

Higher-Order Networks Representation and Learning: A Survey

Hao Tian and Reza Zafarani
Data Lab, EECS Department, Syracuse University
{haotian,reza}@data.syr.edu

ABSTRACT

Network data has become widespread, larger, and more complex over the years. Traditional network data is *dyadic*, capturing the relations among pairs of entities. With the need to model interactions among more than two entities, significant research has focused on *higher-order networks* and ways to represent, analyze, and learn from them. There are two main directions to studying higher-order networks. One direction has focused on capturing *higher-order patterns* in traditional (dyadic) graphs by changing the basic unit of study from nodes to small frequently observed subgraphs, called *motifs*. As most existing network data comes in the form of pairwise dyadic relationships, studying higher-order structures within such graphs may uncover new insights. The second direction aims to directly model higher-order interactions using new and more complex representations such as *simplicial complexes* or *hypergraphs*. Some of these models have long been proposed, but improvements in computational power and the advent of new computational techniques have increased their popularity. Our goal in this paper is to provide a succinct yet comprehensive summary of the advanced higher-order network analysis techniques. We provide a systematic review of the foundations and algorithms, along with use cases and applications of higher-order networks in various scientific domains.

1. INTRODUCTION

Networks are natural representations of relationships between entities using nodes and edges [8]. Real-world networks are observed everywhere: the structure of chemical substances, ecological systems, communication networks, air and land transportation networks, power grids, among many other examples. Many problems can be naturally described and solved using networks. For instance, epidemic models on networks [45] can help predict the spread of pandemics by analyzing the interaction network among infected individuals; link-analysis methods such as PageRank [86] can help assess the importance of Websites, which in turn can be used to fine-tune searching engine results; Shortest paths algorithms [71] calculate the most efficient driving route between two locations on the transportation networks. With the rise in demand to model systems with more complex information, researchers have enriched network models by adding additional attributes to nodes and edges: weights, signs, labels, timestamps, and even metadata. Some mod-

els have even changed or extended the network basics. An example is a *heterogeneous network* [125], where nodes and edges can be of different types. Another example is a *dynamic network* [93], where each node or edge can exist only for a specific period of time. While models for networks have been enriched from various aspects, in most network models, edges still represent *dyadic* relationships, that is, relationships among two entities.

Dyadic relationships are insufficient in many real-world scenarios, specifically when there is an interaction involving more than two entities. For example, a social event may include more than two people. This is not equivalent to social interactions among all pairs of people in the event. However, such a *higher-order* pattern frequently occurs in social networks due to *triadic closure* [41], where a triangle's formation is often dependent on three edges. Another example is the frequent appearance of some specific small subgraphs (with more than two nodes) in real-world networks [10; 62; 115]. In 2002, Shen-Orr et al. [102] introduced the term *network motifs* to represent such frequent subgraphs as the building blocks of transcriptional regulation networks. The idea was further explored for various types of graphs by Milo et al. [76], where they showed that different types of networks can be distinguished using motif counts as features [75]. Such discoveries clearly indicate that *it is insufficient to only model dyadic networks*.

As a result, modeling higher-order networks has a long history. Some higher-order representations are proposed as early as the 1960s-1970s, e.g., *simplicial complexes* [107]. However, there were many obstacles to utilizing such higher-order representations at that time. First, the computational power was insufficient to compute using such representations, even for counting simple motifs. Second, without the development of the internet, data collection was inefficient and small-scale. Hence, there was an earlier decline in the demand to model and represent higher-order networks.

With the fast development of computational resources and algorithmic tools, higher-order network analysis is now widely used across various fields leading to various discoveries. For instance, in brain networks, some motifs with high functionality are generated more than other motifs to increase neural efficiency [108; 31]; In biological networks, motifs are commonly found and capture evolution [53]; In social networks, motifs help understand group interactions [63; 47]. As most existing network data is already collected as dyadic graphs, it is often impossible to recover the original higher-order interactions (if they exist). In recent years, more higher-order network data is collected, e.g., in collabo-

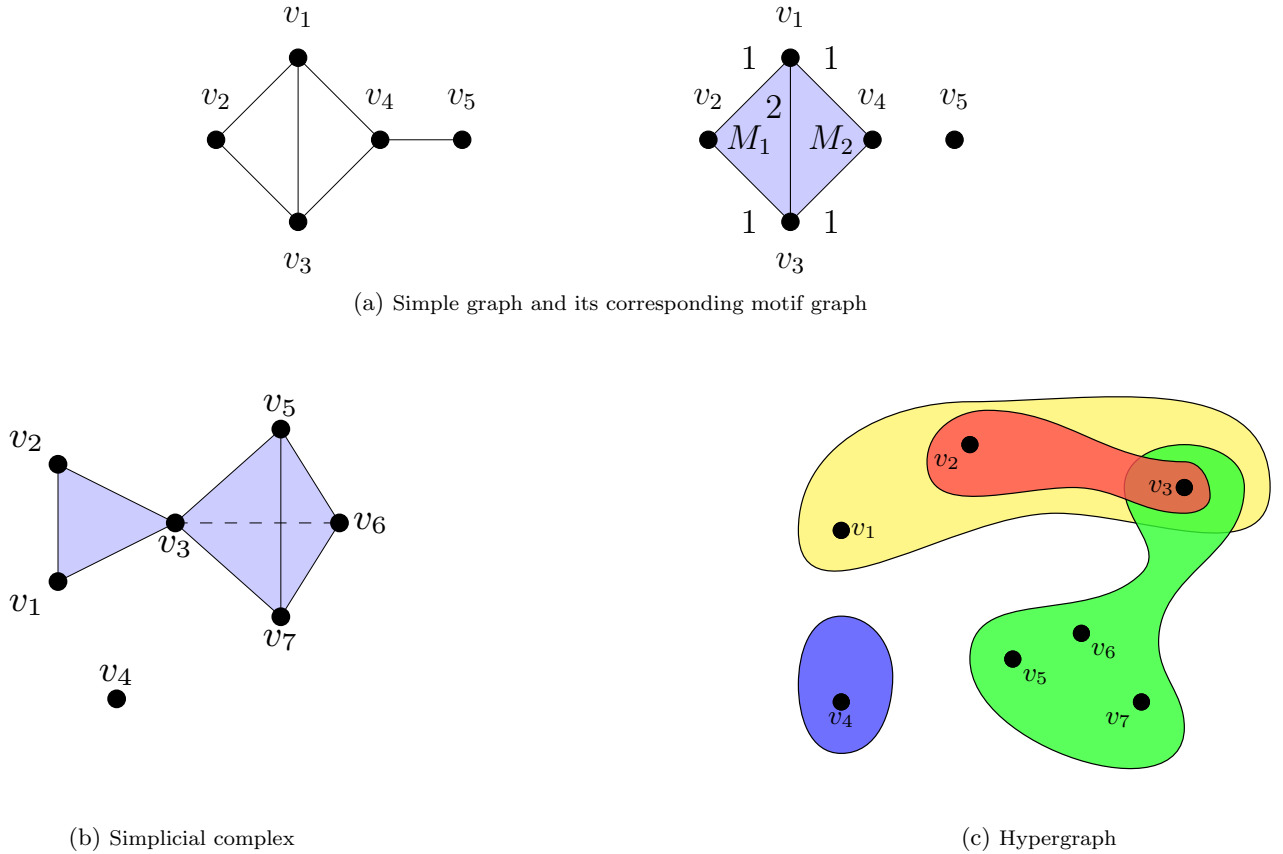


Figure 1: Comparison among Higher-Order Network Representations. (a). A simple graph (left) along with its *motif graph* (right) of 3-Cliques (triangles). With the same set of nodes, motif graphs transform edges into “membership in given motifs.” From the example, the given motif is a triangle, so there are two triangles detected and only one edge (v_1, v_3) shared by both triangles. The edge (v_4, v_4) that does not belong to any motif will be ignored. (b). A Simplicial Complex, including a 0-simplex (single node), a 2-simplex (triangle), and a 3-simplex (tetrahedron). Note that any *face* (sub-simplex) of an existing simplex is also included in the simplicial complex, for example, (v_5, v_6, v_7) . (c). A Hypergraph. Unlike simplicial complexes, any subedge of hyperedges does not have to appear in the edge set, i.e., (v_1, v_2, v_3) and (v_2, v_3) are two different hyperedges.

ration networks or on hashtags [12]. As a result, more complex higher-order algorithms and representations can now be directly used. Our goal is to broadly survey such ways to represent, analyze, and learn from higher-order networks.

Higher-Order Network Representations: there are three main branches of network abstractions that are widely studied or utilized for higher-order networks — *motifs*, *simplicial complexes*, and *hypergraphs*. Before diving into detailed mathematical representations of them, we briefly describe their differences using concrete examples involving three individuals A, B, and C [13]:

1. **Network Motifs:** A, B, and C have contact information of **each other** in their contact list. They form a triangle (an example of a *motif*).
2. **Simplicial Complexes:** A, B, and C are in the same class in high school. **Any subset of $\{A, B, C\}$** indicates a classmate relationship. A, B, and C form a *simplex*, a basic unit of a *simplicial complex*.
3. **Hypergraphs:** A, B, and C publish a paper together. We create a new *hyperedge* representing the collection

of **all authors** of this paper, i.e., hyperedge $\{A, B, C\}$. A, B, and C form a *hypergraph* with a single hyperedge. Note that $\{A, B\}, \{A, C\}$ and $\{B, C\}$ are not included as hyperedges in the hypergraph.

Figure 1 shows a comparison of these network representations. We still often study motifs in dyadic graphs, but instead of looking at pairwise relationships, we extract higher-order relationships. For example, we can identify whether an edge is part of some prespecified motif and set edge weights to the number of times the edge belongs to such motif. As shown in Figure 1 (a), one has to prespecify a motif to study, for example, a triangle. Then the original graph can be transformed by only keeping its edges that belong to such a given motif. Figure 1 (b) and (c) show a simplicial complex and a hypergraph with similar structures. Both edges (simplexes) can be represented by sets. However, in simplicial complexes, any subset simplex, by definition, also exists; however, in hypergraphs, a hyperedge only acknowledges the existence of the exact set. For example (v_1, v_2, v_3) and (v_2, v_3) exist in Figure 1 (c), while their subedges (v_1, v_2) and (v_1, v_3) do not exist.

Topics not covered in this survey. Some studies are outside of the scope of this survey and presenting them might obfuscate some of the formalism used in this survey. These studies either (1) apply a higher-order abstraction but outside the field of networks; or (2) use similar terminologies as those of higher-order networks but with a different meaning across various domains. Here, we list some such areas:

- *Network of Networks (NoN)*, or *multilayer networks* [65; 56; 23], are basically heterogeneous networks. Such models combine networks from different sources and merge into a larger, more complex network. Here, nodes and edges can be different kinds of entities. The basic unit of such networks is one type of network, but not subgraphs.
- *Higher-Order Markov Models*: These studies investigate higher-order dependencies in random walks on networks [58; 73; 123]. In some cases, first-order random walks are not able to describe network flows, so it becomes necessary to utilize higher-order Markov chains to fit the real-world observations in networks.
- *Higher-order Graph Signal Processing*: Graph signal processing [84] has also been extended from dyadic networks to simplicial complexes [96; 103; 95]. In graph signal processing, vertices carry samples of signals, and edges capture linear transformations on such signals. Such a system is denoted by a *graph filter*, which aims to model complex signal transformations based on a graph structure. For a comprehensive review of higher-order graph signal processing, refer to [98].
- *Higher-order Dynamical Systems on Networks*: A network dynamical system models pairwise node interactions based on a dynamic network structure [19]. For example, robots can form real-time shapes together by observing their neighbors. Such pairwise interactions are extended to higher-order interactions [17], using either simplicial complexes [74] or hypergraphs [80].

Scope and Organization. Compared to other surveys, our goal is to provide a succinct yet broad and comprehensive survey that focuses on higher-order networks. Abstract network topologies are categorized on the basis of the data type and applications. The surveys aim to assist researchers in identifying the appropriate methods, resources, and higher-order tools for their research tasks.

The rest of this paper is structured as follows. Section 2 introduces necessary dyadic graph basics. Sections 3, 4, and 5 introduce foundations, algorithms, and applications of network motifs, simplicial complexes, and hypergraphs, respectively. Section 6 summarizes existing datasets and tools that have been used in higher-order network studies. We present challenges and future direction and conclude in Section 7.

2. PRELIMINARY OF GRAPHS

We briefly introduce some basic dyadic graph concepts, which are being applied and generalized in various higher-order representations.

- *Undirected/Directed Graph*: An *undirected graph* is an ordered pair $G = (V, E)$. Set $V = \{v_1, v_2, \dots, v_n\}$ is the set of vertices and set $E \subseteq \{\{v_i, v_j\} | v_i, v_j \in V\}$

is the set of edges, where $\{v_i, v_j\}$ is an unordered pair of nodes. In contrast, a *directed graph* is an ordered pair $G = (V, E)$, where V is the set of vertices and $E \subseteq \{(v_i, v_j) | (v_i, v_j) \in V^2\}$ are ordered pairs of nodes.

- *Simple Graph*: A *self-loop* is an edge that starts and ends at the same vertex, for example, (v_i, v_i) . Duplicate edges in an edge set are called *multiple edges*, *multi-edges*, or *parallel edges*. A graph without any self-loop or multiple edges is a *simple graph*. Most graph-based studies focus on simple graphs.
- *Weighted Graph*: A *weighted graph* $G = (V, E, w)$ is a graph with assigned weights $w : E \rightarrow \mathbb{R}$ to its edges.
- *Graph Isomorphism*: Graphs G and H are *isomorphic*, denoted as $G \simeq H$, if there is a bijection between the vertex sets of G and H , denoted by $V(G)$ and $V(H)$,

$$f : V(G) \rightarrow V(H),$$

such that any two vertices u and v of G are adjacent in G if and only if $f(u)$ and $f(v)$ are adjacent in H .

- *Walk, Trail, Path and Cycle*: A *walk* is a sequence of vertices and edges of a graph. For example, one can traverse from one vertex to another once there is an edge between them. A walk is said to be *open* when the starting and ending nodes are different, and *closed*, otherwise. A *trail* is an open walk where no edge is repeated. A *path* is a trail in which neither a vertex nor an edge is repeated. A *cycle* is a closed walk that neither a vertex nor an edge is repeated.
- *Cut, Volume, and Conductance*: A *cut* is a partition of the vertices of a graph into two disjoint subsets, marked as (S, \bar{S}) . The *volume* of a node set $S \subseteq V$ is defined as the total number (or weight) of the edges incident with S , denoted as $\text{vol}(S)$. The *conductance* for set S , denoted as $\varphi(S)$, measures the ‘goodness’ of a cut separating a graph,

$$\varphi(S) = \frac{\sum_{i \in S, j \in \bar{S}} A_{ij}}{\min(\text{vol}(S), \text{vol}(\bar{S}))},$$

where A_{ij} are the entries of the adjacency matrix for G . Lower conductance ensures a balanced cut with fewer cross-edges (between S and \bar{S}).

- *Adjacency Matrix*: a square matrix used to represent a finite graph. Assume a graph has n nodes. Its corresponding adjacency matrix A is a matrix of size $n \times n$, where $A_{i,j} = 1$ when nodes i and j are connected and $A_{i,j} = 0$, otherwise. Adjacency matrices of undirected graphs are symmetric. Adjacency matrices of simple graphs are binary, with all zeros on the main diagonal.
- *Laplacian Matrix*: Given a simple graph G , the *Laplacian Matrix* of G is defined as $L = D - A$, where D is the diagonal degree matrix and A is the adjacency matrix of G . If a graph G is undirected, its *Laplacian Matrix* is also a symmetric matrix. It can be normalized to matrix of unit vectors, usually denoted as $\mathcal{L} = D^{-1/2} L D^{1/2} = I - D^{-1/2} A D^{1/2}$.

3. NETWORK MOTIFS

Compared to more complex higher-order network representations, motifs have been studied for a relatively longer period. There are two main reasons: (1) motifs are studied directly on dyadic graphs, which are widely used in various research fields; (2) some specific subgraphs already have some special meanings in the real world, so it is natural to study them in networks.

3.1 Foundation and Algorithms

A network motif is generally defined as a **highly significant subpattern or subgraph in the network** [76]. The term ‘‘significant’’ indicates that the number of times the motif appears in the graph is higher than what is expected or normal, where such expected numbers are often understood within *random graphs* (e.g., the *Erdős-Rényi model* [32]). A motif can be some fixed-size subgraph such as a triangle; or can have a variable size representing some general conceptual pattern such as a star or a loop [102]. For brevity, we focus on fixed-size subgraphs. Formally, a network motif is

DEFINITION 3.1 (NETWORK MOTIF). *Motif M of graph G is a subgraph of G that has multiple isomorphic graphs that are also subgraphs of G . That is, $G_1 \subseteq G, G_2 \subseteq G, \dots, G_n \subseteq G$ and G_1, \dots, G_n are isomorphic to M .*

3.1.1 Motif Frequency

For motif M , its number of appearances in graph G can be denoted as $F_G(M)$. By comparing this frequency with the mean (expected) count in random graphs with the same size as G , we obtain the Z -score of motif frequency

$$Z(M) = \frac{F_G(M) - \mu_R(M)}{\sigma_R(M)}, \quad (1)$$

where $\mu_R(M), \sigma_R(M)$ are the expected mean and standard deviation of frequencies of M in a random graph. This Z -score is an important statistic for measuring the significance of motifs.

3.1.2 Motif Matrix and Motif Cuts

Similar to the adjacency matrix, Benson et al. [15] define *motif adjacency matrix* based on the memberships of edges in the given motifs. Formally, given a specific motif M , the motif adjacency matrix W_M is defined as

$$(W_M)_{ij} = \left| \{M | i \in M, j \in M\} \right|, \quad (2)$$

where $(W_M)_{ij}$ is the number of instances of M that contain nodes i and j . When the given motif is an edge (two connected nodes), the motif adjacency matrix is simply the adjacency matrix. Note that when the given motif is not a complete graph, the nodes i and j can belong to the same motif even if they are not connected.

The motif adjacency matrix is also of size $|V| \times |V|$, so most algorithms designed for the adjacency matrix are also suitable for the motif adjacency matrix.

With the motif adjacency matrix in place, a *motif cut* can be defined consequently for motif M and motif adjacency matrix W_M :

$$cut_M(S, \bar{S}) = \sum_{i \in S, j \in \bar{S}} W_{Mij}. \quad (3)$$

Similarly, *motif conductance* is defined as

$$\varphi_M(S) = \frac{cut_M(S, \bar{S})}{\min(\text{vol}_M(S), \text{vol}_M(\bar{S}))}, \quad (4)$$

where $\text{vol}_M(S)$ denotes the volume (total sum of edge weights) of set S in the motif matrix. If a clustering algorithm minimizes motif conductance, the result leads to preserving the structure of the given motif in the graph while generating a balanced split by the clustering algorithm.

3.1.3 Motif Clustering Coefficient

The *clustering coefficient* is a measure of the degree to which nodes in a graph tend to cluster together [121]. In dyadic graphs, the clustering coefficient can be calculated by the fraction of length-2 paths (wedges) that are involved in triangles. From this perspective, the *global clustering coefficient* can be defined as

$$C = \frac{6|K_3|}{|W|}, \quad (5)$$

where $|K_3|$ is the number of triangles (3-cliques), and $|W|$ is the number of wedges. Each triangle is counted six times since it contains six different wedges, considering the order. The *local clustering coefficient* of node u is defined as

$$C(u) = \frac{2|K_3(u)|}{|W(u)|}, \quad (6)$$

where $C(u)$ is the fraction of triangles that node u belongs to over the number of wedges in which u is the center node.

Based on the above interpretation of the clustering coefficients, Yin et al. [126] introduce a generalized higher-order clustering coefficient for motifs of higher-order cliques. For order $l \geq 2$, clustering coefficient of order l can be defined as

$$C_l = \frac{(l^2 + l)|K_{l+1}|}{|W_l|}, \quad (7)$$

where K and W are higher-order cliques and wedges, and $(l^2 + l)$ is basically $(l + 1)l$, which is the number of wedges that an $(l + 1)$ -clique closes. Consequently, the local higher-order clustering coefficient is generalized as

$$C_l(u) = \frac{l|K_{l+1}(u)|}{|W_l(u)|}. \quad (8)$$

3.2 Use Cases and Applications

Motifs are utilized across many scientific fields. Here, we survey use cases and applications of motifs.

3.2.1 Capturing Functionalities in Networks

Highly frequent motifs are often related to specific functionalities that networks capture, especially in biology. We provide some examples in biological and brain networks.

Biological networks. In 2002, Shen-Orr et al. [102] distinguished three families of motifs in *Escherichia coli* (E. coli) directed transcriptional network. These families are closely related to specific functionalities. They are: feedforward loop (a directed acyclic graph), single input module (one to many transactions) and dense overlapping regulons (many to many transactions). Each motif relates to a specific function in the determination of gene expression. Such frequent motifs can be detected using a brute-force approach on some sub-matrix of the adjacency matrix.

The functionality of motifs in biological networks is further validated through a simulation of *spontaneous evolution process* [53]. First, an electronic combinatorial logic circuit is initiated by random wiring. The goal of network evolution is to increase the fraction of correct output under given logical functions. The baseline, a *fixed goal* given by $G_1 = (X \oplus Y)AND(Z \oplus W)$, is very slow to converge with a low rate of reaching the perfect solution. Networks that evolved under fixed goal have less significant motifs and lower modularity (the separability of the design into units that perform independently). Addressing this issue, the authors introduce *modularly varying goals*, where the goals switch between G_1 and $G_2 = (X \oplus Y)OR(Z \oplus W)$ every 20 epochs. Surprisingly, such evolved networks could always find perfect solutions of the current goal (either G_1 and G_2) quickly within a few epochs. Similar findings are discovered in neural networks, which explain adaptiveness and robustness of motifs in real-world biological networks.

Brain Networks. Due to the complexity of brain networks, relationships between subareas of the brain and their functionality are of significant research interest. In brain networks, motifs are often divided into two groups: *functional* and *structural*. Structural motifs are those currently presented in this survey and capture the anatomical building blocks of the brain network, whereas functional motifs capture patterns of elementary processing within such structural motifs. In other words, functional motifs can be considered as all possible subgraphs of the structural motif with the same number of nodes but different edges. For example, a triangle structural motif consists of three different functional motifs (all paths of length two). One hypothesis suggests that the number and variety of *functional motifs* are maximized in the brain to increase effectiveness [108].

Functional motifs vary significantly over time. Duclos et al. [31] investigate motif appearances in the brain network over time. They count all connected subgraphs of size three in directed brain networks, and as a result, show that anesthetic-induced unconsciousness is associated with a topological re-organization of the brain network. Specifically, the frequency of chain-like and loop-like motifs change significantly when people transition from a responsive state to an unresponsive state. Such observations demonstrate links between motifs and functionalities.

In terms of the shape of motifs in brain networks, research has been more interested in specific motifs that are easier to count, such as cliques and *cavities* (enclosed spaces). For instance, all maximal cliques are counted, and their frequency is compared to that of what is expected in some null model. The results indicate the spatial distributions of maximal cliques are more than expected in different brain regions. Similarly, cavities can be studied (ranging from minimum cycles to incomplete cliques). Research shows that, unlike cliques, cavities are less than expected in different areas of the brain [105].

3.2.2 Network Classification

In the seminal work of Milo et al. [76], network motifs are defined as specific fixed subgraphs whose frequencies are higher than what one would expect in random graphs. It turns out motifs found and their frequencies from various types of networks (including food webs, biological networks, electronic circuits, World Wide Web, and the like) can be utilized to classify types of networks. In particular,

graphs from the same category exhibit significant overlap in terms of motifs observed and their frequencies, which can be used as features to distinguish graphs. For example, on food webs, a three-node chain is frequently observed. However, this motif is not frequently observed in any other category of networks. The feed-forward loop is popular on most information-processing networks, including brain networks, electronic circuits, and the World Wide Web.

Milo et al. further investigate the number of motifs across various categories [75]. Z -scores (see Equation 1) are calculated for subgraph frequencies and are compared with those expected in random graphs. The results show that the graphs from the same category have highly similar subgraph frequencies. However, some graphs from different categories also show similarities in their frequencies, which captures the intrinsic similarity between graphs that are from similar categories. As a result of this discovery, networks from different categories can be classified into superfamilies using motif frequencies.

3.2.3 Network Models

As motif frequencies are different in real-world graphs compared to what one would expect in random graphs, there is a need for *network models* that can generate realistic random graphs with motif frequencies similar to those of real-world graphs.

Pržulj et al. [91] propose a new series of network models called geometric random graphs, which uniformly generate nodes (i.e., points) in 2D/3D/4D Euclidean space and form links between nodes based on a threshold on their distances. By counting all possible motifs under size five, they found that the random graphs generated by geometric models are more similar in motif counts to the original protein-protein networks than those generated by other network models.

There is also a significant need to develop non-random network models that can generate motifs with similar frequencies to those observed in real-world data. One example is the network model designed by Leskovec et al. [63]. The work examines all triangles in *signed* networks, where edges can have a sign: + or -. For instance, a + edge may indicate that two nodes are “friends,” and a - edge may indicate that two nodes are “enemies.” They discovered that network models that simulate the *balance theory* (colloquially stated as “an enemy of an enemy is my friend”) might be insufficient to explain the frequency of triangles appearing in real-world networks. Therefore, they propose another network model for directed links, inspired by *status theory*, where a positive link from node a to b indicates that a has a higher “status” than b . Given a three-node directed cycle with two positive signs, these two theories will predict opposite signs for the third link. Balance theory explains this as three pairwise friends (friend of a friend is a friend), while status theory considers $a \rightarrow b \rightarrow c$ as a pattern of increasing status, so c should link to a with a negative sign (high to low status). In directed graphs, research shows that graphs generated based on status theory could more realistically replicate the frequencies of signed motifs compared to graphs generated based on balance theory.

3.2.4 Clustering

Motifs are closely related to clusters in higher-order networks. Benson et al. [15] develop a framework of network clustering based on higher-order properties. More specifi-

cally, cuts on edges are generalized to cuts on motifs, and the adjacency matrix is generalized to the motif membership matrix. Their results show that such clustering accurately preserves higher-order structures. Two real-world examples are presented. For clustering based on a *bi-fan* motif, the clustering clearly distinguishes between the role of source and sink by assigning them to different clusters in neuronal networks. In the airline network, transportation hubs are clustered together by using a bi-directional 2-path motif.

Building upon the ideas of motif matrix and motif cuts, Yin et al. [127] generalize clustering methods to the motif level. They first propose a motif-based approximate personalized PageRank (MAPPR), which performs an approximation of the Personalized PageRank using the motif matrix. The method can quickly find a cluster that contains a given seed node that has the minimum *motif conductance*. To enhance the performance in case the clustering is performed on the whole graph, they introduce an efficient method to identify good seed nodes to be used as input to MAPPR. The proposed method is validated by performing cluster recovery tasks on synthetic and real-world graphs. Experimental results show that the proposed techniques could preserve higher-order clustering coefficient (as detailed in Section 3.1.3).

Furthermore, as we also showed in Section 3.1.3, some traditional measurements of clusters (or clusterability) are generalized to the motif-level. One important generalized graph measurement was the *clustering coefficient*, which reflects the degree of cohesiveness of communities. Yin et al. [126] introduced higher-order variants of the local and global clustering coefficients. For order-3, the global clustering coefficient is defined as the ratio of cliques to wedges (length-2 paths); the local clustering coefficient is defined as the ratio of cliques that a node involved over wedges. Interested readers can refer to Section 3.1.3 for extensions of clustering coefficients to orders greater than three.

Another important measure for a clustering is its *modularity*. Modularity is a quantitative measure to evaluate the significance of clusters, which is also generalized to the motif level [56], specifically two special motifs – cycles and paths. In the general case, motif modularity is defined as the fraction of motifs laying fully inside the community subtracted by that of expected in the random graphs. Higher-order modularity can distinguish differences in higher-order structures, such as cycles and cliques/hubs and leaves. For example, in a multipartite network roles can be easily distinguished simply by applying higher-order modularity.

3.2.5 Representation Learning

Representation Learning aims at encoding specific network properties into fixed-length vectors. In order to capture higher-order structures, Rossi et al. [94] propose a network embedding method based on motifs. First, they build several weighted motif adjacency matrices based on the nodes' occurrences in specific motifs. Then they define a series of functions over these weighted motif matrices, such as k -step paths, the transition matrix, and various Laplacians. By minimizing the distances between the motif-based matrix formulation and the embedding matrix, each motif matrix learns a local embedding. Finally, they concatenate the local embedding to calculate a global embedding for the network. Experimental results show that the proposed higher-order network embeddings outperform other embedding methods

in link prediction tasks.

Another higher-order representation learning method using motifs is *LEMOM* [101]. First, *LEMOM* converts the graph by adding *supervertices* for motifs (e.g., triangles), and then links the nodes that are involved in such motifs to the corresponding supervertices. The result is a *two-mode* network that captures the memberships of nodes in motifs, where the edges between motif supervertices to regular nodes are defined as structural edges. The embedding vector is learned through a random walk process that captures the similarities of nodes that share similar motif structures. A parameter q controls the traversal probability from a regular node to supervertices. With larger q , any node will become closer to nodes similar in motif properties rather than its structural neighbors. *LEMOM* has been successfully applied in anomaly detection, link prediction, and node classification.

3.2.6 Link Prediction

Motif counts can be used a powerful feature to predict missing links. Abuoda et al. [1] convert the link prediction problem into a classification problem by counting the motifs involved in the link candidates. All possible connecting motifs within size five are enumerated as features. The performance of several classical classifiers shows that larger motifs can lead to higher performance and that a combination of motifs can further improve the results. The work shows that motif-combined feature classification outperforms most state-of-the-art link prediction methods.

4. SIMPLICIAL COMPLEXES

A simplicial complex can be interpreted as another generalization of a graph. In graphs, there are two different types of entities – nodes and edges. But in simplicial complexes, the concepts of nodes and edges are merged into a generalized basic unit, the *simplex*, where 0-simplex represents the single vertex and 1-simplex represents the edge. Furthermore, a simplicial complex can contain any *order* of interactions, such as k -simplices ($k \geq 0$).

The main difference between simplicial complexes and hypergraphs is the requirement of being *inclusive*; that is, a simplicial complex also contains all *faces* (sub-simplices) of its current simplices. For example, if three people A , B , and C belong to the same university, then all pairs: AB , AC , and BC have the same relationship (being part of the university). When modeling higher-order data with simplicial complexes, ensuring the relationship being modeled is inclusive is the first requirement.

4.1 Foundation and Algorithms

In mathematics, a simplicial complex is a set composed of points, line segments, triangles, and their n -dimensional counterparts.

4.1.1 Simplex

A *simplex* is the basic unit of a simplicial complex and is the generalization of the notion of a triangle or a tetrahedron to higher dimensions. More specifically, a k -*simplex* is a k -dimensional polytope that is the convex hull of its $k + 1$ vertices.

The convex hull of any nonempty subset of the $k + 1$ points that define a k -simplex is called a *face* of the simplex.

Faces are also simplices. Any $k - 1$ -face of a k -simplex is called a *facet*.

4.1.2 Simplicial Complex

DEFINITION 4.1 (SIMPLICIAL COMPLEX). *A simplicial complex \mathcal{X} is a set of simplices that satisfies the following:*

1. *Every face of a simplex from \mathcal{X} is also in \mathcal{X} ; and*
2. *The non-empty intersection of any two simplices $\sigma_1, \sigma_2 \in \mathcal{X}$ is a face of both σ_1 and σ_2 .*

Roughly speaking, simplicial complexes are simplices that are (1) closed under taking faces and (2) have no inner intersections other than faces. A *simplicial k -complex \mathcal{X}* is a simplicial complex where the largest dimension of any simplex in \mathcal{X} equals k .

4.1.3 Homology

First, we define the *orientation* of a simplex. The orientation of a k -simplex is given by an ordering of the vertices (v_0, v_1, \dots, v_k) . There are exactly two orientations – even and odd permutations, and switching any two vertices in the ordering leads to a change of the orientation. For example, in a two-dimensional space, we have clockwise and counterclockwise ordering for a triangle. The orders $(v_1, v_2, v_3), (v_2, v_3, v_1), (v_3, v_1, v_2)$ indicate one orientation, and the orders $(v_1, v_3, v_2), (v_3, v_2, v_1), (v_2, v_1, v_3)$ indicate the opposite one.

Let \mathcal{X} be a simplicial complex. A *simplicial k -chain* is a finite formal sum

$$\sum_{i=1}^N c_i \sigma_i, \quad (9)$$

where c_i is an integer and σ_i is an oriented simplex. For each simplex, the sum includes a sign based on the orientation. One way of assigning orientations is to order all vertices of the simplicial complex and give each simplex the orientation corresponding to it. The group of k -chains on \mathcal{X} is written $C_k(\mathcal{X})$, and for simplicity we write C_k . Note that C_k is a vector space with the number of k -simplices as its dimension.

Based on the group of k -chains C_k , we define *boundaries* and *cycles*. First, we define the boundary operator.

DEFINITION 4.2 (BOUNDARY OPERATOR). *Let $\sigma = (v_0, \dots, v_k)$ be an oriented k -simplex, viewed as a basis element of C_k . The boundary operator $\partial_k : C_k \rightarrow C_{k-1}$ is the homomorphism defined by:*

$$\partial_k(\sigma) = \sum_{i=0}^k (-1)^i (v_0, \dots, \widehat{v}_i, \dots, v_k),$$

where $(v_0, \dots, \widehat{v}_i, \dots, v_k)$ is the i^{th} face of σ , which deletes v_i from σ .

The boundary of each k -simplex is the collection of all its $(k - 1)$ -faces. In C_k , elements of the subgroup $Z_k := \ker \partial_k$ are referred to as *cycles*, which is the collection of k -simplices whose boundary is zero. While subgroup $B_k := \text{Im } \partial_{k+1}$ denotes the boundaries, i.e., boundaries of $(k + 1)$ -simplices. Note that using definition 4.2, it is easy to prove that the boundary of boundaries is empty.

Cycles are essential entities for detecting holes. However, some simplices in Z_k are just boundaries of $(k + 1)$ -simplices, which are not holes themselves. Hence, we remove the boundaries of the $(k + 1)$ -simplices from the cycles. This can be defined as the *quotient abelian group*

$$H_k = Z_k / B_k = \ker \partial_k / \text{Im } \partial_{k+1}, \quad (10)$$

where the remaining simplices in H_k represent k -dimensional holes in the complex. H_k is called the *homology group*.

4.1.4 Cohomology and Hodge Laplacian

Remember the group of k -chain C_k is a vector space over \mathbb{R} . Hence, it is possible to give an inner product structure to each C_k to make the basis (oriented simplices) orthogonal. We denote this dual space of C_k as C^k [79], called the group of *k -cochains*.

We denote the dual operator of the boundary map ∂_k as δ_k . Operator $\delta_k : C^k \rightarrow C^{k+1}$ is called the *co-boundary operator*, which is the adjoint of boundary map ∂_k . Consequently, the *cohomology group* is defined over *cochains*

$$H^k = \ker \delta_k / \text{Im } \delta_{k-1}, \quad (11)$$

which is exactly a dual group of the homology group H_k . Note that the cohomology groups are defined more algebraically with less geometric meaning. The main purpose of introducing the cohomology group is to derive the *Hodge Laplacian* (see Definition 4.3). For a detailed explanation, interested readers can refer to [66].

Given a simplicial complex \mathcal{X} , its boundary map ∂_k can be represented as a matrix B_k . B_k has the dimension $n_{k-1} \times n_k$, where n_{k-1} and n_k are the number of $(k - 1)$ -simplices and k -simplices, respectively. For example, $B_0 = 0$ and B_1 is a matrix of dimension $|V| \times |E|$.

Similarly, the co-boundary map δ_k can also be represented as the adjoint matrix B_k^* . In a finite real space, it is equal to the transpose of B_k , so we can also write it as B_k^T .

DEFINITION 4.3 (HODGE LAPLACIAN). *The k th Hodge Laplacian of a simplicial complex \mathcal{X} is defined as*

$$\mathcal{L}_k = B_k^T B_k + B_{k+1} B_{k+1}^T. \quad (12)$$

When $k = 0$, $\mathcal{L}_0 = B_1 B_1^T$ is exactly the Laplacian matrix in dyadic graphs, with dimension $|V| \times |V|$. Matrix \mathcal{L}_1 has dimension $|E| \times |E|$, capturing relationships among basic units of edges [97].

4.1.5 Degrees and Random Simplicial Complex

Degree is generalized in simplicial complex as follows [27]:

DEFINITION 4.4 (DEGREE OF A SIMPLEX). *For any simplex $\sigma \in \mathcal{X}$, the degree $k_{d,\lambda}(\sigma)$ is the number of d -dimensional simplices adjacent with σ in λ -faces.*

When we are only interested in the degree of vertices, we let $\lambda = 0$. Then $k_d(v)$ becomes the number of d -simplex incident to v , i.e., those that v belongs to.

As an analog to the *Erdős–Rényi* model [32] in dyadic graphs, the generative model of 2-complexes can be defined as follows:

DEFINITION 4.5 (RANDOM 2-COMPLEX). *The $\mathcal{X}(n, p)$ model of a simplicial complex is defined to have vertex set $[n]$, edge set $\binom{[n]}{2}$, and each of the $\binom{[n]}{3}$ possible triangle faces is included independently with probability p .*

Note that for a 2-complex both nodes and edges (e.g. a complete graph) have to be specified, and the random process only occurs on random triangles [50]. One can define random simplicial complexes of higher order in similar ways.

4.2 Use Cases and Applications

Applications of the simplicial complex have mainly focused on two directions. One direction is focused on topological properties, where a simplicial complex is used to represent a space. In many fields, the real-world information can be abstracted to pure topology entities through a process called *filtration*, which transforms real distances into topological edges. Such techniques are widely used in sensor coverage problems, biological networks, mobility analysis, robotics, and the like. The second direction is to model real-world interactions of more than three individuals as simplices.

4.2.1 Sensor Coverage

Sensor coverage problem aims at measuring a “coverage” area and detecting locations that are uncovered: also known as *holes*. Ghrist and Muhammad [36] modeled the sensor coverage problem using simplicial complexes. A coordinate-free sensor network is formed by relative distances between any pair of nodes, without any specific coordinates. This is simpler to obtain through the strength of signals sent by the sensors, especially in dynamic systems. Based on simplicial homology theory, coverage holes are what remain after removing boundaries from cycles. The theoretical results are also verified by practical simulations in computational homology software.

The sensor cover can be further linked to the homology of the diagram of complexes [28]. In particular, the sensor cover can be defined as a collection of discs of radius r_c , and the radius of strong and weak signals of pairwise distances can be represented as r_s and r_w . *Rips complex* is defined as a simplicial complex whose simplices are tuples of nodes whose pairwise Euclidean distances are within a certain threshold. Each node can detect the existence of the boundary of the domain within another radius r_f . By forming the simplicial complexes of all these graphs, one can derive the sensor coverage.

Under similar settings with previous studies, Tahbaz-Salehi and Jadbabaie [109] present a distributed algorithm for coverage verification without any metric information. The goal of coverage verification is basically three aspects – detecting coverage holes, calculating their locations, and detecting redundancies in the network. The main novel contribution of this work is to solve the homology problem through a linear programming relaxation.

4.2.2 Disease/Abnormality Detection

Point cloud is one of the classic data formats often used in biology, where points are substances such as proteins. Similar to the sensor coverage problem, a simplicial complex can be constructed over a point cloud through the *filtration* process [81]. Points agglomerate together and become simplices when their distances fall under a specific threshold, specified by some distance function. As a result, a simplicial complex can be used for preprocessing to enhance the clustering performance [82]. One notable usage is to identify subtypes of breast cancer. Simplicial complexes can also help distinguish between recurrent and non-recurrent sub-

types [29; 6].

In brain networks, a new topology called *homological scaffold* can be defined to represent low-connection areas in the network [90]. First, a brain network can be seen as a weighted network, where larger weights indicate longer distances. Then a filtration process is applied to generate sparse structures (larger weights), followed by the detection of a homology group. The remaining structure in the homology group contains cycles with larger distances, which captures areas in the network that exhibit extremely low connections.

In neuroscience, for amplifying the differences in network topologies, a novel matrix signature is proposed to facilitate forming the homology groups [38; 37]. Instead of the absolute distances, the orders of distance are used. For example, if the distance of v_0 and v_1 is the minimum of all pairwise distances, then the entry of A_{01} will be encoded as 0. Such a non-linear transformation obscures the distances but focuses more on the intrinsic structure of the network. The order matrix captures a more robust relationship with the topological structure, for example, the number of non-contractible cycles. Experiments on pyramidal neurons in the rat hippocampus show that the proposed signature is capable of detecting geometric organization.

4.2.3 Mobility Analysis

To study mobility, topological signatures have been proposed that represent trajectories as points in k -dimensional space, where k is the number of obstacles [35]. The goal of this mapping is to characterize the differences in traces when passing by obstacles. First, obstacles are represented as simplicial complexes, and any motion toward faces can be recorded by sensors. Based on homology, these faces (edges) are encoded as real numbers. These values are added to the entries of the related obstacles as trajectories records. Then, trajectory traces can be distinguished by this signature. For example, one coming across an obstacle from the left is encoded as 1, while as -1 if coming from the right; if one loops clockwise around an obstacle, we can encode that as 2, and -2 for counterclockwise loops.

4.2.4 Network Modeling

The configuration model is a method for generating random networks from a given degree sequence. For a simplicial complex, the configuration model is also generalized along with the *canonical ensemble* [27]. In short, the canonical ensemble aims to derive the probability of the simplicial complex that maximizes the entropy defined by it. The configuration model is the uniform distribution of all possible simplicial complexes with the same degree sequence.

Based on such generalizations, Young et al. [128] further develop efficient sampling algorithms for the simplicial complex configuration model. First, they elaborate the numerical constraint of the configuration model by switching the simplicial complex to its equivalent *graphical* representation. That is, to introduce extra nodes to represent adjacent relationships between nodes and simplices. In this way, the simplicial complex is transformed into a dyadic bipartite graph, which can yield a solution [33].

The *social contagion* can be modeled as a propagation network, where people get infected through social interactions (edges). In terms of a simplicial complex, a contagion model could also consider an infection being caused by a

group, called *Simplicial Contagion Model* [48]. Similar to the dyadic contagion model, any pairwise interaction can lead to an infection with a uniform probability (β_1). In addition, higher-order interactions have unique contagion probabilities if there are multiple infecteds involved. For example, in a simplicial complex, if the other two nodes are infected, the candidate will have a probability of (β_Δ) being infected. The behavior of the infection pattern is discussed by simulating the process over both real-world and synthetic graphs. The *Simplicial Contagion Model* is a more flexible fit for more complex diseases with varying infection probabilities of different orders.

In quantum physics, Bianconi and Rahmede [16] propose a model of emergent geometry that is based on a growing simplicial complex. The model is simple as it just keeps including simplices with fixed dimension d and gluing them to the existing simplicial complex on one of its $d-1$ faces. Many advantages of such a model are validated and discussed under certain settings, including scale-free degree distribution, small-world properties, and modular structure.

4.2.5 Network Analysis Tools

Instead of studying motifs in dyadic networks, Benson et al. [12] directly collect higher-order relationships in the real world, such as co-authorships, event participation, drug instances, among other similar interactions. Such coappearances are modeled as simplicial closures (timestamped vertex sets). For example, a closed triangle indicates relationships among three nodes, while an open triangle just represents pairwise relationships between any two nodes. This representation enriches the network information and can be used in dynamic graph-evolving models or link predictions. In the link prediction task, the goal is to predict whether the open triangles will become close in the future. Results show that even simple local features such as the mean of weights on three edges perform pretty well and are comparable with state-of-the-art methods.

Based on the 1st normalized Hodge Laplacian, Schaub et al. [97] discuss random walks on basic units of edges in a simplicial complex. This work enriches the traditional field of network analysis, which is mostly node-based. Two applications are performed to verify the usage. One is representation learning of edge-flows and trajectory data, as a higher-order generalization of diffusion maps and Laplacian eigenmaps. Another is the edge-based generalization of PageRank [39], which focuses on the importance of edges rather than nodes.

Regarding clustering, Osting et al. [85] applied a sparsification process on a simplicial complex, which downgrades the maximum dimension under a given threshold. It is proved that such a sparsification preserves the *up Laplacian*. The authors also generalize *Cheeger inequality* to a simplicial complex. The preservation of the spectrum is verified through experiments, and spectral clustering is performed as one application.

Advanced network representation learning methods have also been extended to the simplicial complex. Hajij et al. [43] use the autoencoder to perform simplex-level embedding. The encode function ($X \rightarrow \mathbb{R}^d$) maps each simplex to a fixed vector. The decode function ($\mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}^+$) maps each pair of simplices to a similarity score that reflects the relationship between the two simplices. An example of a user-defined similarity is the simplex-level adjacency matrix. Finally, the

representation of the whole simplicial complex is obtained by a weighted sum of the simplex-level representations.

5. HYPERGRAPHS

Another natural generalization of a graph is a *hypergraph*, which extends the edge space from $|V|^2$ to $|V|^{|V|}$. In other words, an edge of a hypergraph can be any subset of the vertices. The main difference with the simplicial complex is that any subset of an edge can exist independently from others.

Though easy to understand, the extremely sparse and high-dimensional edge space causes difficulty in computations. So, relatively more studies on hypergraphs have focused on the theoretical aspects rather than applications. For a more comprehensive review of concepts and measurements in hypergraphs, we refer interested readers to the survey by Lee et al. [59].

5.1 Foundation and Algorithms

Many definitions and properties of hypergraphs are directly inherited from graphs, such as *node degrees*, *hyperedge weights*, *simple/multi hypergraphs*, *hypergraph isomorphism*, and so on. Here, for brevity, we mainly focus on the definitions unique to hypergraphs.

5.1.1 Definition and Basics

We use $\mathcal{H} = (V, \mathcal{E})$ to distinguish a hypergraph from a graph, where only vertex set V remains the same. The edge set $\mathcal{E} = \{e | e \subseteq V\}$ is the set of subsets of V . We define *size* of edge $|e|$ as the number of nodes that belong to edge e .

We can have special kinds of hypergraphs:

- *d*-regular: each vertex has degree d ;
- *k*-uniform: each edge has size k ;
- *k*-partite: vertices belong to one of k different classes, and each edge has exactly one node from each class.

A hypergraph can be always represented by a bipartite graph of vertices and edges. The *biadjacency matrix* of this bipartite graph is a $|V| \times |E|$ matrix, which is also called the *incidence matrix* of the hypergraph.

5.1.2 Tensor Representation

The natural generalization of the adjacency matrix for the hypergraph is a *tensor*, often denoted by \mathbf{T} . For example, a 3-uniform hypergraph (or the subset of order-3 edges) can be represented as an order-3 tensor $\mathbf{T} \in \mathbb{R}^{|V| \times |V| \times |V|}$, i.e., the entry $(i, j, k) = 1$ when $(v_i, v_j, v_k) \in \mathcal{E}$. If the hypergraph is undirected, the corresponding tensor is *symmetric*, where its value at any permutation of (i, j, k) remains the same. A *simple tensor* can be written as the outer product of the vectors. The *rank* of a tensor is the minimum number of simple tensors whose linear combination equals that tensor.

(Tensor Decomposition) Due to tensor dimensionality, it is expensive and inconvenient to perform calculations directly on it. Decomposition techniques are widely used to decrease the dimension and preserve the graph characteristics. Here, we introduce two popular decomposition methods – CP decomposition and Tucker decomposition. For a detailed reference on the decomposition techniques of the tensor, interested readers can refer to [104].

CP (CANDECOMP/PARAFAC) Decomposition [22; 46] is also called *tensor rank decomposition* or *Canonical Polyadic Decomposition (CPD)*. The CP decomposition factorizes a tensor into a sum of component vectors. For example, a tensor $\mathbf{T} \in \mathbb{R}^{I \times J \times K}$ can be decomposed as

$$\mathbf{T} \approx \sum_{i=1}^R \mathbf{a}_i \otimes \mathbf{b}_i \otimes \mathbf{c}_i, \quad (13)$$

where R is an integer and $\mathbf{a}_i \in \mathbb{R}^I$, $\mathbf{b}_i \in \mathbb{R}^J$, $\mathbf{c}_i \in \mathbb{R}^K$, and \otimes denotes the outer product sign. This decomposition is often solved by some minimization algorithm.

Tucker Decomposition [114] is the generalization of *Singular Value Decomposition (SVD)* for the tensors. Tucker decomposition of a tensor $\mathbf{T} \in \mathbb{R}^{I \times J \times K}$ is represented as

$$\mathbf{T} \approx \mathcal{G} \times \mathbf{A} \times \mathbf{B} \times \mathbf{C} = \llbracket \mathcal{G}; \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket, \quad (14)$$

where $\mathbf{A} \in \mathbb{R}^{I \times P}$, $\mathbf{B} \in \mathbb{R}^{J \times Q}$, $\mathbf{C} \in \mathbb{R}^{K \times R}$. Here, tensor $\mathcal{G} \in \mathbb{R}^{P \times Q \times R}$ is called the *core tensor*.

5.1.3 Hypergraph Cuts

In dyadic graphs, a cut is defined as partitioning the vertices into two disjoint subsets, where the cut edge connects two nodes, one from each subset. However, for a hyperedge, there is no such fair split where we cannot assign some nodes to one side and the rest to the other side. Among various cut functions, there is one that might be more reasonable to minimize when solving cut-based hypergraph problems, called *all-or-nothing* [118].

Specifically, a hyperedge is in the middle of the cut when it is assigned to both sides of the cut. One can present the set of *cut hyperedges* by

$$\partial S = \{e \in \mathcal{E} : e \cap S \neq \emptyset \text{ and } e \cap \bar{S} \neq \emptyset\}. \quad (15)$$

The cut function is

$$\text{all-or-nothing}(s) = \sum_{e \in \partial S} w_e, \quad (16)$$

where w_e is the edge weight in the case of weighted hypergraphs.

5.1.4 Random Walks and Laplacian

In dyadic graphs, a sequence of vertices is sufficient to define a walk as there is only one way to traverse from one node to another in one step. However, due to the flexibility of hypergraphs, one has to specify the order of both edges and walks.

DEFINITION 5.1 (*s*-WALK). *Let \mathcal{H} be an r -uniform hypergraph, for $1 \leq s \leq r-1$, an s -walk of length k is defined as a sequence of vertices*

$$v_1, v_2, \dots, v_j, \dots, v_{(r-s)(k-1)+r}$$

together with a sequence of edges F_1, F_2, \dots, F_k such that

$$F_i = \{v_{(r-s)(i-1)+1}, v_{(r-s)(i-1)+2}, \dots, v_{(r-s)(i-1)+r}\}.$$

Basically, any two adjacent edges of an s -walk have exactly s vertices in their intersection. For the vertex set V , let V^s be the set of all ordered s -tuples consisting of s distinct elements in V . For example when $s = 2$,

$$V^2 = \{(v_1, v_2), (v_1, v_3), \dots, (v_2, v_1), (v_2, v_3), \dots\}.$$

To compute the Laplacian, we consider the following two cases [70]:

(1) In the case of $1 \leq s \leq r/2$, for any F_i , there will not be any intersection with F_{i+2} or F_{i-2} . So, the s -walk can be interpreted as the walk on a weighted dyadic graph. We define a weighted undirected graph $G^{(s)}$ over V^s as follows. Let the weight $w(x, y) = |\{F \in \mathcal{E} : [x] \cup [y] \subseteq F\}|$. Here, $[x] \cup [y]$ is the disjoint union of $[x]$ and $[y]$, and x, y are vertices of s -tuple of the original vertices. Then, we define the s -th *Laplacian* $\mathcal{L}^{(s)}$ of hypergraph \mathcal{H} to be the Laplacian of graph $G^{(s)}$.

(2) Another case is $r/2 < s \leq r-1$, where F_i also intersects with F_{i+2} or F_{i-2} (if it exists). We define a directed graph $D^{(s)}$ over the vertex set V^s as follows. With x, y are still the s -tuples of the original vertices, let (x, y) be a directed edge if $x_{r-s+j} = y_j$ for $1 \leq j \leq 2s-r$ and $[x] \cup [y]$ is an edge of \mathcal{H} . Then, we define the s -th *Laplacian* $\mathcal{L}^{(s)}$ of hypergraph \mathcal{H} to be the Laplacian of Eulerian directed graph $D^{(s)}$.

5.1.5 Downgrading a Hypergraph to a Dyadic Graph

To downgrade hypergraphs to dyadic graphs, one of the most straightforward ways is to perform (*clique*) *expansion*. For each hyperedge e , we enumerate all its size-two subedges to form a dyadic graph. Depending on the application, one can either inherit weights from the original hypergraph for these new edges or only keep the structural information [113]. While it is often much more convenient to perform dyadic graph algorithms, the expanded new graph indeed loses higher-order information.

A more general expansion is the *Multi-Level decomposition* [30]. In Multi-Level decomposition in addition to enumerating size-two subedges to form a dyadic graph, we construct hypernodes based on the coexistences within the original hyperedges. For example, assume there is a hyperedge of size $|e|$. At each layer k , we generate all possible $\binom{|e|}{k}$ hypernodes of size k to form a clique. By selecting k ranging from 2 to the maximum order of the hypergraph, we construct a corresponding dyadic graph for each k . The most appealing feature of such a decomposition is that one can easily reconstruct the hypergraph from its expansions across all layers.

5.2 Use Cases and Applications

Here, we summarize studies utilizing hypergraphs. One branch focuses on developing tools for network analysis, mostly extending graph theories to hypergraph-level. Another branch applies hypergraph to model higher-order interactions for network analysis.

5.2.1 Network Measurements

Many network properties and measurements are generalized from dyadic graphs to hypergraphs. As an analogy to walks on dyadic graphs, s -walk is proposed and applied to hypergraphs [3]. As mentioned in Definition 5.1, an s -walk is a series of hyperedges where the intersection nodes between any adjacent hyperedges have size greater than s . More specifically, a larger s indicates a component filled by denser overlapped hyperedges. Consequently, s -connected components and s -distance are also defined based on s -walks, which form a series of hypergraph analysis tools. Such measurements can help distinguish real-world networks from random hypergraphs generated by network models.

Centrality measures assess how important a node is in terms of its position and how it connects to other nodes in the graph. Three eigenvector centralities for uniform hypergraphs are defined by Benson [11]. Similar to dyadic networks, the centrality of a node in hypergraph can be influenced by centralities of all its neighbors. For each specific hyperedge to which it belongs, there could be a weighting function based on the centrality scores of its neighbors on that edge. First version is called *Clique motif Eigenvector Centrality*, where it simply refers to the total sum of centrality of all neighbors; Second is *Z-Eigenvector Centrality*, which multiplies the neighbors’ centralities on each hyperedge and then sums them up; Third is *H-Eigenvector Centrality*, which is the square root of the *Z-Eigenvector Centrality*. Experimental results show that none of these three centralities is consistently superior to others. So, one has to consider a specific objective to make an informed selection.

As motifs are significant subgraph patterns in dyadic graphs, *higher-order motifs* are defined in a similar way [69]. Given a set of nodes of size k , a higher-order motif is formed by a collection of hyperedges consisting of only these k nodes. As k increases, the motif variations can exponentially increase, which makes it impossible to enumerate and detect all motifs in the hypergraph. Addressing this problem, Lee et al. [61] propose *hypergraph motifs*, a special kind of higher-order motifs. Hypergraph motifs have a fixed structure that consists of three connected hyperedges. Based on overlapping nodes in hyperedges, all nodes are classified into one of the seven possible areas in a Venn diagram. Such a structure significantly simplifies motif detection and isomorphism checks, which are required to count motifs.

5.2.2 Generative Models

Graph generation algorithms aim to generate realistic graphs similar to those observed in the real world. Chodrow generalizes two variants (*vertex-labeled* and *stub-labeled*) of the *configuration model*, a well-known graph generation algorithm, to hypergraphs [24]. Configuration model in dyadic graphs requires the knowledge of the degree sequence. In hypergraphs, in addition to degree sequences, dimension sequence (sizes of edges) is also needed to generate a random graph. The vertex-labeled hypergraph configuration model is just a uniform distribution over the space of hypergraphs defined by degree and dimension. The stub-labeled hypergraph configuration model simply copies nodes as many times as their degrees and places them into a multiset. The algorithm then uniformly samples hyperedges based on the dimension sequence, where each node can only appear once in a specific hyperedge.

Lee et al. [60] extend the *Chung-Lu* model [26] to hypergraphs (*HyperCL*) by ensuring to preserve the distribution from the given degree sequence and the edge-size sequence as input. However, real-world hypergraphs exhibit stronger communities than random graphs. Addressing this issue, the authors further propose *HyperLap*, a multilevel HyperCL that introduces a group parameter L at each level, which aims to help reconstruct community patterns of real-world hypergraphs. New hyperedges generated within each group are expected to have high number of overlapping nodes, especially when the group is small.

Furthermore, a later study extends the *degree-corrected stochastic blockmodel* [51], which is a generative model of graphs with both community structure and degree sequences,

to hypergraphs [25]. For hyperedge candidates, the authors introduce an affinity function to compute the wiring possibility based on the group memberships of their nodes. Basically, more nodes in the same group have a higher probability of forming hyperedges. Three estimates—the affinity function, node labels, and node degrees, are alternatively learned by optimizing a likelihood function. To solve this objective, the authors propose an ‘All-or-Nothing’ (AON) *Louvain*-type algorithm [18] under the assumption that hyperedges are expected to lie fully within the cluster. Experimental results on both synthetic and empirical data validate the efficiency and accuracy of the framework.

5.2.3 Hypergraph Partitioning and Clustering

Hypergraph partitioning methods are generalized from classical graph cut problems. Veldt et al. [118] propose a comprehensive set of steps for solving hypergraph $s-t$ cuts problem. The first step is to select a splitting function, which maps the hyperedge that is going to be cut to a real number penalty (this has to be defined specifically). The authors specify a property for the splitting functions called *cardinality-based*, where the penalty only correlates with the sizes of split clusters. The hypergraph $s-t$ cuts problem is defined as minimizing the total splitting penalty for all crossing hyperedges. Various splitting functions are analyzed and tested over real-world datasets. Based on the results, a new clustering framework for hypergraphs is proposed [117], which minimizes the localized *ratio cut* objective. The algorithm requires a set of input nodes and returns a well-connected cluster that highly overlaps with the inputs. The running time of this algorithm only depends on the size of input set, and guarantees cuts or conductance under a specific bound.

The hypergraph cut problem can also be solved with tensor representations. By extending matrix-based methods, a tensor spectral clustering method is developed for partitioning higher-order networks [14]. For example, an order-3 undirected network can be represented as an order-3 symmetric tensor. As an analog to random walk on dyadic networks, a second-order Markov process is applied to express state changes on order-3 networks. Clustering of higher-order networks can be achieved by recursively partitioning the graph by minimizing *sweep cuts*. Such a clustering method preserves higher-order structures rather than just edges, as shown through experiments on both synthetic and real-world networks.

The clustering method can be further extended when introducing additional information. As for labeled networks, Amburg et al. [4] propose a novel hypergraph clustering framework based on given edge labels. The objective simultaneously minimizes (1) edges across clusters and (2) edges that do not belong to the assigned cluster. These two requirements can be combined by simply counting the number of nodes whose labels are inconsistent with their connected edges. In the case of two categories, this problem reduces to an $s-t$ cut problem by forming a dyadic graph and adding a terminal node to it. In the case of more than two categories, multiple approximation algorithms for such an NP-hard problem are developed, such as an LP relaxation and multiway cuts. Experimental results on synthetic and real-world graphs show that the proposed method outperforms baselines including *Majority Vote*, *Chromatic Balls* and *Lazy Chromatic Balls*.

Table 1: Tools for Discovering Motifs

Package Name	Year	Description	Official Link (if exist)
MFinder [54]	2005	motif detection, enumeration/edge sampling	https://www.weizmann.ac.il/mcb/UriAlon/download/ParTI
MAVisto [100]	2005	motif detection, enumeration	https://kim25.wwdns.kim.uni-konstanz.de/vanted/addons/mavisto/
FANMOD [122]	2005	motif detection, enumeration/node sampling	https://github.com/gabgabe/fanmod-cmd (unofficial)
Grochow–Kellis [42]	2007	motif detection, mapping	https://github.com/jptboy/CSCI3104_GC2 (unofficial)
MODA [83]	2009	motif detection, mapping/sampling, undirected only	https://github.com/smbadiwe/ParaMODA (unofficial)
Kavosh [52]	2009	motif detection, enumeration	https://github.com/shmohammadi86/Kavosh
G-Tries [92]	2010	motif detection, enumeration/mapping, undirected only	https://www.dcc.fc.up.pt/gtries/
TemporalMotif [87]	2016	temporal motif count	http://snap.stanford.edu/temporal-motifs/
MODET [88]	2019	motif detection, mapping, undirected only	https://github.com/sabyasachipatra/modet

Table 2: Tools for Learning Simplicial Complexes

Package Name	Environment	Description	Official Link
Simplicial	Python	topology, homology, filtrations	https://simplicial.readthedocs.io/en/latest/
Javaplex [111]	Matlab/Java	persistent homology, filtrations	https://github.com/appliedtopology/javaplex
Ripser [9]	C++	persistent homology, Vietoris–Rips filtrations	https://github.com/Ripser/ripser
simplextree	R	topology	https://github.com/peekxc/simplextree
Simplicial.jl	Julia	simplicial complexes, directed complexes	https://github.com/nebneuron/Simplicial.jl
simplicial-complex	JavaScript	structural and topological operations	https://www.npmjs.com/package/simplicial-complex
Dionysus 2	C++	persistent homology	https://mrzv.org/software/dionysus2/
DIPHA	C++	distributed, persistent homology	https://github.com/DIPHA/dipha
Perseus [78]	C++	persistent homology	https://people.maths.ox.ac.uk/nanda/perseus/
Moise	Maple	homology groups	https://www.math.drexel.edu/~ahicks/Moise/
TopoEmbedX [44]	Python	representation learning	https://github.com/pyt-team/TopoEmbedX

As discussed in Section 5.1.5, downgrading hypergraph to dyadic is also an effective way to utilize existing algorithms. Liu et al. [67] propose a hypergraph clustering method, called *Local Hypergraph Quadratic Diffusions (LHQD)*. The first step of LHQD reduces the hypergraph to a directed graph that preserves the conductance property of the original graph. The equality of conductance is achieved by introducing auxiliary nodes for each node. The second step of LHQD creates a source and a sink node in the directed graph, whose weights to the auxiliary nodes are equal to their degrees. Such a conceptual transformation ensures that the objective function becomes the same as the original problem. The performance is validated by performing clustering on two real-world networks.

5.2.4 Modeling Higher-Order Interactions

Many real-world interactions can be modeled directly as hypergraphs. The entities often have different types, but in terms of hypergraphs, research rarely emphasizes on node heterogeneity. Tan et al. [110] model the music recommendation problem as a hypergraph ranking problem. At the beginning, different objects (users, groups, tags, tracks, albums, and artists) are represented as nodes and pairwise relationships among them are identified. Hyperedges are created by combining edges based on the intrinsic connections among them, e.g. tracks in the same album. Given a query (set of nodes), the recommendation functions based on rank scores of all other nodes on the hypergraph. This scoring process is trained by the ground truth of node labels under the smoothing constraint based on which close nodes should have similar scores.

Similar to music recommender systems, the image retrieval problem can be modeled as the hypergraph ranking problem, where images are nodes that are assigned to hyperedges based on similarities. Liu et al. [68] applied a *soft hypergraph* model in which the entries on the incidence matrix are calculated using some similarity functions instead of arbitrary values. Hyperedges are created by selecting any

node as centroid and adding its k nearest neighbors. When an image query comes, the recommender system solves a linear system based on a cost function capturing hypergraph partitioning, which ensures vertices sharing many incidental hyperedges obtain similar labels.

Rather than built on the basis of similarities, hypergraphs are also directly used to model data in biology. Patro and Kingsford [89] model *network history inference* using a hypergraph structure. Generally speaking, *network history inference* is to find a small set of tuples that record the historical interactions between leaves on the protein network. The mapping of different states is represented as a hypergraph, where current state and correlated historical states are connected by order-3 hyperedges. The problem is solved by minimizing total cost over the network. Such a model can be applied to reconstruct the ancestral networks or to predict missing links.

5.2.5 Hypergraph Neural Networks

Graph neural networks have been generalized beyond pairwise interactions modeled as hypergraphs. Hypergraphs can be applied to either (1) model higher-order graph data as input matrices, such as the incidence matrix of hypergraphs; or to (2) build multilayer neural network structures but using higher-order forward- and back-propagation instead. Many state-of-the-art models, especially for graph neural networks, are extended to hypergraphs. Examples include hyper-models such as HGNN [34], HGAT [7; 120], MHCN [129], and HGCN [124; 119]. For a comprehensive survey on graph neural networks, readers are referred to the survey by Thomas et al. [112].

6. DATASETS AND TOOLS

We summarize some available datasets and tools for studying higher-order networks. As some of these official links might disappear in the future, we provide a comprehensive

Table 3: Tools for Learning Hypergraphs

Package Name	Environment	Description	Official Link
HyperG	R	Hypergraph Modeling	https://cran.r-project.org/web/packages/HyperG/
HyperNetX [49]	Python	Hypergraph Modeling, Visualization	https://github.com/pnnl/HyperNetX
GraphML [20]	XML	File Format	http://graphml.graphdrawing.org/index.html
hypergraph	R	Hypergraph Modeling	https://bioconductor.org/packages/3.15/bioc/html/hypergraph.html
SimpleHypergraphs.jl [106]	Julia	Hypergraph Modeling, Visualization	https://github.com/pszufe/SimpleHypergraphs.jl
halp	Python	Hypergraph Modeling, Algorithms	https://murali-group.github.io/halp/
kahypar [99]	Python	Hypergraph Partitioning	https://pypi.org/project/kahypar/
Tensorly [57]	Python	Tensor Learning	http://tensorly.org/stable/index.html
Tensors.jl [21]	Julia	Tensor Learning	https://juliahub.com/ui/Packages/Tensors/F7rKl/1.11.0
rTensor [64]	R	Tensor Learning	https://cran.r-project.org/web/packages/rTensor

Table 4: Applications of Higher-Order Networks

	Network Motifs	Simplicial Complexes	Hypergraphs
Statistical Significance	[102][53][108][31][105]		
Graph Classification	[76][75]		
Network Modeling	[91][63]	[27][128][48][16]	[24][60][25]
Clustering	[15][127][126][15][12]	[85]	[118][117][14][4][67]
Representation Learning	[94][101]	[43]	[116][40]
Link Prediction	[1]	[12]	[89][116]
Persistent Homology		[36][28][109][81][82][29][6][90][38][37][35]	
Analysis Tools and Measurements	[126][15]	[97]	[3][11][69][61]
Recommender System			[110][68]
Neural Networks			[34][7][120][129][124][119]

backup of all tools and datasets in our own repository.¹

6.1 Higher-Order Network Tools

For most network motif based studies, the first step is to find motifs. In Table 1, we list some scalable algorithms for enumerating or counting motifs. Among these methods, *enumeration* indicates an exhaustive search through the whole graph. *Sampling* indicates that the method calculates an estimated frequency of a given motif by sampling the node/edge and exploring its neighborhood. The *mapping* strategy is a reverse process of enumeration, which maps the given motif onto the whole network.

Table 2 collects packages and software for studying simplicial complexes. Some are tagged as ‘topology’, which are comprehensive packages that build the data structure from the lower level information. Some are software that are easy-to-use for most popular applications such as persistent homology and filtrations.

Table 3 collects packages for modeling hypergraphs and tensors. Most packages provide a data structure and implement the most basic algorithms using it; some packages support network visualization.

6.2 Higher-Order Datasets

We summarize some dataset resources with higher-order interactions:

- ARB Data²: A dataset repository (19 datasets—4 with millions of nodes, 10 with thousands of nodes, and 5 with hundreds of nodes) collected by Austin R. Benson. Most are higher-order networks from various fields, including temporal and labeled hypergraphs.
- LINQS³: A collection of 11 relational datasets (1 with a million nodes, others are thousands of nodes or less).

¹<https://github.com/haotian-syr/HON-tools>

²<https://www.cs.cornell.edu/~arb/data/>

³<https://linqs.soe.ucsc.edu/>

Many of them are collaboration networks, which naturally include higher-order interactions.

- Twitter Data⁴: Tweets can be modeled as higher-order networks by taking hashtags as nodes and co-appearances as edges (require preprocessing). Besides these, one can search twitter datasets online for any specific interest or collect data using APIs.
- Temporal Co-authorship⁵: Three large-scale (sizes: 27 million / 13 million / 41 thousand nodes) hypergraph datasets in both static and temporal forms [55].

6.3 Expected Time Complexities

While time complexity of exploring higher-order networks can vary across topologies and algorithms, some time complexities are typically expected for some basic higher-order algorithms. Here, we briefly list some expected time complexities for analyzing higher-order networks.

Motif counting: The time complexity of counting motifs depends highly on the structural complexity of given motifs as the essential algorithm for finding motifs involves checking for subgraph isomorphism, which is known to be NP-complete. Most fast motif-finding algorithms focus on size 3 or 4 motifs. For example, counting motifs of size 3 (triangles) can be solved in $O(|E|d_{max})$ [2] and counting motifs of size 4 can be solved in $O(|E|d_{max} + |E|^2)$ [72], where d_{max} is the maximum degree in the graph.

Homology groups: The time complexity of computing homology groups of the simplicial complex is $O(n^\omega)$, where n is the number of simplices and the exponent $\omega \leq 2.4$ [77]. Such an acceptable and stable complexity facilitates the wide usage of homology methods.

⁴<https://data.world/datasets/twitter>

⁵<https://github.com/kswoo97/pc1>

7. CONCLUSIONS AND FUTURE DIRECTIONS

We survey essential algorithms and applications in the literature of higher-order network modeling and analysis. In Table 4, we summarize the applications collected in this survey. Due to the scope of this survey, we have highlighted representative studies from each field. To explore each area comprehensively, we have directed readers to other surveys focusing on each domain.

However, research on higher-order networks has significant future potential. Here, we list some open problems or under-explored research directions:

7.1 Data Source and Modeling

Unlike dyadic graphs, not many repositories of higher-order network data are available. Building tools to collect, store, and model higher-order data is of significant interest for various academic use cases. Below, we list some essential, yet under developed, tools for studying higher-order data.

Recovering Higher-order Data in Dyadic Graphs:

Most existing network data is collected in dyadic form, which has already lost higher-order interactions. In motif-based studies, one cannot distinguish whether a specific motif is indeed a higher-order interaction or formed by a combination of dyadic interactions. It is therefore a challenge worth addressing to build tools that can distinguish higher-order interaction or that can rebuild higher-order networks from a dyadic graph.

Dynamic Higher-Order Graph: Most higher-order interactions are associated with time, such as protein interactions, hashtags, and group chats. Rather than a snapshot analysis of a network during some small interval, there is a demand for tools that can analyze the whole or partial network structure of a temporal higher-order network. Such tools enable real-time fast algorithms for various tasks including link prediction, community detection, anomaly detection, and the like.

Matrix/Tensor Representation: Topologies representing higher-order interactions are always associated with matrix or tensor representations, such as motif matrix [15], tensors [14] and incidence matrix. However, due to complexity and lack of mathematical support, algorithms on these matrix representations are not explored as extensively as matrix-based methods for dyadic graphs.

Interpretability and Causality: Most higher-order networks studies develop algorithms and applications that explore collected higher-order data. Current research rarely aims to interpret findings in higher-order networks or understand why some high-order interactions or patterns exist. As mentioned, some patterns might reflect important real-world information beyond network structures, such as the small functional unit in brain networks [108]. Interpretability is especially crucial as more black-box techniques (e.g., deep neural networks) or embedding methods are designed for higher-order networks.

7.2 Machine Learning Applications

Many real-world applications have focused on higher-order graphs. Below, we list three general application domains for higher-order networks that have a significant potential for future research.

Representation Learning: Representation learning is a powerful tool for transforming high-dimensional data into fixed-size vectors and has been extremely successful for downstream machine learning tasks, such as node classification, link prediction, among others. However, not many representation learning methods are introduced for higher-order networks. Hence, there is a significant demand for representation learning methods that can embed higher-order graphs both at the node-level [116; 40] and the graph-level.

Recommender Systems: One of the most direct uses of higher-order networks is in recommender systems. For example, nodes involved in same hyperedge can share some similarities. As discussed in Section 5.2.4, some applications such as music recommendation [110] and image retrieval [68] are developed. However, for individual recommendations, the advantage of higher-order graphs over simple or heterogeneous graphs needs to be further studied. Similarly, group recommendation [5] is another direction that has the potential to be studied.

Graph Neural Networks: Similar to modeling real-world interactions, the structure of neural networks can be designed to accept higher-order interactions as input when necessary. Due to insufficient studies on hypergraphs, there is only limited work that directly utilizes hypergraphs in neural networks. One main challenge is to extend the adjacency matrix, where some studies have considered the incidence matrix as a solution to this challenge [34].

8. REFERENCES

- [1] G. AbuOda, G. D. F. Morales, and A. Aboulnaga. Link prediction via higher-order motif features. *CoRR*, abs/1902.06679, 2019.
- [2] N. K. Ahmed, J. Neville, R. A. Rossi, and N. Duffield. Efficient graphlet counting for large networks. In *2015 IEEE International Conference on Data Mining*, pages 1–10, 2015.
- [3] S. G. Aksoy, C. A. Joslyn, C. O. Marrero, B. Praggastis, and E. Purvine. Hypernetwork science via high-order hypergraph walks. *EPJ Data Sci.*, 9(1):16, 2020.
- [4] I. Amburg, N. Veldt, and A. Benson. *Clustering in Graphs and Hypergraphs with Categorical Edge Labels*, page 706–717. Association for Computing Machinery, New York, NY, USA, 2020.
- [5] S. Amer-Yahia, S. B. Roy, A. Chawlat, G. Das, and C. Yu. Group recommendation: semantics and efficiency. *Proc. VLDB Endow.*, 2(1):754–765, aug 2009.
- [6] J. Arsuaga, N. A. Baas, D. DeWoskin, H. Mizuno, A. Pankov, and C. Park. Topological analysis of gene expression arrays identifies high risk molecular subtypes in breast cancer. *Appl. Algebra Eng. Commun. Comput.*, 23(1-2):3–15, 2012.
- [7] S. Bai, F. Zhang, and P. H. S. Torr. Hypergraph convolution and hypergraph attention. *CoRR*, abs/1901.08150, 2019.
- [8] A.-L. Barabási and M. Pósfai. *Network science*. Cambridge University Press, Cambridge, 2016.

- [9] U. Bauer. Ripser: efficient computation of Vietoris-rips persistence barcodes. *Journal of Applied and Computational Topology*, 2021.
- [10] A. Bavelas. Communication Patterns in Task-Oriented Groups. *Acoustical Society of America Journal*, 22(6):725, Jan. 1950.
- [11] A. R. Benson. Three hypergraph eigenvector centralities. *CoRR*, abs/1807.09644, 2018.
- [12] A. R. Benson, R. Abebe, M. T. Schaub, A. Jadbabaie, and J. M. Kleinberg. Simplicial closure and higher-order link prediction. *CoRR*, abs/1802.06916, 2018.
- [13] A. R. Benson, D. F. Gleich, and D. J. Higham. Higher-order network analysis takes off, fueled by classical ideas and new data. *CoRR*, abs/2103.05031, 2021.
- [14] A. R. Benson, D. F. Gleich, and J. Leskovec. Tensor spectral clustering for partitioning higher-order network structures. *CoRR*, abs/1502.05058, 2015.
- [15] A. R. Benson, D. F. Gleich, and J. Leskovec. Higher-order organization of complex networks. *Science*, 353(6295):163–166, 2016.
- [16] G. Bianconi and C. Rahmede. Emergent hyperbolic network geometry. *Scientific Reports*, 7(1), feb 2017.
- [17] C. Bick, E. Gross, H. A. Harrington, and M. T. Schaub. What are higher-order networks? *CoRR*, abs/2104.11329, 2021.
- [18] V. D. Blondel, J.-L. Guillaume, R. Lambiotte, and E. Lefebvre. Fast unfolding of communities in large networks. *Journal of statistical mechanics: theory and experiment*, 2008(10):P10008, 2008.
- [19] S. Boccaletti, V. Latora, Y. Moreno, M. Chavez, and D.-U. Hwang. Complex networks: Structure and dynamics. *Physics Reports*, 424(4):175–308, 2006.
- [20] U. Brandes, M. Eiglsperger, J. Lerner, and C. Pich. Graph markup language (graphml). In *Handbook of Graph Drawing and Visualization*, 2013.
- [21] K. Carlsson and F. Ekre. Tensors.jl — tensor computations in julia. *Journal of Open Research Software*, 7, 03 2019.
- [22] J. D. Carroll and J. J. Chang. Analysis of individual differences in multidimensional scaling via an n-way generalization of “eckart-young” decomposition. *Psychometrika*, 35:283–319, 1970.
- [23] C. Chen, J. He, N. Bliss, and H. Tong. Towards optimal connectivity on multi-layered networks. *IEEE Transactions on Knowledge and Data Engineering*, 29(10):2332–2346, 2017.
- [24] P. S. Chodrow. Configuration models of random hypergraphs. *Journal of Complex Networks*, 8(3), 08 2020. cnaa018.
- [25] P. S. Chodrow, N. Veldt, and A. R. Benson. Generative hypergraph clustering: From blockmodels to modularity. *Science Advances*, 7(28):eabh1303, 2021.
- [26] F. Chung and L. Lu. The average distances in random graphs with given expected degrees. *Proceedings of the National Academy of Sciences*, 99(25):15879–15882, 2002.
- [27] O. T. Courtney and G. Bianconi. Generalized network structures: The configuration model and the canonical ensemble of simplicial complexes. *Physical Review E*, 93(6), jun 2016.
- [28] V. De, Silva, and R. Ghrist. Coverage in sensor networks via persistent homology.
- [29] D. DeWoskin, J. Climent, I. Cruz-White, M. Vázquez, C. C. Park, and J. Arsuaga. Applications of computational homology to the analysis of treatment response in breast cancer patients. *Topology and its Applications*, 157:157–164, 2010.
- [30] M. Do, S.-e. Yoon, B. Hooi, and K. Shin. Structural patterns and generative models of real-world hypergraphs. pages 176–186, 08 2020.
- [31] C. Duclos, D. Nadin, Y. Mahdid, V. Tarnal, P. Picton, G. Vanini, G. Golmirzaie, E. Janke, M. S. Avidan, M. B. Kelz, G. A. Mashour, and S. Blain-Moraes. Brain network motifs are markers of loss and recovery of consciousness. *bioRxiv*, 2020.
- [32] P. Erdős and A. Rényi. On random graphs i. *Publicationes Mathematicae Debrecen*, 6:290, 1959.
- [33] P. L. Erdős, I. Miklós, and L. Soukup. Towards random uniform sampling of bipartite graphs with given degree sequence. 2010.
- [34] Y. Feng, H. You, Z. Zhang, R. Ji, and Y. Gao. Hypergraph neural networks, 2018.
- [35] A. Ghosh, B. Rozemberczki, S. Ramamoorthy, and R. Sarkar. Topological signatures for fast mobility analysis. In *Proceedings of the 26th ACM SIGSPATIAL International Conference on Advances in Geographic Information Systems*, SIGSPATIAL ’18, page 159–168, New York, NY, USA, 2018. Association for Computing Machinery.
- [36] R. Ghrist and A. Muhammad. Coverage and hole-detection in sensor networks via homology. In *IPSN ’05: Proceedings of the 4th international symposium on Information processing in sensor networks*, Piscataway, NJ, USA, 2005. IEEE Press.
- [37] C. Giusti, R. Ghrist, and D. S. Bassett. Two’s company, three (or more) is a simplex: Algebraic-topological tools for understanding higher-order structure in neural data, 2016.
- [38] C. Giusti, E. Pastalkova, C. Curto, and V. Itskov. Clique topology reveals intrinsic geometric structure in neural correlations. *Proceedings of the National Academy of Sciences*, 112(44):13455–13460, 2015.
- [39] D. F. Gleich. Pagerank beyond the web. *CoRR*, abs/1407.5107, 2014.
- [40] X. Gong, D. J. Higham, and K. Zygalakis. Generative hypergraph models and spectral embedding, 2023.

- [41] M. Granovetter. The strength of weak ties. *The American Journal of Sociology*, 78(6):1360–1380, May 1973.
- [42] J. A. Grochow and M. Kellis. Network motif discovery using subgraph enumeration and symmetry-breaking. In *Proceedings of the 11th Annual International Conference on Research in Computational Molecular Biology*, RECOMB’07, page 92–106, Berlin, Heidelberg, 2007. Springer-Verlag.
- [43] M. Hajij, G. Zamzmi, and X. Cai. Simplicial complex representation learning. *CoRR*, abs/2103.04046, 2021.
- [44] M. Hajij, G. Zamzmi, T. Papamarkou, N. Miolane, A. Guzmán-Sáenz, K. N. Ramamurthy, T. Birdal, T. K. Dey, S. Mukherjee, S. N. Samaga, N. Livesay, R. Walters, P. Rosen, and M. T. Schaub. Topological deep learning: Going beyond graph data, 2023.
- [45] T. Harko, F. S. Lobo, and M. Mak. Exact analytical solutions of the susceptible-infected-recovered (sir) epidemic model and of the sir model with equal death and birth rates. *Applied Mathematics and Computation*, 236:184–194, 2014.
- [46] R. A. Harshman. Foundations of the parafac procedure: Models and conditions for an "explanatory" multi-model factor analysis. 1970.
- [47] X. Hong Lin, Y. Han-bing, G. Cui-fang, and Z. Ping. Social network analysis based on network motifs. *Journal of Applied Mathematics*, 2014:1–6, 02 2014.
- [48] I. Iacopini, G. Petri, A. Barrat, and V. Latora. Simplicial models of social contagion. *Nature Communications*, 10(1), Jun 2019.
- [49] C. A. Joslyn, S. Aksoy, T. J. Callahan, L. E. Hunter, B. A. Jefferson, B. Praggastis, E. A. H. Purvine, and I. J. Tripodi. Hypernetwork science: From multidimensional networks to computational topology. *CoRR*, abs/2003.11782, 2020.
- [50] M. Kahle. Topology of random simplicial complexes: a survey, 2013.
- [51] B. Karrer and M. E. J. Newman. Stochastic block-models and community structure in networks. *Physical Review E*, 83(1), Jan 2011.
- [52] Z. R. M. Kashani, H. Ahrabian, E. Elahi, A. Nowzari-Dalini, E. Ansari, S. Asadi, S. Mohammadi, F. Schreiber, and A. Masoudi-Nejad. Kavosh : a new algorithm for finding network motifs. *BMC Bioinformatics*, 10, 2009. Article Number: 318.
- [53] N. Kashtan and U. Alon. Spontaneous evolution of modularity and network motifs. *Proceedings of the National Academy of Sciences*, 102(39):13773–13778, 2005.
- [54] N. Kashtan, S. Itzkovitz, R. Milo, and U. Alon. Efficient sampling algorithm for estimating subgraph concentrations and detecting network motifs. *Bioinformatics (Oxford, England)*, 20:1746–58, 08 2004.
- [55] S. Kim, D. Lee, Y. Kim, J. Park, T. Hwang, and K. Shin. Datasets, tasks, and training methods for large-scale hypergraph learning. *Data Mining and Knowledge Discovery*, 37:1–39, 07 2023.
- [56] M. Kivelä, A. Arenas, M. Barthelemy, J. P. Gleeson, Y. Moreno, and M. A. Porter. Multilayer networks. *Journal of Complex Networks*, 2(3):203–271, 07 2014.
- [57] J. Kossaifi, Y. Panagakis, A. Anandkumar, and M. Pantic. Tensorly: Tensor learning in python. *Journal of Machine Learning Research*, 20(26):1–6, 2019.
- [58] R. Lambiotte, M. Rosvall, and I. Scholtes. Understanding complex systems: From networks to optimal higher-order models, 2018.
- [59] G. Lee, F. Bu, T. Eliassi-Rad, and K. Shin. A survey on hypergraph mining: Patterns, tools, and generators, 2024.
- [60] G. Lee, M. Choe, and K. Shin. How do hyperedges overlap in real-world hypergraphs? - patterns, measures, and generators. *CoRR*, abs/2101.07480, 2021.
- [61] G. Lee, J. Ko, and K. Shin. Hypergraph motifs: Concepts, algorithms, and discoveries. *CoRR*, abs/2003.01853, 2020.
- [62] S. Leinhardt and J. Berger. The structure of positive interpersonal relations in small groups. *Sociological Theories in Progress. Boston: Houghton Mifflin*, 1971.
- [63] J. Leskovec, D. Huttenlocher, and J. Kleinberg. Signed networks in social media, 2010.
- [64] J. Li, J. Bien, and M. T. Wells. rtensor: An r package for multidimensional array (tensor) unfolding, multiplication, and decomposition. *Journal of Statistical Software*, 87(10):1–31, 2018.
- [65] J. Li, C. Chen, H. Tong, and H. Liu. *Multi-Layered Network Embedding*, pages 684–692.
- [66] L. Lim. Hodge laplacians on graphs. *CoRR*, abs/1507.05379, 2015.
- [67] M. Liu, N. Veldt, H. Song, P. Li, and D. F. Gleich. Strongly local hypergraph diffusions for clustering and semi-supervised learning. *Proceedings of the Web Conference 2021*.
- [68] Q. Liu, Y. Huang, and D. N. Metaxas. Hypergraph with sampling for image retrieval. *Pattern Recognition*, 44(10):2255–2262, 2011. Semi-Supervised Learning for Visual Content Analysis and Understanding.
- [69] Q. F. Lotito, F. Musciotto, A. Montresor, and F. Battiston. Higher-order motif analysis in hypergraphs. *Communications Physics*, 5(1), Apr. 2022.
- [70] L. Lu and X. Peng. High-ordered random walks and generalized laplacians on hypergraphs, 2011.
- [71] A. Madkour, W. G. Aref, F. U. Rehman, M. A. Rahman, and S. M. Basalamah. A survey of shortest-path algorithms. *CoRR*, abs/1705.02044, 2017.
- [72] D. Marcus and Y. Shavitt. Rage – a rapid graphlet enumerator for large networks. *Computer Networks*, 56(2):810–819, 2012.

- [73] N. Masuda, M. A. Porter, and R. Lambiotte. Random walks and diffusion on networks. *Physics Reports*, 716-717:1–58, 2017. Random walks and diffusion on networks.
- [74] A. P. Millán, J. J. Torres, and G. Bianconi. Explosive higher-order kuramoto dynamics on simplicial complexes. *Physical Review Letters*, 124(21), may 2020.
- [75] R. Milo, S. Itzkovitz, N. Kashtan, R. Levitt, S. Shen-Orr, I. Ayzenshtat, M. Sheffer, and U. Alon. Superfamilies of evolved and designed networks. *Science*, 303(5663):1538–1542, 2004.
- [76] R. Milo, S. Shen-Orr, S. Itzkovitz, N. Kashtan, D. Chklovskii, and U. Alon. Network motifs: simple building blocks of complex networks. *Science*, 298(5594):824–827, October 2002.
- [77] N. Milosavljević, D. Morozov, and P. Skraba. Zigzag persistent homology in matrix multiplication time. In *Proceedings of the Twenty-Seventh Annual Symposium on Computational Geometry*, SoCG '11, page 216–225, New York, NY, USA, 2011. Association for Computing Machinery.
- [78] K. Mischaikow and V. Nanda. Morse theory for filtrations and efficient computation of persistent homology. *Discrete Comput. Geom.*, 50(2):330–353, sep 2013.
- [79] A. Muhammad and M. Egerstedt. Control using higher order laplacians in network topologies. In *Proc. of 17th International Symposium on Mathematical Theory of Networks and Systems*, Kyoto, pages 1024–1038, 2006.
- [80] R. Mulas, C. Kuehn, and J. Jost. Coupled dynamics on hypergraphs: Master stability of steady states and synchronization. *Phys. Rev. E*, 101:062313, Jun 2020.
- [81] V. Nanda and R. Sazdanovic. Simplicial models and topological inference in biological systems. 2014.
- [82] M. Nicolau, R. Tibshirani, A.-L. Børresen-Dale, and S. S. Jeffrey. Disease-specific genomic analysis: identifying the signature of pathologic biology. *Bioinformatics*, 23(8):957–965, 02 2007.
- [83] S. Omid and F. Schreiber. Moda: An efficient algorithm for network motif discovery in biological networks. *Genes & genetic systems*, 84:385–95, 10 2009.
- [84] A. Ortega, P. Frossard, J. Kovačević, J. M. F. Moura, and P. Vandergheynst. Graph signal processing: Overview, challenges, and applications. *Proceedings of the IEEE*, 106(5):808–828, 2018.
- [85] B. Osting, S. Palande, and B. Wang. Towards spectral sparsification of simplicial complexes based on generalized effective resistance. *CoRR*, abs/1708.08436, 2017.
- [86] L. Page, S. Brin, R. Motwani, and T. Winograd. The pagerank citation ranking: Bringing order to the web. Technical Report 1999-66, Stanford InfoLab, November 1999. Previous number = SIDL-WP-1999-0120.
- [87] A. Paranjape, A. R. Benson, and J. Leskovec. Motifs in temporal networks. *CoRR*, abs/1612.09259, 2016.
- [88] S. Patra and A. Mohapatra. Application of dynamic expansion tree for finding large network motifs in biological networks. *PeerJ*, 7:e6917, 05 2019.
- [89] R. Patro and C. Kingsford. Predicting protein interactions via parsimonious network history inference. *Bioinformatics*, 29(13):i237–i246, 06 2013.
- [90] G. Petri, P. Expert, F. Turkheimer, R. Carhart-Harris, D. Nutt, P. J. Hellyer, and F. Vaccarino. Homological scaffolds of brain functional networks. *Journal of The Royal Society Interface*, 11(101):20140873, 2014.
- [91] N. Pržulj, D. G. Corneil, and I. Jurisica. Modeling interactome: scale-free or geometric? *Bioinformatics*, 20(18):3508–3515, 07 2004.
- [92] P. Ribeiro and F. Silva. G-tries: An efficient data structure for discovering network motifs. pages 1559–1566, 01 2010.
- [93] G. Rossetti and R. Cazabet. Community discovery in dynamic networks: A survey. *ACM Comput. Surv.*, 51(2), feb 2018.
- [94] R. A. Rossi, N. K. Ahmed, E. Koh, S. Kim, A. Rao, and Y. A. Yadkori. Hone: Higher-order network embeddings, 2018.
- [95] A. Sandryhaila and J. Moura. Big data analysis with signal processing on graphs: Representation and processing of massive data sets with irregular structure. *Signal Processing Magazine, IEEE*, 31:80–90, 09 2014.
- [96] A. Sandryhaila and J. M. F. Moura. Discrete signal processing on graphs. *CoRR*, abs/1210.4752, 2012.
- [97] M. T. Schaub, A. R. Benson, P. Horn, G. Lippner, and A. Jadbabaie. Random walks on simplicial complexes and the normalized hodge laplacian. *CoRR*, abs/1807.05044, 2018.
- [98] M. T. Schaub, Y. Zhu, J.-B. Seby, T. M. Roddenberry, and S. Segarra. Signal processing on higher-order networks: Livin’ on the edge... and beyond. *Signal Processing*, 187:108149, 2021.
- [99] S. Schlag. *High-Quality Hypergraph Partitioning*. PhD thesis, Karlsruhe Institute of Technology, Germany, 2020.
- [100] F. Schreiber and H. Schwöbbermeyer. Frequency concepts and pattern detection for the analysis of motifs in networks. In C. Priami, E. Merelli, P. Gonzalez, and A. Omicini, editors, *Transactions on Computational Systems Biology III*, pages 89–104, Berlin, Heidelberg, 2005. Springer Berlin Heidelberg.
- [101] P. Shao, Y. Yang, S. Xu, and C. Wang. Network embedding via motifs. *ACM Trans. Knowl. Discov. Data*, 16(3), oct 2021.
- [102] S. S. Shen-Orr, R. Milo, S. Mangan, and U. Alon. Network motifs in the transcriptional regulation network of escherichia coli. *Nature Genetics*, 31:64–68, 2002.

- [103] D. I. Shuman, S. K. Narang, P. Frossard, A. Ortega, and P. Vandergheynst. Signal processing on graphs: Extending high-dimensional data analysis to networks and other irregular data domains. *CoRR*, abs/1211.0053, 2012.
- [104] N. D. Sidiropoulos, L. De Lathauwer, X. Fu, K. Huang, E. E. Papalexakis, and C. Faloutsos. Tensor decomposition for signal processing and machine learning. *IEEE Transactions on Signal Processing*, 65(13):3551–3582, 2017.
- [105] A. Sizemore, C. Giusti, A. Kahn, J. Vettel, R. Betzel, and D. Bassett. Cliques and cavities in the human connectome. *Journal of Computational Neuroscience*, 44:1–31, 02 2018.
- [106] C. Spagnuolo, G. Cordasco, P. Szufel, P. Pralat, V. Scarano, B. Kaminski, and A. Antelmi. Analyzing, exploring, and visualizing complex networks via hypergraphs using SimpleHypergraphs.jl. *Internet Mathematics*, apr 2020.
- [107] E. H. Spanier. *Algebraic topology*. McGraw-Hill Book, New York, 1966. Includes index.
- [108] O. Sporns and R. Kötter. Motifs in brain networks. *PLoS biology*, 2:e369, 12 2004.
- [109] A. Tahbaz-Salehi and A. Jadbabaie. Distributed coverage verification in sensor networks without location information. *IEEE Transactions on Automatic Control*, 55(8):1837–1849, 2010.
- [110] S. Tan, J. Bu, C. Chen, B. Xu, C. Wang, and X. He. Using rich social media information for music recommendation via hypergraph model. *ACM Trans. Multimedia Comput. Commun. Appl.*, 7S(1), nov 2011.
- [111] A. Tausz, M. Vejdemo-Johansson, and H. Adams. JavaPlex: A research software package for persistent (co)homology. In H. Hong and C. Yap, editors, *Proceedings of ICMS 2014*, Lecture Notes in Computer Science 8592, pages 129–136, 2014. Software available at <http://appliedtopology.github.io/javaplex/>.
- [112] J. Thomas, A. Moallem-Oureh, S. Beddar-Wiesing, and C. Holzhüter. Graph neural networks designed for different graph types: A survey, 04 2022.
- [113] H. Tian, S. Jin, and R. Zafarani. Exploiting cross-order patterns and link prediction in higher-order networks. In *2022 IEEE International Conference on Data Mining Workshops (ICDMW)*, pages 1–9, 2022.
- [114] L. R. Tucker. Some mathematical notes on three-mode factor analysis. *Psychometrika*, 31:279–311, 1966c.
- [115] J. Ugander, L. Backstrom, and J. M. Kleinberg. Subgraph frequencies: Mapping the empirical and extremal geography of large graph collections. *CoRR*, abs/1304.1548, 2013.
- [116] M. Vaida and K. Purcell. Hypergraph link prediction: Learning drug interaction networks embeddings. In *2019 18th IEEE International Conference On Machine Learning And Applications (ICMLA)*, pages 1860–1865, 2019.
- [117] N. Veldt, A. R. Benson, and J. Kleinberg. *Minimizing Localized Ratio Cut Objectives in Hypergraphs*, page 1708–1718. Association for Computing Machinery, New York, NY, USA, 2020.
- [118] N. Veldt, A. R. Benson, and J. M. Kleinberg. Hypergraph cuts with general splitting functions. *CoRR*, abs/2001.02817, 2020.
- [119] J. Wang, K. Ding, L. Hong, H. Liu, and J. Caverlee. *Next-Item Recommendation with Sequential Hypergraphs*, page 1101–1110. Association for Computing Machinery, New York, NY, USA, 2020.
- [120] J. Wang, K. Ding, Z. Zhu, and J. Caverlee. Session-based recommendation with hypergraph attention networks. *CoRR*, abs/2112.14266, 2021.
- [121] D. J. Watts and S. H. Strogatz. Collective dynamics of ‘small-world’ networks. *Nature*, 393(6684):440–442, 1998.
- [122] S. Wernicke. Efficient detection of network motifs. *IEEE/ACM Trans. Comput. Biol. Bioinformatics*, 3(4):347–359, oct 2006.
- [123] J. Xu, T. L. Wickramaratne, and N. V. Chawla. Representing higher-order dependencies in networks. *Science Advances*, 2(5):e1600028, 2016.
- [124] N. Yadati, M. Nimishakavi, P. Yadav, V. Nitin, A. Louis, and P. Talukdar. Hypergcnn: A new method of training graph convolutional networks on hypergraphs, 2018.
- [125] C. Yang, Y. Xiao, Y. Zhang, Y. Sun, and J. Han. Heterogeneous network representation learning: Survey, benchmark, evaluation, and beyond. *CoRR*, abs/2004.00216, 2020.
- [126] H. Yin, A. R. Benson, and J. Leskovec. Higher-order clustering in networks. *Physical Review E*, 97(5), May 2018.
- [127] H. Yin, A. R. Benson, J. Leskovec, and D. F. Gleich. Local higher-order graph clustering. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD ’17*, page 555–564, New York, NY, USA, 2017. Association for Computing Machinery.
- [128] J.-G. Young, G. Petri, F. Vaccarino, and A. Patania. Construction of an efficient sampling from the simplicial configuration model. *Physical Review E*, 96(3), sep 2017.
- [129] J. Yu, H. Yin, J. Li, Q. Wang, N. Q. V. Hung, and X. Zhang. Self-supervised multi-channel hypergraph convolutional network for social recommendation. *CoRR*, abs/2101.06448, 2021.