simple graph and a CRM according to rightmost extension. Both the simple graph and the CRM are expanded by adding a frequent edge (v_2, v_3) . In case of the simple graph (Figure 5(a)), the original DFS code consisted sequence of two edges (v_0, v_1) and (v_1, v_2) represented as $(0, 1, a, X, b)$ and $(1, 2, b, Y, c)$. After the one edge extension, the new edge $(2, 3, c, Z, d)$ connected the new vertex v_3 to the rightmost vertex v_2 . When the CRM is considered (Figure 5(c)), the original DFS code consisted the same sequence of edges $\{(v_0, v_1), (v_1, v_2)\}$ represented as: $(0, 1, \langle a, a, b \rangle, \langle X, Y, X \rangle, \langle b, c, c \rangle)$,

$(1, 2, \langle b, c, c \rangle, \langle Y, \omega, Y \rangle, \langle c, \omega, f \rangle)$. Based on the combination of the label vectors of vertices v_2 and v_3 , and edge (v_2, v_3) , we have the following options for the (v_2, v_3) edge: $(2, 3, \langle c, c, f \rangle, \langle \omega, Z, Z \rangle, \langle \omega, d, g \rangle)$, $(2, 3, \langle c, \omega, f \rangle, \langle \omega, \omega, Z \rangle, \langle \omega, \omega, g \rangle)$, and $(2, 3, \langle c, c, f \rangle, \langle \omega, Z, \omega \rangle, \langle \omega, d, \omega \rangle)$. Note that we allow a CRM to grow by an edge that may not be present or frequent in all motifs of that CRM to ensure complete set of the results.

To efficiently determine the frequent candidate edges during rightmost extension, we apply frequent itemset mining technique [15]. Even though there can be significantly more number of child CRMs from one edge extension of a CRM than a simple graph, the extended lexicographic ordering is able to order all candidates and enable us to perform pre-order search on the pattern lattice. Since traversal of the pattern lattice of a CRM remains similar to a simple graph, it allows CRMminer to prune the pattern with non-minimum DFS codes and their descendants similar to gSpan without impacting the completeness of the results.

5.3.5 Algorithm Completeness

To eliminate redundancy during the CRM expansion process, we use minimum DFS code as the canonical label and construct the pattern lattice to ensure that every node (i.e., a CRM) is connected to a unique parent and grown via single edge addition. This process ensures that each potential CRM is only explored once.

To ensure that all discovered CRMs contain at least one anchor, we show how we can identify an anchor of a valid CRM. Given a valid CRM c and its canonical label (i.e., the minimum DFS code) $min_code(c) = \langle e_1, e_2, \dots, e_k \rangle$, where e_i is an edge in the lexicographically ordered edge sequence. We claim that the first edge (e_1) of the canonical label of a CRM (c) is an anchor of that CRM. To prove this claim, assume e_1 is not an anchor and e_x is an anchor, where $1 < x \leq k$. Given the graph in CRM c , we can construct a DFS code for c that starts with e_x as the first edge. Assume, the new DFS code is represented as $code(c) = \langle e_x, e_{p_1}, \dots, e_{p_{k-1}} \rangle$, where $\langle e_{p_1}, \dots, e_{p_{k-1}} \rangle$ is an edge sequence containing a permutation of the edges $\{e_1, \dots, e_k\} \setminus \{e_x\}$. Since e_1 is not an anchor, it contains some ω labels. Based on the DFS lexicographic ordering and e_x being an anchor, $e_x < e_1$. Hence, we can state that $\langle e_x, e_{p_1}, \dots, e_{p_{k-1}} \rangle < \langle e_1, e_2, \dots, e_k \rangle$. This is a contradiction, since $\langle e_1, e_2, \dots, e_k \rangle$ is the minimum DFS code of c . Thus, e_1 is an anchor of CRM c . Given that the first edge is an anchor, then CRMminer will generate that CRM by starting from the anchor and then following its right-most extension rule to add the rest of the edges one by one.

5.4 Search space pruning

One of the challenges that any graph mining algorithm needs to handle is the exponential growth of the search space during enumeration. Traditionally, user specified constraints are used to prune the search space. To ensure discovery of the complete set of patterns, the pruning constraints need to have the anti-monotonicity property. For CRMminer, we use both support measure and minimum overlap constraints to prune the search space.

5.4.1 Minimum Support (ϕ)

To efficiently search patterns in a single large graph using a minimum support constraint, the support measure needs to guarantee the anti-monotonicity property. Bringmann [8] presented the minimum image based support measure to prune the search space in a single large graph. Given a pattern $p = (V_p, E_p, L_p)$ and a single large graph $G = (V_G, E_G, L_G)$, this measure identifies the vertex in p which is mapped to the least number of unique vertices in G and uses this number as the frequency of p in G . To formally define the support measure, let each subgraph g of G that is isomorphic to p be defined as an *occurrence* of p in G . For each occurrence g , there is a function $\varphi : V_p \rightarrow V_G$ that maps the nodes of p to the nodes in G such that (i) $\forall v \in V_p \Rightarrow L_p(v) = L_G(\varphi(v))$ and (ii) $\forall (u, v) \in E_p \Rightarrow (\varphi(u), \varphi(v)) \in E_G$. The minimum image based support of a pattern p in G is defined as:

$$\sigma(p, G) = \min_{v \in V_p} | \{ \varphi_i(v) : \varphi_i \text{ is an occurrence of } p \text{ in } G \} |. \quad (1)$$

This minimum image based support is anti-monotonic [8].

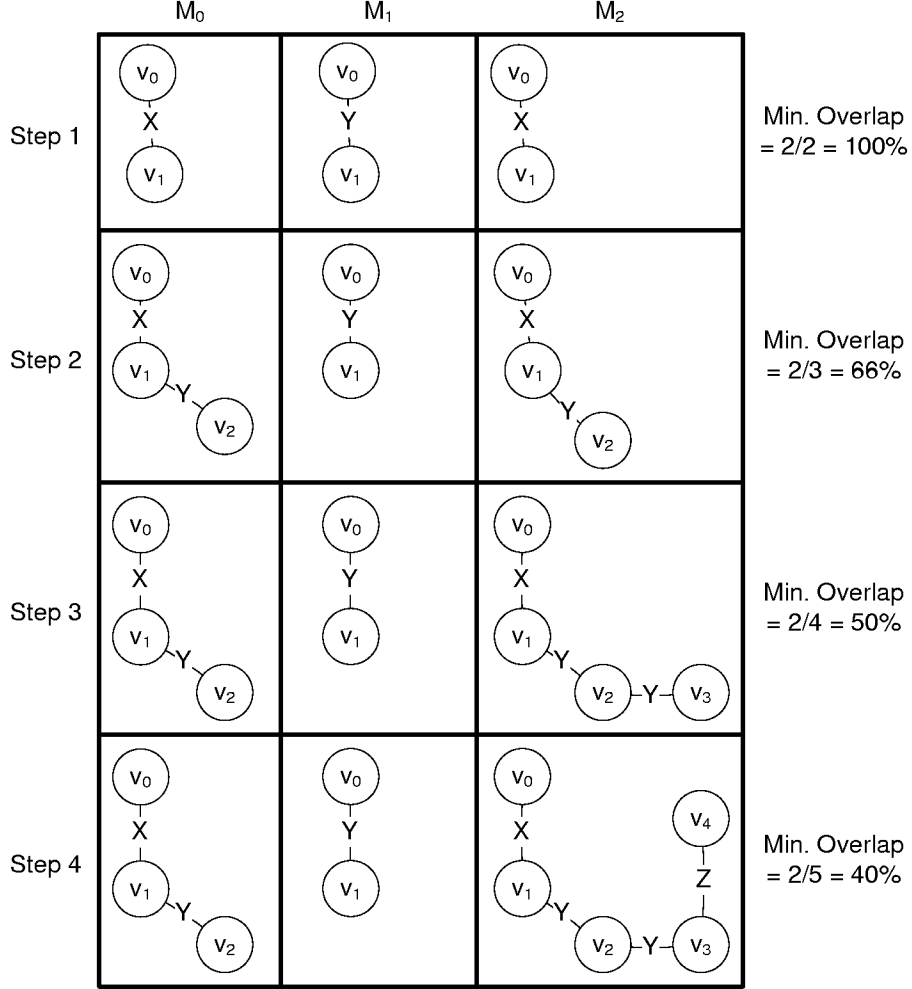


Figure 6: Minimum overlap calculation during a CRM enumeration.

We adopted the minimum image based support measure to calculate the minimum support of a CRM in a dynamic network. As defined in Section 2, a dynamic network \mathcal{N} can be represented as a single large graph where the nodes \mathcal{N} are considered as the vertices of the large graph. Hence, it is possible to calculate the least number of unique vertices of the dynamic network that are mapped to a particular vertex of a CRM. Given a CRM $c = \{N_c, \langle M_1, M_2, \dots, M_m \rangle\}$ in a dynamic network $\mathcal{N} = \{V_{\mathcal{N}}, \langle G_0, G_1, \dots, G_T \rangle\}$ where $m \leq T$, the minimum image based support of c is defined as:

$$\sigma(c, \mathcal{N}) = \min_{v \in V_c} | \{ \varphi_i(v) : \varphi_i \text{ is an occurrence of } c \text{ in } \mathcal{N} \} |. \quad (2)$$

Similar to the support measure of a pattern in a single large graph, by selecting the support of the vertex in c that has the least number of unique mapping in \mathcal{N} , we maintain the anti-monotonicity property.

5.4.2 Minimum Overlap (β)

Each motif of a CRM needs to contain at least a minimum percentage of the nodes from all the nodes of the CRM. This minimum node overlap threshold (defined as β in Section 4) controls the degree of change that is allowed between the sets of nodes in each motif of a CRM. Given a CRM $c = \{N_c, \langle M_1, M_2, \dots, M_m \rangle\}$

containing m motifs, the *minimum overlap* of a CRM is defined as:

$$\rho(c) = \frac{\min_{1 \leq i \leq m} \{|V_{M_i}|\}}{|N_c|}, \quad (3)$$

where V_{M_i} is the set of nodes in motif M_i . Even though the minimum overlap constraint is a reasonable approach to ensure that the motifs that make the CRM are coherent, it is not anti-monotonic [?]. Thus, to generate a complete set of CRMs that meet user specified thresholds of support and overlap, we cannot prune CRMs that do not satisfy this constraint as CRMs derived from it can satisfy the constraint. Hence, we need to enumerate all CRMs that meet the support threshold and then search the output space for CRMs that meet the minimum overlap requirement.

To improve the performance, we developed an approximate version of our algorithm, named CRMminer_x, that discovers a subset of the valid CRMs (i.e., meet the constraints of Definition 1) by pruning the pattern lattice using the minimum overlap threshold. We first check whether a CRM meets the overlap threshold. If it does, we continue the enumeration process. If it does not, then we check whether any of the patterns at the previous level of the pattern lattice from which the current pattern was derived, referred as parent CRMs, meet the overlap threshold. If at least one does, we continue enumeration. If none of the parent CRMs meet the overlap threshold, we prune that CRM.

To approximately calculate the minimum overlap of the parent CRMs, we do not generate all possible parent patterns. The parent is defined to contain one less node than the current CRM; thus, we remove a node $v \in N_c$. In such case, the parent pattern will contain $(|N_c| - 1)$ nodes and each motif M_i may contain $(|V_{M_i}| - 1)$ nodes if $v \in V_{M_i}$ or $(|V_{M_i}|)$ nodes if $v \notin V_{M_i}$. For at least one parent CRM to meet the overlap threshold, we consider the best case when $v \notin V_{M_i}$. Therefore, the minimum overlap threshold of a parent CRM c_p of the CRM $c = \{N_c, \langle M_1, M_2, \dots, M_m \rangle\}$ is defined as:

$$\rho(c_p) = \frac{\min_{1 \leq i \leq m} \{|V_{M_i}|\}}{|N_c| - 1}. \quad (4)$$

In Figure 6, we illustrate the minimum overlap calculation during search space pruning. Assume the user specified $\beta = 60\%$. We start with the anchor at step 1 when the minimum overlap threshold is 100%. Next the edge (v_1, v_2) is added for motif M_0 and M_2 and the minimum overlap based on motif M_1 is $(2/3) = 66\%$. Since the β threshold is met, we continue the enumeration process. At step 3, an edge (v_2, v_3) is added to motif M_2 . Since motif M_1 contains the lowest number of nodes, the minimum overlap is $(2/4) = 50\%$, which does not meet the β threshold. At this point, we check whether any of the parent CRM meets the β threshold and the minimum threshold for it's parent is $(2/(4 - 1)) = 66\% > \beta$. Thus, we continue the enumeration by adding an edge (v_3, v_4) to motif M_2 at step 4. The minimum overlap is $(2/5) = 40\%$ and it's parent overlap threshold is $(2/(5 - 1)) = 50\%$. Both thresholds are lower than β , hence we stop enumerating this CRM any further.

6 Experimental Design

All experiments are conducted on a 64-bit Linux server with 8-core AMD[©] Opteron[©] 6220 processor at 3.00GHz and 192GB of RAM.

6.1 Datasets

We have used two different types of datasets to evaluate CRMminer. The DBLP co-authorship network is a real world dynamic network that captures yearly co-authorship relations. The bioprocess network (GT) and the sales network (Sales) datasets are based on multivariate time-series data. To characterize the relations among different variables and understand changes over time, we represent the time-series data as a dynamic network. The CRMs discovered from these networks can be used to characterize the overall network, as shown in Section 7.2.3.

Table 1: Dynamic Network Datasets.

Dataset	#Vertices	#Edges	Span	Avg. #Edges	Avg. Degree
DBLP	1,057,524	4,841,084	55	88,020	0.08
GT	3458	83,454	47	1776	0.51
Sales	2697	138,044	66	2092	0.78

#Vertices denotes the total number of vertices in the dynamic network. #Edges denotes the total number of edges in the dynamic network. Span denotes the total number of snapshots in the dynamic network. Avg. #Edges denotes the average number of edges per snapshot in the dynamic network. Avg. Degree denotes the avg. number of edges per node in a snapshot.

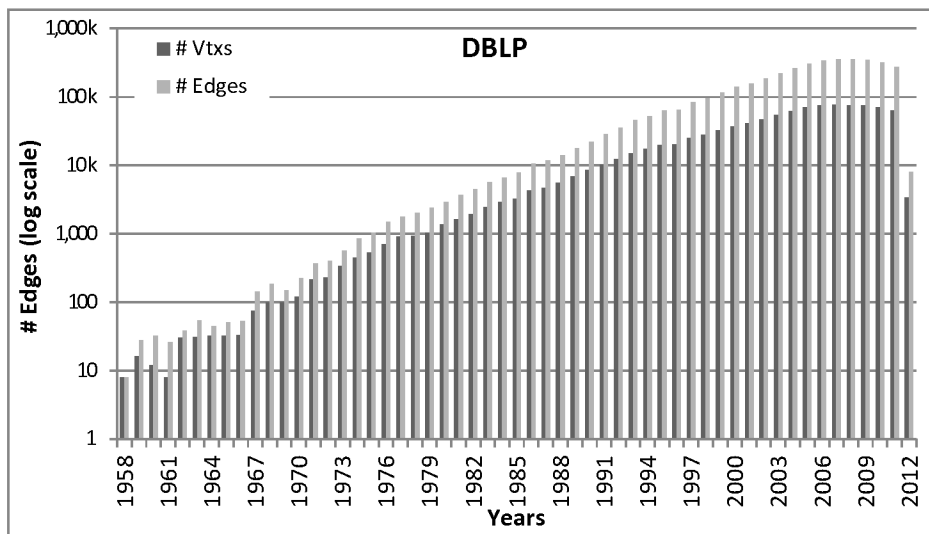


Figure 7: DBLP Dataset.

6.1.1 Co-Authorship Network (DBLP)

This is a co-authorship network that models the yearly co-authorship relations from 1958 to 2012 based on the DBLP Computer Science Bibliography Data [30]. The snapshots of the dynamic network that we created corresponds to the co-authorship network of each year, leading to a dynamic network consisting of $2012 - 1958 = 55$ snapshots. The nodes model the authors of the publications from various conferences, journals, books, lecture series, etc. and the undirected edges between the authors model the collaboration between two scientists at a certain year. If an author A co-authors a publication with an author B in a given year, the edge (A, B) is added to that year’s co-authorship graph. Multiple publications by the same authors in a year are counted as a single relation.

The edge distribution among the snapshots is skewed. In Figure 7, we show the edge distribution of the DBLP network. To determine co-authorship relations between authors with significant contribution, we removed authors who published less than 5 different years. We also removed all the relations of an author if he/she had co-authored with more than 50 other authors in a single year. To assign edge labels, we used the CLUTO software¹ to cluster the publication titles into 50 groups. These title groups can be thought as the research areas or publication topics that the authors focused on their collaboration.

6.1.2 Bioprocess Network (GT)

This is a cell culture bioprocess data [29] captured from the production bioreactors at Genentech’s manufacturing facility. This multivariate time-series data tracks the dynamics of various process parameters at every minute over 11 days period for a 247 production runs. In Table 2, we show the parameters that are used to construct the GT network. Note that some of the related variables are assigned the same vertex label. To

¹<http://www.cs.umn.edu/~cluto>

Table 2: Bioprocess Network Dataset

Vertex label	Variable(s)	Description
ASX	Air sparge rate, Air sparge total	
BT	Base totalized	
CO2	CO ₂ sparge rate, CO ₂ sparge total	
DOX	Dissolved oxygen controller output, Dissolved oxygen primary	
FRO	Flowrate overlay	
O2X	O ₂ sparge rate, O ₂ sparge total	
PHX	pH Controller output, pH Online	
ST	Sparge total	
WLC	Weight load cell	

represent this data as a dynamic network, we computed correlation matrices for 14 of the process parameters using a sliding window of 12 hours interval with 50% overlap. This process resulted in 47 correlation matrices per production run. To construct a network snapshot based on a correlation matrix, we use each parameter as a vertex and two parameters/vertices are connected with an edge labeled as positive/negative if their correlation is above +0.9 or below -0.8 threshold. These threshold values are chosen to select equal number of positive and negative labeled edges. This way we construct 47 network snapshots where each snapshot contains a 3458 vertices (14 parameter * 247 runs) and the edges between them.

6.1.3 Store Transaction Network (Sales)

This dynamic network is constructed using Dominicks Finer Foods store sales data collected from the James M. Kilts Center, University of Chicago Booth School of Business². The data captures weekly store-level sales from 93 stores collected over a period of more than seven years (400 weeks). We considered the sales data as multivariate time-series data for each store where 29 variables are the different categories of the product sold. In Table 3, we show the different categories of the product considered in the Sales network. Note that some of the related variables are assigned the same vertex label. To represent this data as a dynamic network, we considered the total number of items sold per category in a week at a store and computed correlation matrices between 29 categories using a sliding window of 12 weeks interval with 50% overlap. This process resulted in 66 correlation matrices per store. To construct a network, we use each category as a vertex and two vertices are connected with an edge labeled as positive or negative if their correlation is above +0.85 or below -0.4 respectively. These threshold values are chosen to select equal number of positive and negative labeled edges. This way we construct 66 network snapshots where each snapshot contains a 2697 vertices (29 parameters * 93 stores) and the edges between them.

6.2 Metrics

In order to assess the scalability and performance of the algorithm, we collect two sets of results for the discovered CRMs. The first set of results is collected by running CRMminer, that does not perform any overlap based pruning during CRM enumeration. This process generates the complete set of CRMs (T_{CRM}) for the specified support threshold and then applies the overlap threshold to identify the valid CRMs (Q_{CRM}) from the output space. The second set of results are collected by running CRMminer_x, that applies the overlap threshold to prune the search space during enumeration. This process uses Equation (4) to perform the overlap threshold check for search space pruning to collect all CRMs (T_{CRM}) and then performs a final check using Equation (5) to select the valid CRMs (Q_{CRM}). As discussed in Section 5.4.2, overlap based pruning does not guarantee complete set of results. Thus, the set of CRMs in T_{CRM} and Q_{CRM} collected using the CRMminer_x are subsets of T_{CRM} and Q_{CRM} collected by using the CRMminer correspondingly.

7 Results

The evaluation consists of two parts. The first focuses on assessing the performance of the algorithm that we developed for finding CRMs and to assess how the different parameters associated with the definition of

²<http://research.chicagobooth.edu/kilts/marketing-databases/dominicks/dataset>

Table 3: Sales Dataset

Vertex label	Variable(s)	Description
ANA		Analgesics
BER		Beer
CEX		Cereals, Oatmeal
CHE		Cheeses
CIG		Cigarettes
CRX		Cookies, Crackers, Snack Crackers
CSO		Canned Soup
DID		Dish Detergents
DTX		Laundry Detergents, Fabric Softeners
FEC		Front-end-candies
FRX		Frozen Dinners, Entree, Juices
GRO		Grooming Products
JUX		Bottled Juices, Refrigerated Juices
PTW		Paper Towels
SDR		Soft Drinks
SPX		Bath Soap, Soaps, Shampoos
TEX		Toothbrushes, Toothpastes
TNA		Canned Tuna
TTI		Bathroom Tissues

CRMs impacts the performance. The second focuses on assessing the information that can be extracted from the discovered CRMs by analyzing some of the patterns of coevolving relational motifs that were identified in the three datasets.

7.1 Performance Results

7.1.1 Minimum Support & Overlap

In Table 4, we present the performance of the algorithm for different values of the minimum support and overlap threshold using three different datasets. For all datasets, as the support threshold decreases, the number of discovered CRMs ($\#Q_{\text{CRM}}$) increases, consequently the runtime increases to discover those CRMs. The increase in the number of CRMs is expected due to the anti-monotonic property of the support threshold. Similar to the support threshold, we observed that the number of discovered CRMs increases as the overlap threshold decreases. This is expected as the lower overlap requirement allows a greater number of different nodes in the CRM. Based on the results from all three datasets, the runtime increase between experiments seems to be proportional to the rate of increase in the size of the output space. The runtime increase is slightly higher than linear compared to the increase in the size of output space. We noticed that the runtime is dependent on the size (i.e. the number of edges) of the discovered CRMs and the number of embeddings supporting each CRM. As the support decreases, CRMminer discovers CRMs with larger size and also the number of embeddings increases. For the *DBLP* dataset, as the support threshold decreases from 120 to 100, the total number of CRMs ($\#T_{\text{CRM}}$) discovered by both CRMminer and CRMminer_x increases by 3 – 4 times for different values of the overlap threshold when the runtimes increase by 3 – 7 times. For the *GT* dataset, as the support threshold decreases from 80 to 70 ($\beta = 0.80$), both the total number of valid CRMs and the runtimes seem to increase by 6 – 7 times. We also observe similar relationship between runtime and output space increase for the *Sales* dataset.

The performance gains of CRMminer_x against CRMminer is significant for higher values of the overlap threshold. This is expected, since a high overlap threshold prunes more CRMs during the enumeration phase and resulting in less time to complete. As the overlap threshold decreases, the performance gain also decreases and for a low value of the overlap threshold, the gain is minimal. For example, in case of the *Sales* dataset with $\phi = 40$, for $\beta = 0.60$, the performance gain is 10.7 times (the runtime of CRMminer is 1268.14 seconds and CRMminer_x is 118.47 seconds). However, for $\beta = 0.40$, the performance gain is reduced by 1.14 times. Overall, a high value of the support threshold along with a high value of the overlap threshold is most efficient. Such condition might limit the number of valid CRMs, but a higher value of overlap threshold even for a lower support threshold value still reduces the runtime significantly.

In Figure 8, we display the size distribution of the discovered CRMs for all datasets for different values

Table 4: Minimum Support (ϕ) & Overlap (β) study.

Data	ϕ	β	#S	CRMminer			CRMminer _x			
				# T_{CRM}	# Q_{CRM}	Time	# T_{CRM}	# Q_{CRM}	Time	
DBLP	120	0.60	28	524	157	533.01	510	157	463.78	
		0.50	"	"	510	"	524	510	528.02	
		0.40	"	"	524	524	"	524	531.23	
	110	0.60	29	1,073	286	1,440.55	966	966	996.50	
		0.50	"	"	966	"	1,073	1,073	1,364.97	
		0.40	"	"	1,073	"	1,073	1,073	1,331.50	
	100	0.60	41	2,194	466	3,712.50	1,742	466	1,569.87	
		0.50	"	"	1,742	"	2,194	1,742	3,630.97	
		0.40	"	"	2,194	"	2,194	2,194	3,655.45	
GT	90	0.80	19	105,575	210	189.07	5,407	210	19.72	
		0.70	"	"	6,981	"	47,430	6,981	104.52	
		0.60	"	"	40,798	"	52,541	34,331	107.19	
	80	0.80	25	4,806,955	1,005,576	6,082.24	911,199	592,413	1,245.47	
		0.70	"	"	1,504,588	"	3,517,043	1,491,761	4,679.79	
		0.60	"	"	3,305,433	"	3,882,708	3,144,734	4,890.63	
	70	0.80	37	38,127,755	6,103,951	40,438.68	5,878,178	3,421,267	6,428.28	
		0.70	"	"	13,000,807	"	24,016,376	12,292,965	28,743.44	
		0.60	"	"	23,668,959	"	28,983,885	22,255,246	28,183.17	
	Sales	50	0.60	143	28	2	1.38	24	2	1.29
			0.50	"	"	24	"	28	24	1.30
			0.40	"	"	28	"	28	28	1.30
45		0.60	196	3,644	116	34.31	1007	116	12.53	
		0.50	"	"	1007	"	2080	1007	21.43	
		0.40	"	"	3054	"	3543	3054	32.93	
40		0.60	286	169553	1017	1268.14	11547	1017	118.47	
		0.50	"	"	11547	"	37869	11547	346.61	
		0.40	"	"	84487	"	134159	84487	1114.35	

Run times are in seconds. ϕ denotes the minimum support. #S denotes the number of discovered anchors. # T_{CRM} denotes the total number of discovered CRMs. # Q_{CRM} denotes the number of CRMs that meet the β threshold out of # T_{CRM} . Time denotes the amount of time spent in seconds expanding the anchors to discover Q_{CRM} . For all datasets, the k_{\min} is 4, k_{\max} is 10, and m_{\min} is 4.

of the support threshold (ϕ) and the overlap threshold (β). Both the *DBLP* and the *Sales* datasets display similar CRM size distribution. As the CRM size increases, the number of CRMs decreases. For the *GT* dataset, the size and the number of CRMs are not similarly correlated. For $\phi = 90$ and $\beta = 60$, the number of CRMs increases as the size increases from 4 to 7 edges and then decreases for 7 to 10 edges. This increase in the number of CRMs shows that some of the CRMs of larger size were not valid CRMs at the earlier stages of enumeration. However, as more edges are included, the CRMs become valid CRM as they meet all minimum threshold requirements.

Understanding the missing CRMs The results presented in Table 4 show that CRMminer_x is able to find all the valid CRMs for the *DBLP* and the *Sales* datasets in less time than CRMminer. For the *GT* dataset, even though the runtime decreased significantly, only a subset of the valid CRMs were found by CRMminer_x. As we have seen in the size distribution characteristic of the CRMs of the *GT* dataset presented in Figure 8, the CRMminer_x fails to discover all CRMs, since the CRMs that become valid at later stage will get eliminated by the overlap threshold based pruning. The lack of anti-monotonicity property for the overlap threshold based pruning strategy eliminated a large number of potential CRMs at the early stages of the enumeration process.

To understand the characteristics of the missing CRMs, we calculated the *minimum inter-motif similarity* of all CRMs based on the following equation:

$$\psi(c) = \frac{\min_{1 < i < j < m} \{|V_{M_i} \cap V_{M_j}|\}}{|N_c|}, \quad (5)$$

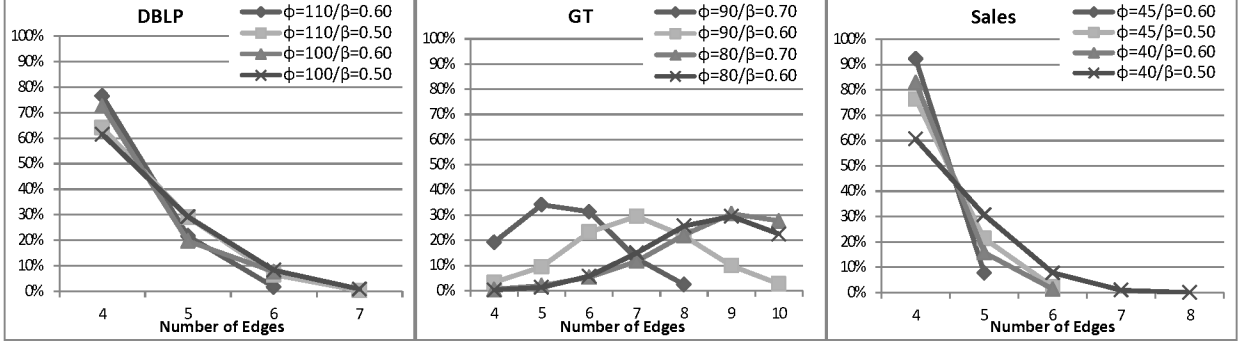


Figure 8: CRMs size distribution of different datasets for different minimum support and overlap thresholds. For all datasets, the $k_{min} = 4$, $k_{max} = 10$, and $m_{min} = 4$.

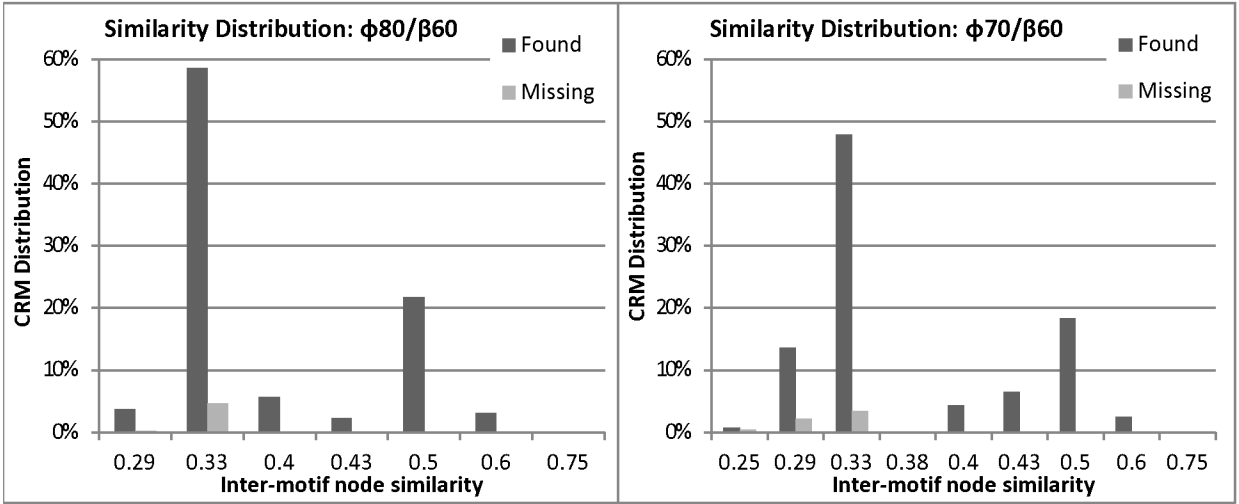


Figure 9: CRMs distribution based on the minimum inter-motif similarity using the *GT* dataset.

where $|V_{M_i} \cap V_{M_j}|$ is the number of common vertices between motif M_i and M_j and $|N_c|$ is the total number of vertices in CRM c .

We analyzed the *GT* dataset results presented in Table 4 for two different experiments ($\phi = 90, \beta = 60$ and $\phi = 80, \beta = 60$). In Figure 9, we present the minimum inter-motif similarity distribution of all identified CRMs from two separate experiments using *GT* datasets. The *found* bars represent the CRMs that were identified by CRMminer_x and the *missing* bars represent the CRMs that were missed due to overlap based pruning during enumeration step. These plots show that the minimum inter-motif similarity among the missing CRMs are low. Hence, these missing CRMs contain motifs that are mostly different from each other in terms of their nodes and resulting in CRMs that may not be capturing interesting relational changes.

7.1.2 Minimum Span

The performance of the algorithm for different values of the minimum span (m_{min}) is shown in Table 5. The value of m_{min} represents the minimum number of motifs per CRM. From the reported results in Table 5, we observe that as the value of m_{min} decreases, the number of discovered CRMs and the runtime increases. The value of m_{min} directly impacts the number of anchors and as the number of anchors increases the total number of CRMs increases.

For the *DBLP* dataset, the number of anchors increases from 20 to 575 for $m_{min} = 4$ to $m_{min} = 3$. As a result, the number of valid CRMs discovered by CRMminer increases from 91 to 68,617 and the number of valid CRMs discovered by CRMminer_x increases from 91 to 63,875. Consequently, the CRMminer runtime

Table 5: Minimum Span (m_{min}) study.

Data	m_{min}	#S	CRMminer			CRMminer _x		
			$\#T_{CRM}$	$\#Q_{CRM}$	Time	$\#T_{CRM}$	$\#Q_{CRM}$	Time
DBLP	5	0	–	–	0.04	–	–	0.04
	4	20	290	91	279.45	288	91	273.56
	3	575	177,454	68,617	163,072.77	104,998	63,875	76,834.26
GT	5	21	486,163	1,032	404.48	44,248	1,032	118.82
	4	49	9,015,625	628,657	5,770.49	2,913,699	587,035	2,853.73
	3	93	12,691,391	1,851,284	8,093.83	4,769,713	1,728,670	4,369.37
Sales	5	77	2,351	58	29.69	801	58	15.30
	4	286	53,403	1,017	229.94	11,435	1,017	103.43
	3	702	532,146	5,937	1,792.02	84,306	5,937	660.93

m_{min} denotes the minimum number of motifs per CRM, and rest of the column labels are described in Table 4. For all datasets, the k_{min} is 4 and k_{max} is 6. For searching the DBLP dataset, the β is 0.60, and ϕ is 130. For searching the GT dataset, the β is 0.80, and ϕ is 60. For searching the Sales dataset, the β is 0.60, and ϕ is 40.

and the CRMminer_x runtime increased by 585 and 230 times. Note that the algorithm did not identify any anchors for $m_{min} = 5$, thus the CRM enumeration terminated immediately. For the *GT* dataset, the increase in the number of CRMs and the runtime is less significant than *DBLP* dataset. As the number of anchors increased from 49 to 93 for $m_{min} = 4$ to $m_{min} = 3$, the total number of CRMs increased by 3 times and the runtime increased by 1.5 times. For the *Sales* dataset, we observe similar increase in both the number of CRMs and the runtime.

7.2 Qualitative Analysis

7.2.1 DBLP Case Studies

For the *DBLP* dataset, the yearly co-authorship relations among the authors are divided into 50 clusters based on the title of the papers Section 6.1. To rank the discovered CRMs, we use the cosine similarity between the centroids of the clusters, referred as *topic similarity*, that ranges from 0.02 to 0.52. For each CRM, we determine a score by calculating the average topic similarity based on all topic transitions (i.e., edge label changes) between two consecutive motifs of the CRM. The CRMs containing the least score ranks the highest. This ranking is designed to capture the frequent co-authorship relational changes that are thematically the most different. Note that the clusters are based on the publication titles and we use the most frequent words that belong to a cluster to describe the topic it represents.

Two of the high-ranked CRMs are shown in Figure 10. The first CRM shows the periodic changes in research topics represented as 43 and 26 and the topic similarity between these topics is 0.12. The CRM captures the periodic transitions of the relations as author a and b collaborate with other authors c , d , and e over the time. The second CRM shows the periodic changes in research topics represented as 2 and 10 and the topic similarity between these topics is 0.17. The CRM captures the periodic transitions of the relations as author a and b collaborate with other authors c , and d over the time.

7.2.2 Sales Case Studies

For the *Sales* dataset, we ranked the discovered CRMs according to their size (i.e., the number of edges) and larger CRMs are ranked higher. We want to capture the sales pattern where a large group of products either gain or lose their sales correlation (i.e., tendency of being sold together). In Figure 11(a), a CRM of size 8 is presented from the Sales network capturing correlated product groups at a certain period and the changes in their relations. At first, the sale of five different product groups seem to be positively correlated. Based on the relations among the nodes, the itemsets (FEC, CIG, DTX), (FEC, PTW, DTX), and (FEC, PTW, CEX) are strongly correlated. In the next motif, (CIG, CEX, DTX) items have a positively correlated sale. In the last two motifs, the relation between the items changed as (FEC, CIG) became negatively correlated. Based on the details of the different embeddings and motif occurrence period, it seems that during the holiday periods (i.e., thanksgiving, Christmas) all the product groups are positively correlated. It is interesting to

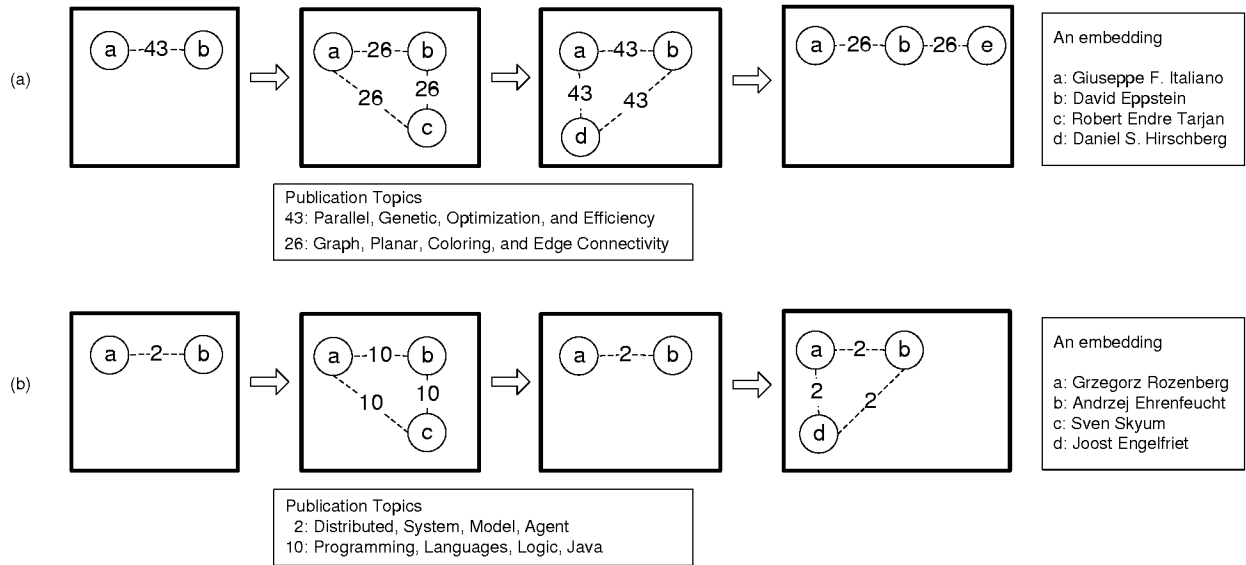


Figure 10: Two CRMs capturing co-authorship patterns. The edge labels represent the domain or subject of the publications the authors were involved together. The vertices are labeled to show the changes in relations between the nodes. These CRMs are collected using $\phi = 50$ and $\beta = 0.40$.

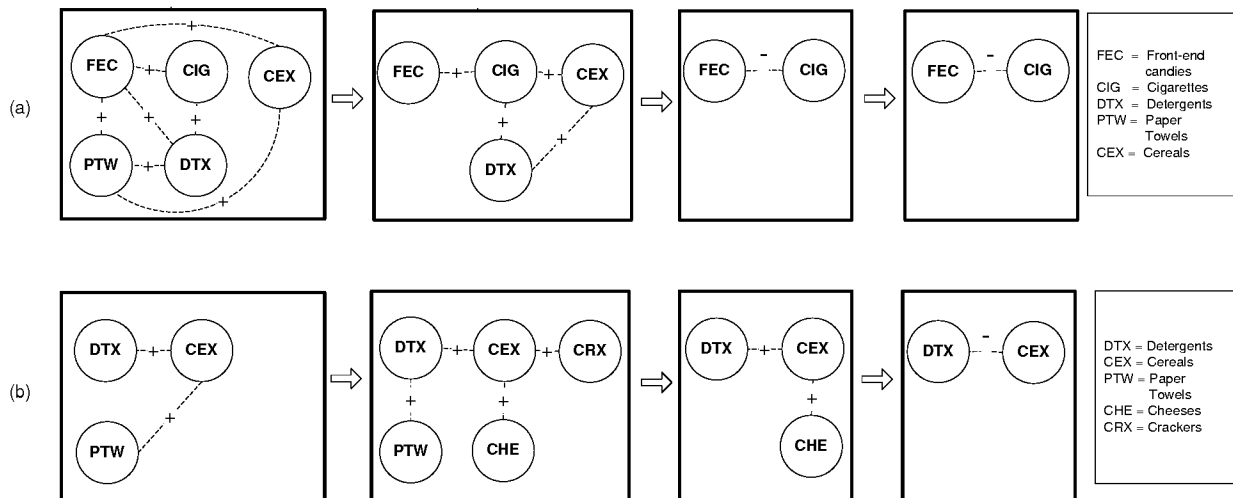


Figure 11: Two CRMs capturing store sales patterns. The edge labels + and - correspond to positive and negative correlation between entities. These CRMs are collected using $\phi = 40$ and $\beta = 0.40$.

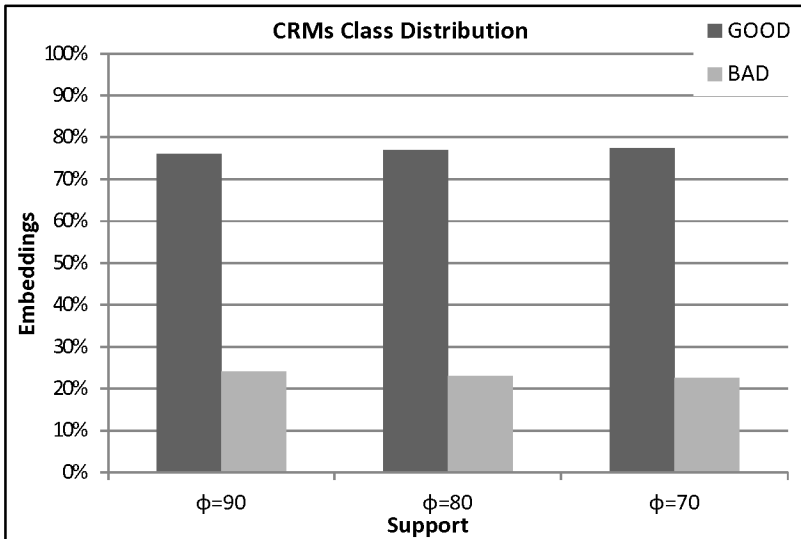


Figure 12: A distribution of the CRM embeddings. The Good class represents the production runs with high yield and the Bad class represents with poor yield. CRMs were collected using $\phi = (90, 80 \text{ and } 70)$, $\beta = 0.60$, $m_{min} = 4$, $k_{min} = 4$, and $k_{max} = 10$.

see in multiple embeddings that FEC and CIG became negatively correlated right after new year and Easter period.

In Figure 11(b), another CRM of size 5 is presented. In this case, product groups are all positively correlated at the beginning. (DTX, CEX) and (CEX, PTW) maintain positive correlation while other groups such as CEX and CHE, and CEX and CRX form positive correlation. However, at the last period, DTX and CEX product groups become negatively correlated. Note that one can focus on understanding why the relation between DTX and CEX changed and try to plan for some actions that stores can take to change the sales of those products.

7.2.3 Genentech Case Studies

For the *GT* dataset, the discovered CRMs display the dynamics of various process parameters during the cell culture process. As it is difficult to understand the relations between the nodes (i.e., parameters) without sufficient domain knowledge in Bio-Chemical processes, we decided to analyze the discovered CRMs by using them as features to discriminate the production runs. Out of the 247 production runs included in the *GT* dataset, based on the quality of the yields, 48 of the runs are labeled as *Good*, 48 of the runs are labeled as *Bad*, and the remaining were not labeled. To understand the class distribution (i.e., Good or Bad) of the discovered CRMs, we analyzed the embeddings of the discovered CRMs from three different experiments using ϕ of 90, 80 and 70. Since we have 96 labeled runs, we do not count the embeddings that belong to unlabeled runs. In Figure 12, we present the class distribution of the embeddings for the discovered CRMs. It shows that the CRMs were present mostly as part of the high yield runs, since more than 75% of the embeddings belong to the Good class. The class representation is consistent for different support parameters. These results suggest that there is a consistency of what makes some thing good but runs can go bad for many reasons. Note that using such information, if we can determine the low yield runs at the early stage of experiment, we can terminate the experiment and save a lot of resources.

In Figure 13, we present a CRM of size 10 from the *GT* network capturing the changes in entity correlation based on their recorded measurements at a particular stage of the cell culture process. Based on the embedding details, we noticed that the motif spans are mostly non overlapping. This indicates that the changes in the entity relations occur at different stages of the cell culture process based the experimental/process environment setup.

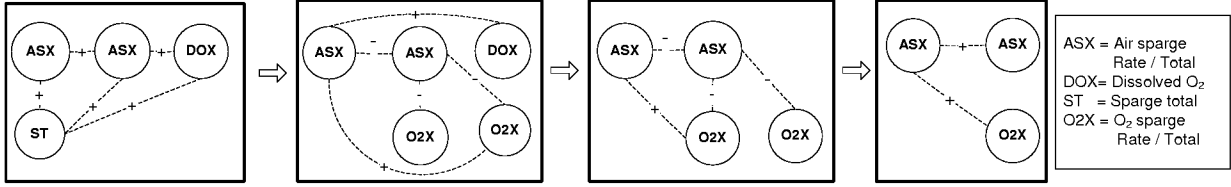


Figure 13: A CRM capturing a cell culture bioprocess pattern. The edge labels + and - correspond to positive and negative correlation between entities. This CRM is collected using $\phi = 90$ and $\beta = 0.50$.

8 Conclusion & Future directions

In this paper, we introduced coevolving relational motifs to represent patterns that change in a consistent way over time in a dynamic network and presented an algorithm to efficiently find all frequent coevolving relational motifs. The algorithm can be used to discover unknown coordination mechanisms in a system by identifying the patterns that evolve and move in a similar and highly conserved fashion in the dynamic networks. The experimental evaluation using multiple real world datasets show that the algorithm is able to discover CRMs from all datasets and is scalable to large and dense dynamic networks. Further, the qualitative analysis shows that the discovered patterns capture important information and can be used as differentiating features for other mining problems.

There are a number of ways the CRM definition can be modified to address important special classes of CRMs. First, require that the set of identified CRMs is non-redundant in the sense that no CRM is a subsequence of another CRM. This can substantially reduce the number of identified CRMs without loss of information. Second, the occurrences of the motifs are temporally synchronized, i.e., the snapshot in which each CRM's motif occurs is the same. Third, mine the CRMs that do not contain an anchor. Finally, require that the set of identified CRMs to be closed; i.e., at least one of its extensions occurs fewer times.

References

- [1] Rezwan Ahmed and George Karypis. Algorithms for mining the evolution of conserved relational states in dynamic networks. *Knowledge and Information Systems*, pages 1–28, 2012.
- [2] Tatsuya Asai, Kenji Abe, Shinji Kawasoe, Hiroki Arimura, Hiroshi Sakamoto, and Setsuo Arikawa. Efficient substructure discovery from large semi-structured data. In *Proc. of the 2nd SIAM Symposium on Data Mining*, pages 158–74, 2002.
- [3] T.Y. Berger-Wolf and J. Saia. A framework for analysis of dynamic social networks. In *ACM KDD*, pages 523–528, 2006.
- [4] M. Berlingerio, F. Bonchi, B. Bringmann, and A. Gionis. Mining graph evolution rules. *Machine Learning and Knowledge Discovery in Databases*, pages 115–130, 2009.
- [5] Karsten M. Borgwardt, Hans-Peter Kriegel, and Peter Wackersreuther. Pattern mining in frequent dynamic subgraphs. In *IEEE ICDM*, pages 818–822, 2006.
- [6] Danah Boyd and Nicole Ellison. Social network sites: Definition, history, and scholarship. *Journal of Computer-Mediated Comm.*, 13(11), October 2007.
- [7] Sergey Brin and Lawrence Page. The anatomy of a large-scale hypertextual web search engine. *Computer Networks and ISDN Systems*, 30(1-7):107–117, 1998.
- [8] Björn Bringmann and Siegfried Nijssen. What is frequent in a single graph? *Advances in Knowledge Discovery and Data Mining*, pages 858–863, 2008.

- [9] L. Cerf, T. Nguyen, and J.F. Boulicaut. Discovering relevant cross-graph cliques in dynamic networks. *Foundations of Intelligent Systems*, pages 513–522, 2009.
- [10] Deepayan Chakrabarti, Ravi Kumar, and Andrew Tomkins. Evolutionary clustering. In *ACM KDD*, pages 554–560, 2006.
- [11] W. W. Cohen. Enron email dataset. Website, 2005. <http://www.cs.cmu.edu/~enron/>.
- [12] Elise Desmier, Marc Plantevit, Céline Robardet, and Jean-François Boulicaut. Cohesive co-evolution patterns in dynamic attributed graphs. In *Discovery Science*, pages 110–124. Springer, 2012.
- [13] D. Duan, Y. Li, Y. Jin, and Z. Lu. Community mining on dynamic weighted directed graphs. In *Proceeding of the 1st ACM international workshop on Complex networks meet information & knowledge management*, pages 11–18. ACM, 2009.
- [14] Thomas L. Friedman. *The world is flat: A brief history of the twenty-first century*. Farrar, Straus & Giroux, 2005.
- [15] Jiawei Han, Jian Pei, and Yiwen Yin. Mining frequent patterns without candidate generation. In *Proc. 2000 ACM-SIGMOD Int. Conf. Management of Data (SIGMOD'00)*, pages 1–12, Dallas, TX, May 2000.
- [16] H. Hu, X. Yan, H. Yu, J. Han, and X.J. Zhou. Mining coherent dense subgraphs across massive biological networks for functional discovery. In *ISMB*, pages 213–221, Ann Arbor, MI, 2005.
- [17] Xiaohua Hu. Mining and analysing scale-free protein protein interaction network. *Int. Journal of Bioinformatics Research and Applications*, 1(1):81–101, 2005.
- [18] Jun Huan, Wei Wang, and Jan Prins. Efficient mining of frequent subgraph in the presence of isomorphism. In *IEEE ICDM*, 2003.
- [19] A. Inokuchi and T. Washio. A fast method to mine frequent subsequences from graph sequence data. In *IEEE ICDM*, pages 303–312, 2008.
- [20] A. Inokuchi and T. Washio. Mining frequent graph sequence patterns induced by vertices. In *Proc. of 10th SIAM Intl. Conf. on Data Mining*, pages 466–477, 2010.
- [21] Akihiro Inokuchi, Takashi Washio, and Hiroshi Motoda. An apriori-based algorithm for mining frequent substructures from graph data. In *Proc. of the 4th European Conf. on Principles of Data Mining and Knowledge Discovery*, pages 13–23. Springer-Verlag, 2000.
- [22] Ruoming Jin, Scott McCallen, and Eivind Almaas. Trend motif: A graph mining approach for analysis of dynamic complex networks. In *IEEE ICDM*, pages 541–546, 2007.
- [23] Yehuda Koren, Stephen C. North, and Chris Volinsky. Measuring and extracting proximity graphs in networks. *ACM Trans. Knowl. Discov. Data*, 1(3):12, 2007.
- [24] M. Koyutürk, Y. Kim, S. Subramaniam, W. Szpankowski, and A. Grama. Detecting conserved interaction patterns in biological networks. *Journal of Computational Biology*, 13(7):1299–1322, 2006.
- [25] S. Kramer, L. De Raedt, and C. Helma. Molecular feature mining in hiv data. In *ACM KDD*, 2001.
- [26] Michihiro Kuramochi and George Karypis. An efficient algorithm for discovering frequent subgraphs. *IEEE TKDE*, 16(9):1038–1051, 2004.
- [27] M. Lahiri and T.Y. Berger-Wolf. Mining periodic behavior in dynamic social networks. In *IEEE ICDM*, pages 373–382, 2008.
- [28] Steve Lawrence, C. Lee Giles, and Kurt Bollacker. Digital libraries and autonomous citation indexing. *Computer*, 32(6):67–71, 1999.

- [29] Huong Le, Santosh Kabbur, Luciano Pollarini, Ziran Sun, Keri Mills, Kevin Johnson, George Karypis, and Wei-Shou Hu. multivariate analysis of cell culture bioprocess data—lactate consumption as process indicator. *Journal of Biotechnology*, 2012.
- [30] Michael Ley. Dblp, computer science bibliography. Website, 2008. <http://www.informatik.uni-trier.de/~ley/>.
- [31] Bing Liu. *Web Data Mining: Exploring Hyperlinks, Contents, and Usage Data*. Springer, 2007.
- [32] Xiaoming Liu, Johan Bollen, Michael L. Nelson, and Herbert Van de Sompel. Co-authorship networks in the digital library research community. *Information Processing and Management*, 41(6):1462–1480, 2005.
- [33] Jian Pei, Jiawei Han, Behzad Mortazavi-Asl, Helen Ping, Qiming Chen, Umeshwar Dayal, and Mei-Chun Hsu. Prefixspan: Mining sequential patterns efficiently by prefix-projected pattern growth. In *Proceedings 2001 International Conference on Data Engineering*, 2001.
- [34] Jian Pei, Jiawei Han, Behzad Mortazavi-Asl, Helen Pinto, Qiming Chen, Umeshwar Dayal, and Meichun Hsu. Prefixspan: Mining sequential patterns by prefix-projected growth. In *ICDE*, pages 215–224, 2001.
- [35] Jian Pei, Daxin Jiang, and Aidong Zhang. On mining cross-graph quasi-cliques. In *ACM KDD*, pages 228–238, 2005.
- [36] C. Robardet. Constraint-based pattern mining in dynamic graphs. In *IEEE ICDM*, pages 950–955, 2009.
- [37] B. Schwikowski, P. Uetz, and S. Fields. A network of protein-protein interactions in yeast. *Nature Biotechnology*, 18(12):1257–1261, December 2000.
- [38] Lei Tang, Huan Liu, Jianping Zhang, and Zohreh Nazari. Community evolution in dynamic multi-mode networks. In *ACM KDD*, pages 677–685, 2008.
- [39] B. Wackersreuther, P. Wackersreuther, A. Oswald, C. B. "ohm, and K.M. Borgwardt. Frequent subgraph discovery in dynamic networks. In *Proc. of the 8th Workshop on Mining and Learning with Graphs*, pages 155–162. ACM, 2010.
- [40] Xifeng Yan and Jiawei Han. gspan: Graph-based substructure pattern mining. In *Data Mining, 2002. ICDM 2002. Proceedings. 2002 IEEE International Conference on*, pages 721–724. IEEE, 2002.
- [41] Xifeng Yan, X. Jasmine Zhou, and Jiawei Han. Mining closed relational graphs with connectivity constraints. In *ACM KDD*, pages 324–333, 2005.
- [42] C.H. You, L.B. Holder, and D.J. Cook. Learning patterns in the dynamics of biological networks. 2009.
- [43] Mohammed J. Zaki. Efficiently mining frequent trees in a forest. In *ACM KDD*, pages 71–80, 2002.