



Published in final edited form as:

*Found Data Sci.* 2022 June ; 4(2): 165–216. doi:10.3934/fods.2022002.

## ASPECTS OF TOPOLOGICAL APPROACHES FOR DATA SCIENCE

**Jelena Grbi** <sup>†</sup>,

School of Mathematical Sciences, University of Southampton, Southampton, UK

**Jie Wu** <sup>\*,†</sup>,

School of Mathematical Sciences, Center of Topology and Geometry based Technology, Hebei Normal University, Yuhua District, Shijiazhuang, Hebei, 050024 China

Yanqi Lake Beijing Institute of Mathematica Sciences, Yanqihu, Huairou District, Beijing, 101408 China

**Kelin Xia** <sup>†</sup>,

School of Physical and Mathematical Sciences, Nanyang Technological University, SPMS-MAS-05-18, 21 Nanyang Link, 1, Singapore 63737

**Guo-Wei Wei** <sup>†</sup>

Department of Mathematics, Department of Computer Science and Engineering, Department of Biochemistry and Molecular Biology, Michigan State University, MI 48824, USA

### Abstract

We establish a new theory which unifies various aspects of topological approaches for data science, by being applicable both to point cloud data and to graph data, including networks beyond pairwise interactions. We generalize simplicial complexes and hypergraphs to super-hypergraphs and establish super-hypergraph homology as an extension of simplicial homology. Driven by applications, we also introduce super-persistent homology.

### Keywords

Topological data analysis; hypergraph; super-hypergraph; persistent homology; super persistent homology; simplicial complex; Delta set; scoring scheme

### 1. Introduction.

Topological data analysis (TDA) is a new-born research area mainly stemming from the pioneering works on persistent homology carried in [57, 148] and the landmark paper of Gunnar Carlsson [36] published in 2009. Under its rapid development, TDA has achieved various successful applications in many areas of sciences and technologies such as material science [81, 119, 91, 103], 3D shape analysis [126, 129], multivariate time series analysis [124], biology [29, 32, 31, 90, 140, 105], sensor networks [48], scientific visualization [128],

<sup>\*</sup>Corresponding author: Jie Wu. wujie@bimsa.cn.

<sup>†</sup>JG, JW, KX and GW should be considered joint first author.

image analysis [33, 109, 125, 17, 64], dynamics systems [100], etc. A great number of TDA softwares have been developed, including JavaPlex [127], Perseus [104], Dionysus [1], jHoles [21], GUDHI [97], Ripser [13], PHAT [15], DIPHA [14], and R-TDA package [61]. The results from TDA can be visualized by persistent diagram [101] and persistent barcode [66], and further transformed into different representations that are suitable for machine learning models, such as Betti curve [118], persistent landscape [25, 24], persistent image [2], persistent path and signature [43], persistent codebook [147], persistent 2D matrices [29], and others [3, 85, 29, 32, 105]. TDA-based machine learning models [40, 111] have been used in various areas, including image analysis [9, 79, 112, 95, 108, 67], shape analysis [22, 145, 92, 78, 143], time-series data analysis [124, 6, 144, 133, 130], computational biology [109, 50, 29, 32], noise data [107], sphere packing [117], language analysis [146], etc. Other than the traditional persistent homology, other TDA models have been proposed for the detailed characterization of the data, such as persistent local homology [18, 62, 16, 19, 5], persistent cohomology [49], multidimensional persistent homology [35, 34, 45, 37], element-specific persistent homology [29, 105], persistent functions [20], persistent spectral [132, 98], persistent Ricci curvatures [134, 135], etc. The wide applications of TDA have made topology as one of the most commonly used mathematical tools in Data Science [120, in Section 1.3]. A summary of TDA and TDA-based learning models can be found in Figure 1.

In a survey paper [41], Chazal and Michel outlined a pipeline that stresses the role of topology and geometry in data science:

- i.** Input data is given in the form of a finite set of points coming with a notion of distance;
- ii.** A “continuous shape” is built on top of the input data: this results in a structure over the data;
- iii.** Topological and geometric information is extracted from the structure;
- iv.** The topological and geometric information is the output of the analysis and forms the new representation of the data, allowing for an in-depth modeling of the original data.

Such an approach can be naturally applied to point cloud data with a drawback that it can not be immediately or directly applied to non-Euclidean data such as graphs for abstract relationship.

The purpose of this article is to provide a new theory that unifies various aspects of topological approaches for data science that is suitable for both point cloud data and graphic data. In our setting, we explore topological structures on graph data with scoring schemes. The popular persistent homology can be obtained as special cases of our more general theory from a natural transformation from point cloud data to graphic data with scoring schemes.

We start with a graph, which is the working graph for the data analytic purpose. Our approach consists of the following steps:

- A. We introduce a homology theory of a collection of subgraphs of the working graph, which is a canonical extension of simplicial homology theory. Briefly, this homology theory will canonically obtain “topological invariants” for collections of subgraphs associated to sample data or experimental data.
- B. We assign a *scoring scheme* on the working graph  $G$ , where a scoring scheme is a function from the set of subgraphs of  $G$  to the set of real numbers. The scoring scheme induces a persistence on homologies in (A), which in this article will be called as *super persistent homology*, as well as its derived *topological features* such as super-persistence diagrams and super-persistence modules for data analytics.
- C. The current persistent homology on point cloud data can be deduced from (A) and (B). Hence our approach is suitable for performing topological data analysis on both graphic data and point cloud data.

The pipeline of our super persistent homology is as follows:

1. The input is assumed to be a finite (or infinite) graph  $G$  with
  - i. A scoring scheme and
  - ii. A selection of subgraphs.

The definition of the scoring scheme on the data is usually given as an input or guided by applications. It is however important to notice that the choice of a scoring scheme may be critical to revealing interesting topological and geometric features of the data. The selection of subgraphs on the data is also usually given as an input or guided by the applications at hand. Again it is important to notice that the selection of subgraphs may be critical to revealing interesting topological and geometric features of the data.

2. An “abstract geometry-like” shape is built on top of the data in order to highlight the underlying topological structure. This is a nested family of *super-hypergraphs* filtered by the scoring scheme that reflects the structure of the data at different scales. Super-hypergraphs can be seen as higher dimensional generalizations of neighboring graphs that are classically built on top of data in many standard data analysis or learning algorithms. The challenge here is to define such structures that reflect relevant information about the structure of the data and that can be effectively constructed and manipulated in practice.
3. The extracted topological and geometric information provides new families of features and descriptors of the data. These can be used to better understand the data or they can be combined with other kinds of features for further analysis and machine learning tasks. Demonstrating the added-value and the complementarity (with respect to other features) of the information provided by super persistent homology is an important issue at this step.
4. Adjust the choice of scoring scheme and the selection of subgraphs to get better features and descriptors of the data.

5. One can repeat the procedure on the choices of scoring scheme and selections of subgraphs to obtain the best suitable features and descriptors of the data.

For analyzing both point cloud data and graph data in a unified way, we convert a point cloud data into a graph data by adding exact one edge to each pair of the points in the data to form a fully connected graph. There have been extensive explorations on topological and categorical structures on graphs. In Appendix A, we give a brief review of simplicial complexes constructed from graphs. In addition to the commonly-used construction of a clique complex, here we give a short discussion on a famous construction called *neighborhood complex* introduced by Lovász in his foundational work [94] in the area of topological combinatorics. Let  $G$  be a graph. Define the *neighborhood complex*  $\mathcal{N}(G)$  to be an abstract simplicial complex whose vertices are the vertices of  $G$  and whose simplices are those subsets of the vertex set  $V(G)$  which has a common neighbor. Lovász Theorem [94, Theorem 2] states that if  $\mathcal{N}(G)$  is  $(k-2)$ -connected, then  $G$  is not  $k$ -colorable<sup>1</sup>. As Lovász remarked [94, p.320], in the case  $k=2$ , the converse statement is also true, and so  $\mathcal{N}(G)$  is connected if and only if  $G$  is not bipartite. Note that a simplicial complex is connected if and only if its 0-th Betti number is 1. If we choose Lovász' neighborhood complex as the construction filtered by a scoring scheme, then the resulting barcodes of the 0-th super-persistent homology introduced in Section 3 would immediately give the bipartite information for the level subgraphs of the filtration.

The geometric realization of neighborhood complex (of a graph) is quite different from that of the clique complex, see Appendix A for details. This indicates that one could have different topological structures from the same working graph.

It is also possible that simplicial complexes model data from practical problems. According to the review article [12] supported by hundreds of references, extensive research has been recently taken on the networks beyond pairwise interactions. A simplicial complex (hypergraph) can be naturally constructed by taking the nodes having a group-interaction as a simplex (hyperedge).

Hypergraph is a preferred model in some practical problems. Although simplicial complexes overcome some of the problems encountered by other lower dimensional representations, they are still quite limited by the requirement on the existence of all subfaces. In some cases, such as group interactions in social systems, this constraint is too restrictive. In other cases, such as author collaborations in scientific papers and gene pathways, the inclusion constraint can be less easily justified. Hypergraphs provide the most general and unconstrained description of higher-order interactions, see [12, paragraphs 2–4, page 7].

For graph data having higher-order interactions, it is natural to take the collections of subgraphs that have group-interactions. This arises a mathematical question: *Given  $\mathcal{H}$  a collection of subgraphs of a working graph, how to introduce topology on  $\mathcal{H}$  with as less constraints as possible on  $\mathcal{H}$ ?*

---

<sup>1</sup>As a consequence of this theorem, he solved the Kneser conjecture in combinatorics [94 Theorem 1].

In this article, we are going to answer this question. Here we give some observations and brief ideas. For graph data, it is likely that two different subgraphs share the same nodes. This observation implies that the notion of simplicial complexes does not work well for the collections of subgraphs due to the fact that any simplex in a simplicial complex must be uniquely determined by its vertices (nodes). The notion of  $\Delta$ -set in algebraic topology, which is a generalization of simplicial complex, can solve this problem. Roughly, a  $\Delta$ -set  $X$  is a graded (multi-layered) set labeled by  $X_0, X_1, X_2, \dots$ , where  $X_n$  can be intuitively viewed as the set of  $n$ -dimensional simplices. The structure of a  $\Delta$ -set  $X$  is given by face operators, that is, we assign  $n + 1$  face operators labeled by  $d_0, d_1, \dots, d_n$  as functions from the set of  $n$ -dimensional simplices to the set of  $(n - 1)$ -dimensional simplices, satisfying the  $\partial$ -identity (the matching rule of faces) that  $d_i d_j = d_j d_{i+1}$  for  $i < j$ . The geometric realization of a  $\Delta$ -set is  $\Delta$ -complex, a notion defined in the popular textbook of algebraic topology [80]. A  $\Delta$ -set can be described in terms of feed-forward neural networks, see Subsection 3.5. Using the notion of  $\Delta$ -sets, we are allowed to select two or more subgraphs treated as  $n$ -simplices sharing the same vertices. For addressing *as less constraints as possible* in the question, we introduce the notion of *super-hypergraph* in Subsection 2.3, which is defined as a graded (multi-layered) subset of a  $\Delta$ -set. If a  $\Delta$ -set is given by an oriented simplicial complex, then our definition of super-hypergraph coincides with a hypergraph. Roughly speaking, a super-hypergraph is an extension of a hypergraph that allows hyperedges to form a multiset. An important aspect in the present article is that simplicial homology can be naturally extended to a homology theory on super-hypergraphs as described in Section 2, and so there are topological invariants (in terms of homology) on super-hypergraphs.

To introduce persistence on graph data, we propose to use the notion of a scoring scheme, which is a real valued function on finite subgraphs of the working graph. As discussed in Subsection 3.2, the classical persistent homology can be converted into super-persistent homology under the notion of scoring schemes. Hence super-persistent homology is a novel topological approach that can be applied to broader objects in data science.

Mathematically, the main results of the article are Theorems 2.7, 2.20 and 3.7. Theorem 2.7 states that hypergraph homology does not depend on the choice of orientation, which shows that the hypergraph homology provides the invariants on the intrinsic structure of hypergraph. The Mayer-Vietoris sequence is one of the fundamental tools for computing simplicial homology. Theorem 2.20 gives an analogue of the classical Mayer-Vietoris sequence for super-hypergraphs. Theorem 3.7 gives the structure theorem on super-persistent homology.

We should point out that it is commonly known (such as the textbook of algebraic topology [80]) that the computation of simplicial homology can be largely simplified using the notion of  $\Delta$ -complex, if the geometric shape can be homotopically deformed. Also the complexity of computing super-hypergraph homology is essentially the same as that of computing the simplicial homology of  $\Delta$ -sets.

It should be pointed out that even though this article is a theoretic framework, it has deep roots in application. One of the main motivations is from the drug design, in particular the analysis of binding affinities between proteins and ligands, i.e., how powerful is a

drug (ligand) for a certain disease (protein-related). Figure 2 illustrates the topological representations for the protein-ligand complex (ID: 3E6Y). Biologically, the interactions between the ligand and the protein are key to the drug potency and efficacy. Graph models (as in Figure 2(B)) is widely used to characterize the atomic interactions. A major drawback for graph models is that they only characterize pair-wise interactions and can not be used in many-body interaction characterization. Recently, we have used hypergraph models to describe various types of interactions (beyond pair-wise ones) for biomolecular interactions (Figure 2(C)). Essentially, many-body interactions can be modeled as hyperedges. It has been found that hypergraph models have great advantages over graph models in molecular representation. More details can be found in Section 4. A further generalization of hypergraphs to super-hypergraphs (Figure 2(D)) provides us more flexibility to characterize the complicated topology within each individual hyperedge. We believe that the super-hypergraph model provides an upgraded topological approach to data science, and can help to foster further interactions between topology and data science. The content of the article is organized in the following way.

## 2. Homology theory on super-hypergraphs.

Recently, homology of hypergraphs has opened new avenues for using topological tools in data analysis. Hypergraphs have been used for data analytics in various areas of sciences from social networks to molecular bioscience. The notion of hypergraph can be generalized as super-hypergraph, see Subsection 2.3. These objects are important for understanding the different explorations of topological structures on spaces of subgraphs. This setting realizes our aim to establish a unified approach to explore data science using topological combinatorics. The purpose of this section is to establish a homology theory of super-hypergraphs as a natural extension of simplicial homology and homology of hypergraphs.

### 2.1. Algebraic lemmas.

The following algebraic tools will be needed to define a homology theory of super-hypergraphs. Although we will only make use of chain complexes of abelian groups, we note that using simplicial group models, homotopy groups can be combinatorially defined using Moore chain complexes, which are chain complexes of possibly non-abelian groups [47, 139]. There are many studies of the homotopy type of topological structures of subgraphs as indicated in the references in Appendix A. Therefore, we consider chain complexes of possibly non-abelian groups so that the results in this subsection may be relevant for future research.

A graded group  $G_* = \{G_n\}_{n \in \mathbb{Z}}$  is a sequence of groups  $G_n$  indexed by the integers. A graded subgroup  $G'_* = \{G'_n\}_{n \in \mathbb{Z}}$  of  $G_* = \{G_n\}_{n \in \mathbb{Z}}$  is a sequence of subgroups  $G'_n$  such that  $G'_n \leq G_n$  for  $n \in \mathbb{Z}$ . A chain complex  $G_*$  of groups is a graded group  $G_*$  with a group homomorphism  $\partial_n = \partial_n^{G_*}: G_n \rightarrow G_{n-1}$  for  $n \in \mathbb{Z}$  such that the composite

$$\partial_{n-1} \circ \partial_n: G_n \rightarrow G_{n-2}$$

is the trivial homomorphism. Let us emphasise that in this definition we do not require  $G_n$  to be abelian. A subcomplex  $C_*$  of  $G_*$  is a graded subgroup  $C_*$  of  $G_*$  such that

$$\partial_n^{G_*}(C_n) \subseteq C_{n-1}$$

for each  $n \in \mathbb{Z}$ . So  $C_*$  together with the restrictions  $\partial_n^{G_*}|_{C_n} : C_n \rightarrow C_{n-1}$  forms a chain complex.

**Definition 2.1.** Let  $G_*$  be a chain complex of groups and let  $D_*$  be a graded subgroup of  $G_*$ . Define

$$\sup_*^{G_*}(D_*) = \cap \{C_* \mid D_n \leq C_n \text{ for } n \in \mathbb{Z}, \text{ and } C_* \text{ is a subcomplex of } G_*\}$$

$$\inf_*^{G_*}(D_*) = \prod \{E_* \mid E_n \leq D_n \text{ for } n \in \mathbb{Z}, \text{ and } E_* \text{ is a subcomplex of } G_*\}.$$

For simplicity, if the embedding of  $D_* \subseteq G_*$  is clear, we denote  $\sup_*^{G_*}(D_*)$  by  $\sup_*(D_*)$  and  $\inf_*^{G_*}(D_*)$  by  $\inf_*(D_*)$ .

**Proposition 2.2.** Let  $G_*$  be a chain complex of groups and let  $D_*$  be a graded subgroup of  $G_*$ . Then

1.  $\sup_*(D_*)$  is the smallest subcomplex of  $G_*$  containing  $D_*$ . Moreover,

$$\sup_n(D_*) = D_n \cdot \partial_{n+1}^{G_*}(D_{n+1})$$

is the product of  $D_n$  and  $\partial_{n+1}^{G_*}(D_{n+1})$ .

2.  $\inf_*(D_*)$  is the largest subcomplex of  $G_*$  contained in  $D_*$ . Moreover,

$$\inf_n(D_*) = D_n \cap \partial_n^{-1}(D_{n-1})$$

is the intersection of  $D_n$  and  $\partial_n^{-1}(D_{n-1})$ .

*Proof.* (1) The first part follows from the definition. Let

$$\tilde{D}_n = D_n \cdot \partial_{n+1}^{G_*}(D_{n+1})$$

for  $n \in \mathbb{Z}$ . Let  $C_*$  be any subcomplex of  $G_*$  such that  $D_n \subseteq C_n$  for each  $n \in \mathbb{Z}$ . Then

$$\partial_{n+1}^{G_*}(D_{n+1}) \leq \partial_{n+1}^{G_*}(C_{n+1}) \leq C_n$$

and so

$$\tilde{D}_n = D_n \cdot \partial_{n+1}^{G_*}(D_{n+1}) \leq C_n.$$

Thus  $\tilde{D}_*$  is a graded subgroup of  $C_*$  for any subcomplex  $C_*$  of  $G_*$  with  $D_n \leq C_n$  for  $n \in \mathbb{Z}$ , and so  $\tilde{D}_*$  is a graded subgroup of  $\text{sup}_*(D_*)$ . Notice that

$$\begin{aligned} \partial_n^{G_*}(\tilde{D}_*) &= \partial_n^{G_*}(D_n \cdot \partial_{n+1}^{G_*}(D_{n+1})) \\ &\leq \partial_n^{G_*}(D_n) \cdot \partial_n^{G_*}(\partial_{n+1}^{G_*}(D_{n+1})) \\ &= \partial_n^{G_*}(D_n) \\ &\leq \tilde{D}_{n-1}. \end{aligned}$$

Hence  $\tilde{D}_*$  is a subcomplex of  $G_*$  containing  $D_*$ , and so  $\text{sup}_*(D_*) = \tilde{D}_*$ .

(2) The first part follows from the definition. Let

$$\check{D}_n = D_n \cap \partial_n^{-1}(D_{n-1}).$$

Let  $x \in \check{D}_n$ . Then  $\partial_n(x) \in D_{n-1}$  because  $x \in \partial_n^{-1}(D_{n-1})$ , and

$$\partial_n(x) \in \partial_{n-1}^{-1}(D_{n-2})$$

because  $\partial_{n-1}(\partial_n(x)) = 1 \in D_{n-2}$ . Thus  $\partial_n(x) \in \partial_{n-1}^{-1}(D_{n-2})$ . It follows that  $\check{D}_*$  is a subcomplex of  $G_*$  contained in  $D_*$ . Hence

$$\check{D}_* \leq \text{inf}_*(D_*).$$

Let  $E_*$  be any subcomplex of  $G_*$  such that  $E_n \leq D_n$  for  $n \in \mathbb{Z}$ . Then

$$E_n \leq \partial_n^{-1}(E_{n-1}) \leq \partial_n^{-1}(D_{n-1}).$$

Thus  $E_n \leq \check{D}_n$  for  $n \in \mathbb{Z}$ . It follows that  $\text{inf}_*(D_*) = \check{D}_*$ . This finishes the proof.  $\square$

Let  $G_*$  be a chain complex of groups. The homology of  $G_*$  is defined as the right cosets

$$H_n(G_*) = \text{Ker}(\partial_n^{G_*}) / \partial_{n+1}^{G_*}(G_{n+1}).$$



**Proposition 2.3.** Let  $G_*$  be a chain complex of groups and let  $D_*$  be a graded subgroup of  $G_*$ .

1. The inclusion

$$\text{inf}_*(D_*) \rightarrow \text{sup}_*(D_*)$$

induces an injective map on homology.

2. Suppose that  $\partial_{n+1}^{G_*}(D_{n+1})$  is contained in the normalizer of  $D_n$  in  $G_n$  for each  $n$ . Then the inclusion

$$\text{inf}_*(D_*) \rightarrow \text{sup}_*(D_*)$$

induces an isomorphism on homology. In particular, if  $D_n$  is normal in  $G_n$  for  $n \in \mathbb{Z}$ , then the inclusion  $\text{inf}_*(D_*) \rightarrow \text{sup}_*(D_*)$  induces an isomorphism on homology.

*Proof.* (1) From Proposition 2.2 (2),

$$H_n(\text{inf}_*(D_*)) = \left( D_n \cap \partial_n^{-1}(D_{n-1}) \cap \text{Ker}(\partial_n^{G_*}) \right) / \partial_{n+1} \left( D_{n+1} \cap \partial_{n+1}^{-1}(D_n) \right)$$

as right cosets. Since  $\text{Ker}(\partial_n^{G_*}) \leq \partial_n^{-1}(D_{n-1})$ , we have

$$D_n \cap \partial_n^{-1}(D_{n-1}) \cap \text{Ker}(\partial_n^{G_*}) = D_n \cap \text{Ker}(\partial_n^{G_*}).$$

We also claim that

$$\partial_{n+1} \left( D_{n+1} \cap \partial_{n+1}^{-1}(D_n) \right) = D_n \cap \partial_{n+1}(D_{n+1}).$$

Clearly,  $\partial_{n+1} \left( D_{n+1} \cap \partial_{n+1}^{-1}(D_n) \right) \leq D_n \cap \partial_{n+1}(D_{n+1})$ .

Let  $x \in D_n \cap \partial_{n+1}(D_{n+1})$  and let  $y \in D_{n+1}$  such that  $\partial_{n+1}(y) = x$ . Then  $y \in D_{n+1} \cap \partial_{n+1}^{-1}(D_n)$ .

Thus  $x \in \partial_{n+1} \left( D_{n+1} \cap \partial_{n+1}^{-1}(D_n) \right)$ . Hence  $\partial_{n+1} \left( D_{n+1} \cap \partial_{n+1}^{-1}(D_n) \right) = D_n \cap \partial_{n+1}(D_{n+1})$

and so

$$H_n(\text{inf}_*(D_*)) = \left( D_n \cap \text{Ker}(\partial_n^{G_*}) \right) / \left( D_n \cap \partial_{n+1}(D_{n+1}) \right).$$

From Proposition 2.2 (1),

$$H_n(\text{sup}_*(D_*)) = \left( \left( D_n \cdot \partial_{n+1}^{G_*}(D_{n+1}) \right) \cap \text{Ker}(\partial_n^{G_*}) \right) / \partial_{n+1}^{G_*} \left( D_{n+1} \cdot \partial_{n+2}^{G_*}(D_{n+2}) \right).$$

Since  $\partial_{n+1}^{G^*}(\partial_{n+2}^{G^*}(D_{n+2})) = \{1\}$ ,  $\partial_{n+1}^{G^*}(D_{n+1} \cdot \partial_{n+2}^{G^*}(D_{n+2})) = \partial_{n+1}^{G^*}(D_{n+1})$ . Thus

$$H_n(\text{sup}_*(D_*)) = \left( (D_n \cdot \partial_{n+1}^{G^*}(D_{n+1})) \cap \text{Ker}(\partial_n^{G^*}) \right) / \partial_{n+1}^{G^*}(D_{n+1}).$$

Let  $w_1, w_2 \in D_n \cap \text{Ker}(\partial_n^{G^*})$  such that  $w_1 \equiv w_2$  in  $H_n(\text{sup}_*(D_*))$ . Then there exists  $y = \partial_{n+1}^{G^*}(D_{n+1})$  such that  $w_2 = w_1 y$ . Note that

$$y = w_1^{-1} w_2 \in D_n \cap \text{Ker}(\partial_n^{G^*}) \leq D_n.$$

We have  $y \in D_n \cap \partial_{n+1}^{G^*}(D_{n+1})$  with  $w_2 = w_1 y$ . Thus  $w_1 \equiv w_2$  in  $H_n(\text{inf}_*(D_*))$ . So

$$H_n(\text{inf}_*(D_*)) \rightarrow H_n(\text{sup}_*(D_*))$$

is injective. This proves (1).

(2) Let  $w \in (D_n \cdot \partial_{n+1}^{G^*}(D_{n+1})) \cap \text{Ker}(\partial_n^{G^*})$ . Then  $w \in D_n \cdot \partial_{n+1}^{G^*}(D_{n+1})$  and so

$$w = x_1 y_1 x_2 y_2 \cdots x_m y_m$$

with  $x_i \in D_n$  and  $y_i \in \partial_{n+1}^{G^*}(D_{n+1})$  for  $1 \leq i \leq m$ . Since  $\partial_{n+1}^{G^*}(D_{n+1})$  is contained in the normalizer of  $D_n$ , the product

$$w = x_1 (y_1 x_2 y_1^{-1}) (y_1 y_2 x_3 y_2^{-1} y_1^{-1}) \cdots (y_1 \cdots y_{m-1} x_m y_{m-1}^{-1} \cdots y_1^{-1}) y_1 \cdots y_m = xy$$

with

$$x = x_1 (y_1 x_2 y_1^{-1}) (y_1 y_2 x_3 y_2^{-1} y_1^{-1}) \cdots (y_1 \cdots y_{m-1} x_m y_{m-1}^{-1} \cdots y_1^{-1}) \in D_n$$

and

$$y = y_1 y_2 \cdots y_m \in \partial_{n+1}^{G^*}(D_{n+1}).$$

Since  $y \in \text{Ker}(\partial_n^{G^*})$  and  $w \in \text{Ker}(\partial_n^{G^*})$ ,

$$x = w y^{-1} \in \text{Ker}(\partial_n^{G^*}).$$

It follows that  $x \in D_n \cap \text{Ker}(\partial_n^{G^*})$ , and so

$$H_n(\inf_*(D_*)) \rightarrow H_n(\sup_*(D_*))$$

is surjective. From (1),  $H_n(\inf_*(D_*)) \rightarrow H_n(\sup_*(D_*))$  is injective and so it is an isomorphism. This finishes the proof.  $\square$

## 2.2. Hypergraphs.

Recall that a *hypergraph*  $\mathcal{H}$  is a pair  $\mathcal{H} = (V_{\mathcal{H}}, \mathcal{E}_{\mathcal{H}})$ , where the vertex set  $V_{\mathcal{H}}$  is a finite or infinite set and the hyperedge set  $\mathcal{E}_{\mathcal{H}}$  is a collection of finite nonempty subsets of  $V_{\mathcal{H}}$ . Let  $\mathcal{P}(V_{\mathcal{H}})$  be the set of all finite subsets of  $V_{\mathcal{H}}$ . The hypothesis in the definition of hypergraph  $\mathcal{H} = (V_{\mathcal{H}}, \mathcal{E}_{\mathcal{H}})$  only requires that  $\mathcal{E}_{\mathcal{H}} \subseteq \mathcal{P}(V_{\mathcal{H}}) \setminus \emptyset$ . This is different from the notion of an abstract simplicial complex as hypergraphs do not require  $\mathcal{E}_{\mathcal{H}}$  to be closed under taking subsets.

The *simplicial closure* (or the *associated simplicial complex* as in [110]) of a hypergraph  $\mathcal{H} = (V_{\mathcal{H}}, \mathcal{E}_{\mathcal{H}})$ , denoted by  $\Delta\mathcal{H}$ , is defined as

$$\Delta\mathcal{H} = \{A \neq \emptyset \mid A \subseteq B \text{ for some } B \in \mathcal{E}_{\mathcal{H}}\}.$$

It is straightforward to check that the simplicial closure of  $\mathcal{H}$  is the minimal simplicial complex containing  $\mathcal{H}$ . The homology of  $\Delta\mathcal{H}$  has been studied previously in [110]. However, it is desirable for a homology theory of  $\mathcal{H}$  to be directly derived from  $\mathcal{H}$  itself rather than the simplicial closure  $\Delta\mathcal{H}$ . Using Proposition 2.3, there is an embedded homology theory of hypergraphs that is an extension of simplicial homology theory.

**Definition 2.4.** Let  $\mathcal{H} = (V_{\mathcal{H}}, \mathcal{E}_{\mathcal{H}})$  be a hypergraph with a total ordering on  $V_{\mathcal{H}}$  and let  $G$  be an abelian group. Let  $C_*(\Delta\mathcal{H}; G)$  be the chain complex with coefficients in group  $G$ . Consider  $\mathbb{Z}(\mathcal{H}) \otimes G$  as a graded subgroup of the chain complex of abelian groups  $C_*(\Delta\mathcal{H}; G)$ . The *embedded homology*  $H_*^{\text{emb}}(\mathcal{H}; G)$  with coefficients in  $G$  is defined by

$$H_*^{\text{emb}}(\mathcal{H}; G) = H_*\left(\inf_*^{C_*(\Delta\mathcal{H}; G)}(\mathbb{Z}(\mathcal{H}) \otimes G)\right) \cong H_*\left(\sup_*^{C_*(\Delta\mathcal{H}; G)}(\mathbb{Z}(\mathcal{H}) \otimes G)\right).$$

The crucial point is that by Proposition 2.3 the inclusion

$$\inf_*^{C_*(\Delta\mathcal{H}; G)}(\mathbb{Z}(\mathcal{H}) \otimes G) \rightarrow \sup_*^{C_*(\Delta\mathcal{H}; G)}(\mathbb{Z}(\mathcal{H}) \otimes G)$$

induces an isomorphism on homology. Hence this homology can be considered as a natural topological invariant of  $\mathcal{H}$ . To detect more subtle information about  $\mathcal{H}$ , one could explore the acyclic chain complex

$$\sup_* C_*(\Delta\mathcal{H}; G)_{(\mathbb{Z}(\mathcal{H}) \otimes G)} / \inf_* C_*(\Delta\mathcal{H}; G)_{(\mathbb{Z}(\mathcal{H}) \otimes G)}.$$

For example, when  $G$  is a field, one can investigate the Hilbert-Poincaré series

$$\xi^{\text{emb}}(\mathcal{H}, t) = \sum_{n=0}^{\infty} \left( \dim \left( \sup_* C_*(\Delta\mathcal{H}; G)_{(\mathbb{Z}(\mathcal{H}) \otimes G)} / \inf_* C_*(\Delta\mathcal{H}; G)_{(\mathbb{Z}(\mathcal{H}) \otimes G)} \right) \right) t^n$$

to detect gaps and get more robust information.

Let  $\delta(\mathcal{H})$  denote the maximal simplicial complex contained in  $\mathcal{H}$ . In general,  $H_*^{\text{emb}}(\mathcal{H}; G)$  is different from  $H_*(\delta(\mathcal{H}); G)$  and  $H_*(\Delta(\mathcal{H}); G)$  as shown in the following example.

**Example 2.5.** Let  $\mathcal{H}$  be the boundary of a 2-simplex with all vertices removed,  $V_{\mathcal{H}} = \{0, 1, 2\}$  and  $\mathcal{E}_{\mathcal{H}} = \{\{0, 1\}, \{0, 2\}, \{1, 2\}\}$  as depicted in Figure 4. Then  $\delta(\mathcal{H})$  is the empty set, and  $\Delta(\mathcal{H})$  is the boundary of the 2-simplex. By definition,  $H_1^{\text{emb}}(\mathcal{H}; \mathbb{Z}) = \mathbb{Z}$  and  $H_0^{\text{emb}}(\mathcal{H}; \mathbb{Z}) = 0$ . Thus  $H_*^{\text{emb}}(\mathcal{H}; \mathbb{Z})$  is different from  $H_*(\delta(\mathcal{H}); \mathbb{Z})$  and  $H_*(\Delta(\mathcal{H}); \mathbb{Z})$ .

This example shows that  $H_*^{\text{emb}}(\mathcal{H}; G)$  may not be the homology of any simplicial complex as  $H_0^{\text{emb}}(\mathcal{H}; \mathbb{Z}) = 0$ , which is not the case for any nonempty simplicial complex. Let us consider another example.

**Example 2.6.** Let  $\mathcal{H} = (V_{\mathcal{H}}, \mathcal{E}_{\mathcal{H}})$  with  $V_{\mathcal{H}} = \{0, 1, 2\}$  ordered by  $0 < 1 < 2$ , and

$$\mathcal{E}_{\mathcal{H}} = \{\{0, 1, 2\}, \{0, 1\}, \{0, 2\}, \{0\}, \{1\}, \{2\}\}$$

see Figure 4. Then  $\Delta\mathcal{H}$  is the abstract simplicial complex of a 2-simplex with vertices labeled by 0, 1, 2. The 1-face  $\{1, 2\}$  is not in  $\mathcal{H}$ . Let  $G = \mathbb{Z}$ . Then the chain complex  $C_*(\Delta\mathcal{H})$  is given by  $C_0(\Delta\mathcal{H}) = \mathbb{Z}^{\oplus 3} = \mathbb{Z}\{\{0\}, \{1\}, \{2\}\}$ ,  $C_1(\Delta\mathcal{H}) = \mathbb{Z}^{\oplus 3} = \mathbb{Z}\{\{0, 1\}, \{0, 2\}, \{1, 2\}\}$ , and  $C_2(\Delta\mathcal{H}) = \mathbb{Z} = \mathbb{Z}\{\{0, 1, 2\}\}$ .

We have  $\inf_0 = C_0(\Delta\mathcal{H}) = \mathbb{Z}\{\{0\}, \{1\}, \{2\}\}$ ,

$$\inf_1 = \mathbb{Z}(\mathcal{E}_1) \cap \partial_1^{-1}(\mathbb{Z}(\mathcal{E}_0)) = \mathbb{Z}(\mathcal{E}_1) \cap C_1(\Delta\mathcal{H}) = \mathbb{Z}(\mathcal{E}_1) = \mathbb{Z}\{\{0, 1\}, \{0, 2\}\}$$

$$\inf_2 = \mathbb{Z}(\mathcal{E}_2) \cap \partial_2^{-1}(\mathbb{Z}(\mathcal{E}_1)) = 0$$

with  $\partial_1(\inf_1) = \mathbb{Z}\{\{1\} - \{0\}, \{2\} - \{0\}\}$ . Thus  $H_0^{\text{emb}}(\mathcal{H}) = \mathbb{Z}$  and  $H_i^{\text{emb}}(\mathcal{H}) = 0$  for  $i = 1$ .

Let  $\mathcal{H}' = (V_{\mathcal{H}'}, \mathcal{E}_{\mathcal{H}'})$  with  $V_{\mathcal{H}'} = V_{\mathcal{H}} = \{0, 1, 2\}$  ordered by  $0 < 1 < 2$ , and

$$\mathcal{E}_{\mathcal{H}'} = \{\{0, 1, 2\}, \{0, 1\}, \{0\}, \{1\}, \{2\}\}.$$

Then  $H_0^{\text{emb}}(\mathcal{H}') = \mathbb{Z} \oplus \mathbb{Z}$  and  $H_i^{\text{emb}}(\mathcal{H}') = 0$  for  $i \geq 1$ . Thus the embedded homology of  $\mathcal{H}'$  can not be realized as the homology of a path-connected topological space.

These examples indicate that embedded homology is a new homology theory with unusual properties and that poses its own questions and challenges.

The definition of embedded homology of a hypergraph  $\mathcal{H}$  depends on the orientation of its simplicial closure  $\Delta\mathcal{H}$ . It is well-known that simplicial homology is independent on the choice of orientation. The following theorem shows that this is also true for the embedded homology of hypergraphs.

**Theorem 2.7.** *The embedded homology  $H_*^{\text{emb}}(\mathcal{H}; G)$  of a hypergraph  $\mathcal{H}$  does not depend on a choice of orientation on  $\Delta\mathcal{H}$ .*

*Proof.* Let  $H_*^{\text{emb}}(\mathcal{H})$  and  $G(\mathcal{H})$  denote  $H_*^{\text{emb}}(\mathcal{H}, G)$  and  $\mathbb{Z}(\mathcal{H}) \otimes G$ , respectively. We assume that  $V_{\mathcal{H}}$  is a finite set  $\{v_1, v_2, \dots, v_m\}$ . Take a linear ordering on  $V_{\mathcal{H}}$  so that  $v_1 < v_2 < \dots < v_m$  as a fixed choice of total order and let  $C_* = C_*(\Delta\mathcal{H}; G)$  denote the oriented chain complex. It suffices to show that the homology stays the same up to isomorphism under the transpositions  $(i, i + 1)$  of the ordering on  $V(\mathcal{H})$  for  $1 \leq i \leq m - 1$ .

Let  $\partial'_n: C_n \rightarrow C_{n-1}$ ,  $n \geq 1$  be the boundary homomorphism defined using the new order on  $V_{\mathcal{H}}$ , that is,  $v_1 < v_2 < \dots < v_{i-1} < v_{i+1} < v_i < v_{i+2} < \dots < v_m$ . For  $n \geq 1$ , the abelian group  $C_n$  admits a direct sum decomposition

$$C_n = C_n^{v_i v_{i+1}} \oplus C_n^{\widehat{v_i v_{i+1}}} \tag{1}$$

where  $C_n^{v_i v_{i+1}}$  is the subgroup of  $C_n$  given by linear combinations with coefficients in group  $G$  of the  $n$ -simplices  $\sigma \in \Delta\mathcal{H}$  whose vertex set contains both  $v_i$  and  $v_{i+1}$ , and  $C_n^{\widehat{v_i v_{i+1}}}$  is the subgroup of  $C_n$  given by linear combinations with coefficients in group  $G$  of the remaining  $n$ -simplices in  $\Delta\mathcal{H}$ . For any chain  $\alpha \in C_n$  there is a corresponding unique decomposition

$$\alpha = \alpha^{v_i v_{i+1}} + \alpha^{\widehat{v_i v_{i+1}}}. \tag{2}$$

Since  $v_i$  and  $v_{i+1}$  are neighbored vertices in the order, we have  $\partial'_n(\sigma) = \partial_n(\sigma)$  if  $\sigma$  does not contain both  $v_i$  and  $v_{i+1}$  in its vertex set.

Therefore

$$\partial'_n|_{C_n^{\widehat{v_i v_{i+1}}}} = \partial_n|_{C_n^{\widehat{v_i v_{i+1}}}} \rightarrow C_{n-1}. \tag{3}$$

Let  $\sigma = [a_1 \cdots a_t v_j v_{j+1} b_1 \cdots b_s]$  be an oriented simplex in  $\Delta \mathcal{X}$  with  $a_1 < \cdots < a_t < v_j < v_{j+1} < b_1 < \cdots < b_s$ . By the definition of  $(\sigma)$ , we have

$$\begin{aligned} \partial(\sigma)^{v_i v_i + 1} &= \sum_{j=1}^t (-1)^{j-1} [a_1 \cdots \widehat{a_j} \cdots a_t v_j v_{j+1} b_1 \cdots b_s] + \\ &\sum_{k=1}^s (-1)^{t+k+1} [a_1 \cdots a_t v_i v_{i+1} b_1 \cdots \widehat{b_k} \cdots b_s] \end{aligned}$$

where  $\cdots \widehat{x} \cdots$  means that  $x$  is deleted, and

$$\partial(\sigma)^{\widehat{v_i v_i + 1}} = (-1)^t [a_1 \cdots a_t v_{i+1} b_1 \cdots b_s] + (-1)^{t+1} [a_1 \cdots a_t v_i b_1 \cdots b_s].$$

By switching the order of  $v_j$  and  $v_{j+1}$ , we have

$$\partial'(\sigma) = (\partial(\sigma))^{v_i v_i + 1} - (\partial(\sigma))^{\widehat{v_i v_i + 1}}.$$

Extending this formula linearly with coefficients in group  $G$ , we obtain the formula

$$\partial'(\alpha) = (\partial(\alpha))^{v_i v_i + 1} - (\partial(\alpha))^{\widehat{v_i v_i + 1}} \text{ for } \alpha \in C_*^{v_i v_i + 1}. \tag{4}$$

Define the group homomorphism

$$\phi_n: C_n = C_n^{v_i v_i + 1} \oplus C_n^{\widehat{v_i v_i + 1}} \rightarrow C_n = C_n^{v_i v_i + 1} \oplus C_n^{\widehat{v_i v_i + 1}}$$

by setting

$$\phi_n(z^{v_i v_i + 1} + z^{\widehat{v_i v_i + 1}}) = z^{v_i v_i + 1} - z^{\widehat{v_i v_i + 1}}.$$

Clearly,  $\phi_n$  is an isomorphism. Let  $z = z^{v_i v_i + 1} + z^{\widehat{v_i v_i + 1}} \in C_n$  be a chain. Then

$$\begin{aligned} \partial(z) &= \partial(z^{v_i v_i + 1}) + \partial(z^{\widehat{v_i v_i + 1}}) \\ &= (\partial(z^{v_i v_i + 1}))^{v_i v_i + 1} + (\partial(z^{v_i v_i + 1}))^{\widehat{v_i v_i + 1}} + \partial(z^{\widehat{v_i v_i + 1}}) \end{aligned}$$

so

$$(\partial(z))^{v_i v_i + 1} = (\partial(z^{v_i v_i + 1}))^{v_i v_i + 1}$$

and

$$(\partial(z))^{\widehat{v_i v_i + 1}} = (\partial(z^{v_i v_i + 1}))^{\widehat{v_i v_i + 1}} + \partial(z^{\widehat{v_i v_i + 1}}).$$

On the other hand, by direct computation

$$\begin{aligned} \partial'(\phi_n(z)) &= \partial(z^{v_i v_i + 1} - z^{\widehat{v_i v_i + 1}}) \\ &= \partial(z^{v_i v_i + 1})^{v_i v_i + 1} - \left( (\partial(z^{v_i v_i + 1}))^{\widehat{v_i v_i + 1}} + \partial(z^{\widehat{v_i v_i + 1}}) \right) \\ &= (\partial(z))^{v_i v_i + 1} - (\partial(z))^{\widehat{v_i v_i + 1}}. \end{aligned}$$

This gives a commutative diagram

$$\begin{array}{ccc} C_n & \xrightarrow[\cong]{\phi_n} & C_n \\ \downarrow \partial_n & & \downarrow \partial'_n \\ C_{n-1} & \xrightarrow[\cong]{\phi_{n-1}} & C_{n-1}. \end{array}$$

Note that the decomposition (1) restricted to  $G(\mathcal{H}_n)$  gives the decomposition

$$G(\mathcal{H}) = G(\mathcal{H})^{v_i v_i + 1} \oplus G(\mathcal{H})^{\widehat{v_i v_i + 1}}$$

with the same rule on simplices. The subgroup  $G(H_n)$  is invariant under  $\phi_n$ . Moreover, there is a commutative diagram

$$\begin{array}{ccc} \inf_n^{C_*(\Delta\mathcal{H};G)}(G(\mathcal{H})) & \xrightarrow[\cong]{\phi_n} & \inf_n^{C_*(\Delta\mathcal{H};G)}(G(\mathcal{H})) \\ \downarrow \partial_n| & & \downarrow \partial'_n| \\ \inf_{n-1}^{C_*(\Delta\mathcal{H};G)}(G(\mathcal{H})) & \xrightarrow[\cong]{\phi_{n-1}} & \inf_{n-1}^{C_*(\Delta\mathcal{H};G)}(G(\mathcal{H})). \end{array}$$

The assertion then follows by taking homology of this commutative diagram.  $\square$

The embedded homology of hypergraphs was introduced in 2019 in [23]. Previously, cohomological aspects on  $k$ -uniform hypergraphs have been studied, see [28, 27, 42, 44, 96, 99, 122, 123, 138, 142], using cohomology introduced in a combinatorial way. Also

Emtander [59] studied the homology of the *independence complex*  $\Delta^c \mathcal{H}$  of a hypergraph  $\mathcal{H} = (V_{\mathcal{H}}, \mathcal{E}_{\mathcal{H}})$  in 2009, where  $\Delta^c \mathcal{H} = \{F \subseteq V_{\mathcal{H}} \mid E \not\subseteq F \text{ for any } E \in \mathcal{E}_{\mathcal{H}}\}$ . The approach of embedded homology is different from the classical research on topological structures related to hypergraphs as it is directly define on the hypergraph.

Although the embedded homology of hypergraphs is a new topic with surprising properties, it inherits many characteristics of simplicial homology. The following theorem is an example of this, the proof of this theorem is similar to that of [102, Theorem 8.2, p.45].

Recall that the *cone*  $CK$  of a simplicial complex  $K$  is defined as a join  $CK = w * K$  with  $w$  a vertex not in  $K$ . Analogously, we can define the join of hypergraphs and the cone  $C\mathcal{H} = w * \mathcal{H}$ .

**Theorem 2.8.** *Let  $\mathcal{H}$  be a hypergraph and let  $G$  be an abelian group. Then*

$$H_n^{\text{emb}(C\mathcal{H}; G)} = \begin{cases} 0 & \text{if } n > 0 \\ G & \text{if } n = 0. \end{cases}$$

□

### 2.3. Super-hypergraphs.

Recall [47, 139] that a  $\Delta$ -set  $X_*$  is a sequence of sets  $X_* = (X_n)_{n \geq 0}$  with maps  $d_i: X_n \rightarrow X_{n-1}$ , for  $0 \leq i < n$  and  $n \geq 1$ , called *face operations*, satisfying the following  $\Delta$ -identity

$$d_i d_j = d_j d_{i+1} \text{ for } i \geq j. \tag{5}$$

**Definition 2.9.** A *super-hypergraph* is a pair  $(\mathcal{H}, X)$ , where  $X$  is a  $\Delta$ -set and  $\mathcal{H}$  is a graded subset of  $X$ . We call  $\mathcal{H}$  a *super-hypergraph born from  $X$* , and  $X$  is called a *parental  $\Delta$ -set* of  $\mathcal{H}$ . The  $\Delta$ -closure of  $\mathcal{H}$  in  $X$  is defined by

$$\Delta^X(\mathcal{H}) = \bigcap \{Y \mid \mathcal{H} \subseteq Y \text{ as a graded subset and } Y \subseteq X \text{ as a } \Delta\text{-subset}\}.$$

A *morphism*  $\phi: (\mathcal{H}, X) \rightarrow (\mathcal{H}', Y)$  of super-hypergraphs is a  $\Delta$ -map  $\phi: X \rightarrow Y$  such that  $\phi(\mathcal{H}) \subseteq \mathcal{H}'$ .

**2.3.1. Homology of super-hypergraphs.**—Using Proposition 2.3, there is an embedded homology on super-hypergraphs.

**Definition 2.10.** Let  $(\mathcal{H}, X)$  be a super-hypergraph and let  $G$  be an abelian group. The *embedded homology*  $H_*^{\text{emb}, X}(\mathcal{H}; G)$  with coefficients in  $G$  of  $(\mathcal{H}, X)$  is defined by

$$H_*^{\text{emb}, X}(\mathcal{H}; G) = H_* \left( \inf_*^{C^*(X; G)} (\mathbb{Z}(\mathcal{H}) \otimes G) \right) \cong H_* \left( \sup_*^{C^*(X; G)} (X; G)(\mathbb{Z}(\mathcal{H}) \otimes G) \right)$$



where  $\mathbb{Z}(\mathcal{H}) \otimes G$  is a graded subgroup of the chain complex of abelian groups  $C_*(X; G)$ .

We want to show that this definition is an extension of the embedded homology of hypergraphs. An *oriented hypergraph* is a hypergraph  $\mathcal{H}$  with a partial order on its vertex set so that the restriction to the vertices of each hyperedge of  $\mathcal{H}$  is linear. If the vertices of a simplex are totally ordered, then the restricted order on the vertices of any of its faces is linear. Thus the simplicial closure  $\Delta\mathcal{H}$  can be oriented with its orientation induced by the order on  $\mathcal{H}$ . From Definition 2.4,

$$H_*^{\text{emb}}(\mathcal{H}; G) = H_*^{\text{emb}, \Delta\mathcal{H}}(\mathcal{H}; G)$$

by considering the oriented simplicial complex  $\Delta\mathcal{H}$  as a  $\mathbb{Z}$ -set. From Theorem 2.7, this definition is independent on the choice of orientation.

We now consider how morphisms of super-hypergraphs, recall Definition 2.9, induce maps on the infimum and supremum chain complexes as well as embedded homology of super-hypergraphs.

**Proposition 2.11.** *Let  $\phi: (\mathcal{H}, X) \rightarrow (\mathcal{H}', Y)$  be a morphism of super-hypergraphs. Then there is a commutative diagram*

$$\begin{array}{ccccccc}
 \text{inf}_*^{C_*(X; G)}(\mathbb{Z}(\mathcal{H}) \otimes G) & \subseteq & \mathbb{Z}(\mathcal{H}) \otimes G & \subseteq & \text{sup}_*^{C_*(X; G)}(\mathbb{Z}(\mathcal{H}) \otimes G) & \hookrightarrow & C_*(X; G) \\
 \downarrow \phi_{\#1} & & \downarrow \phi_{\#1} & & \downarrow \phi_{\#1} & & \downarrow \phi_{\#} \\
 \text{inf}_*^{C_*(Y; G)}(\mathbb{Z}(\mathcal{H}') \otimes G) & \subseteq & \mathbb{Z}(\mathcal{H}') \otimes G & \subseteq & \text{sup}_*^{C_*(Y; G)}(\mathbb{Z}(\mathcal{H}') \otimes G) & \hookrightarrow & C_*(Y; G)
 \end{array}
 \tag{6}$$

which induces a map  $\phi_*: H_*^{\text{emb}, X}(\mathcal{H}; G) \rightarrow H_*^{\text{emb}, Y}(\mathcal{H}'; G)$ . Moreover, if  $\phi(\mathcal{H}) = \mathcal{H}'$ , then there is a short exact sequence of chain complexes

$$\begin{array}{c}
 \text{sup}_*^{C_*(X; G)}(\mathbb{Z}(\mathcal{H}) \otimes G) \cap \text{Ker}(\phi_{\#}) \hookrightarrow \text{sup}_*^{C_*(X; G)}(\mathbb{Z}(\mathcal{H}) \otimes G) \\
 \downarrow \phi_{\#} \\
 \text{sup}_*^{C_*(Y; G)}(\mathbb{Z}(\mathcal{H}') \otimes G)
 \end{array}
 \tag{7}$$

*Proof.* Since  $\text{sup}_*^{C_*(Y; G)}(\mathbb{Z}(\mathcal{H}') \otimes G)$  is a subcomplex of  $C_*(Y; G)$  containing  $\mathbb{Z}(\mathcal{H}') \otimes G$ , its preimage

$$\phi_{\#}^{-1}\left(\text{sup}_*^{C_*(Y; G)}(\mathbb{Z}(\mathcal{H}') \otimes G)\right)$$

is a subcomplex of  $C_*(X; G)$  containing  $\mathbb{Z}(\mathcal{H}) \otimes G$  because  $\phi(\mathcal{H}) \subseteq \mathcal{H}'$ . Thus

$$\text{sup}_*^{C_*(X; G)}(\mathbb{Z}(\mathcal{H}) \otimes G) \subseteq \phi_{\#}^{-1}\left(\text{sup}_*^{C_*(Y; G)}(\mathbb{Z}(\mathcal{H}') \otimes G)\right) \tag{8}$$

as a subcomplex. Similarly

$$\phi_{\#} \left( \inf_*^{C_*(X; G)} (\mathbb{Z}(\mathcal{H}) \otimes G) \right) \subseteq \inf_*^{C_*(Y; G)} (\mathbb{Z}(\mathcal{H}') \otimes G) \tag{9}$$

as a subcomplex. Therefore, there is a commutative diagram 6.

Now, we assume that  $\phi(\mathcal{H}) = \mathcal{H}'$ . Then

$$\phi_{\#} \left( \sup_*^{C_*(X; G)} ((\mathbb{Z}(\mathcal{H}) \otimes G)) \right) \supseteq \phi_{\#} ((\mathbb{Z}(\mathcal{H}) \otimes G)) = \mathbb{Z}(\mathcal{H}') \otimes G$$

is a subcomplex of  $C_*(Y; G)$  containing  $\mathbb{Z}(\mathcal{H}') \otimes G$ . Hence

$$\phi_{\#} \left( \sup_*^{C_*(X; G)} ((\mathbb{Z}(\mathcal{H}) \otimes G)) \right) \supseteq \sup_*^{C_*(Y; G)} ((\mathbb{Z}(\mathcal{H}') \otimes G)).$$

Together with the containment 8, we have

$$\sup_*^{C_*(Y; G)} ((\mathbb{Z}(\mathcal{H}') \otimes G)) \supseteq \phi_{\#} \left( \sup_*^{C_*(X; G)} ((\mathbb{Z}(\mathcal{H}) \otimes G)) \right)$$

and hence the short exact sequence 7.  $\square$

**Corollary 2.12.** *Let  $\phi: (\mathcal{H}, X) \rightarrow (\mathcal{H}', Y)$  be a morphism of super-hypergraphs. Suppose that*

1.  $\phi: X \rightarrow Y$  is an injective  $\mathbb{Z}$ -map and
2.  $\phi(\mathcal{H}) = \mathcal{H}'$ .

Then

$$\phi_{\#}: \sup_*^{C_*(X; G)} ((\mathbb{Z}(\mathcal{H}) \otimes G)) \rightarrow \sup_*^{C_*(Y; G)} ((\mathbb{Z}(\mathcal{H}') \otimes G))$$

is an isomorphism. In particular,  $\phi_*: H_*^{\text{emb}, X}(\mathcal{H}; G) \rightarrow H_*^{\text{emb}, Y}(\mathcal{H}'; G)$  is an isomorphism.

*Proof.* By the assumption (1),

$$\sup_*^{C_*(X; G)} ((\mathbb{Z}(\mathcal{H}) \otimes G) \cap \text{Ker}(\phi_{\#})) = 0$$

and so the assertion follows by Proposition 2.11.  $\square$

**2.3.2. Variations of parental  $\mathbb{Z}$ -sets.**—In recent topological applications to data analytics and machine learning, one of the most common approaches is to use discrete Hodge-Laplacian theory. Mathematically, combinatorial Laplacian operators defined on linear transformations on cochains of simplicial complexes have been studied, for example in [55, 82]. Therefore, in addition to simplicial homology, research on (co)chains of simplicial complexes such as spectral analysis on combinatorial Laplacian operators is also important for potential applications in data science. Similarly, the research on chains

$\inf_*^{C_*(X;G)}(\mathbb{Z}(\mathcal{H}) \otimes G)$  and  $\sup_*^{C_*(X;G)}(\mathbb{Z}(\mathcal{H}) \otimes G)$  could be useful for the applications of super-hypergraphs. Next, we describe some basic properties related to the chain complexes arising from super-hypergraphs.

A super-hypergraph  $\mathcal{H}$  is assumed to have a parental  $\mathbb{Z}$ -set  $X$  that carries geometric structural information about  $\mathcal{H}$ . The embedded homology  $H_*^{\text{emb}, X}(\mathcal{H}; G)$  is defined using the geometric information inherited from  $X$ . On level of chains, there are inclusions of graded groups

$$\inf_*^{C_*(X;G)}(\mathbb{Z}(\mathcal{H}) \otimes G) \hookrightarrow \mathbb{Z}(\mathcal{H}) \otimes G \hookrightarrow \sup_*^{C_*(X;G)}(\mathbb{Z}(\mathcal{H}) \otimes G) \hookrightarrow C_*(X;G)$$

where the right most inclusion is a chain map. By Proposition 2.3, the inclusion

$$\inf_*^{C_*(X;G)}(\mathbb{Z}(\mathcal{H}) \otimes G) \hookrightarrow \sup_*^{C_*(X;G)}(\mathbb{Z}(\mathcal{H}) \otimes G)$$

is a chain homotopy equivalence that defines the embedded homology. The gap complex

$$\sup_*^{C_*(X;G)}(\mathbb{Z}(\mathcal{H}) \otimes G) / \inf_*^{C_*(X;G)}(\mathbb{Z}(\mathcal{H}) \otimes G) \tag{10}$$

which is an acyclic chain complex, gives more robust information about the graded set  $\mathcal{H}$ .

Let  $\mathcal{H}$  be a fixed graded data set. Our aim is to vary the parental  $\mathbb{Z}$ -set  $X$  such that the corresponding infimum and supremum chain complexes reveal different aspects of the topological structure of  $\mathcal{H}$ . One natural way to vary the parental  $\mathbb{Z}$ -set would be to consider morphisms  $\phi: (\mathcal{H}, X) \rightarrow (\mathcal{H}, Y)$  that fix  $\mathcal{H}$  and investigate how these affect the embedded homology. Another important question is whether one can vary the parental  $\mathbb{Z}$ -set so that the gap complex (10) is as small as possible. We consider the following example.

**Example 2.13.** Let  $n$  be an odd positive integer. Let  $X = \mathbb{Z}^+[n]$  be the  $\mathbb{Z}$ -set induced by an  $n$ -simplex with vertices labelled  $0, 1, \dots, n$ . Let  $Y$  be the  $\mathbb{Z}$ -set with  $Y_k = \{a_k\}$ ,  $0 \leq k \leq n$ ,  $Y_k = \emptyset$  for  $k > n$  and  $d_i(a_k) = a_{k-1}$  for  $0 \leq i < k \leq n$ . Let  $\mathcal{H}$  be the graded set given by  $\mathcal{H}_n = \{x_n\}$  and  $\mathcal{H}_k = \emptyset$  for  $k < n$ . Consider the super-hypergraphs  $(\mathcal{H}, X)$ ,  $x_n = [0, 1, 2, \dots, n]$ , and  $(\mathcal{H}, Y)$ ,  $x_n = a_n$ . There is a unique morphism  $\phi: (\mathcal{H}, X) \rightarrow (\mathcal{H}, Y)$  such that  $\phi|_{\mathcal{H}} = \text{id}_{\mathcal{H}}$ . Then, we have  $H_k^{\text{emb}, X}(\mathcal{H}; \mathbb{Z}) = 0$  for  $k \neq 0$  and

$$H_k^{\text{emb}, Y}(\mathcal{H}; \mathbb{Z}) = \begin{cases} \mathbb{Z} & \text{if } k = n, \\ 0 & \text{otherwise.} \end{cases}$$

This shows the embedded homology of super-hypergraphs depends on the parental  $\mathbb{Z}$ -set. The gap complex for  $(\mathcal{H}, X)$  is the same as the acyclic complex

$$\sup_k^{C^*(X; \mathbb{Z})}(\mathbb{Z}(\mathcal{H})) = \begin{cases} \mathbb{Z} & \text{if } k = n, n - 1, \\ 0 & \text{otherwise} \end{cases}$$

with  $\partial: \sup_n^{C^*(X; \mathbb{Z})}(\mathbb{Z}(\mathcal{H})) \rightarrow \sup_{n-1}^{C^*(X; \mathbb{Z})}(\mathbb{Z}(\mathcal{H}))$ . an isomorphism, and the gap complex for  $(\mathcal{H}, Y)$  is 0.

Now we consider the effect of morphisms on super-hypergraphs with a fixed hypergraph more generally. Let  $\phi: (\mathcal{H}, X) \rightarrow (\mathcal{H}, Y)$  be a morphism of super-hypergraphs so that  $\phi|_{\mathcal{H}} = \text{id}_{\mathcal{H}}$ . By 9, we have

$$\inf_*^{C^*(X; G)}(\mathbb{Z}(\mathcal{H}) \otimes G) \hookrightarrow \inf_*^{C^*(Y; G)}(\mathbb{Z}(\mathcal{H}) \otimes G)$$

namely the chain complex of  $Y$  is closer to the graded group  $\mathbb{Z}(\mathcal{H}) \otimes G$  than the chain complex of  $X$ . Together with the inclusion 8, it follows that the gap complex for  $(\mathcal{H}, Y)$  is smaller, as shown in the above example. If  $\phi$  is injective, then both

$$\inf_*^{C^*(X; G)}(\mathbb{Z}(\mathcal{H}) \otimes G) = \inf_*^{C^*(Y; G)}(\mathbb{Z}(\mathcal{H}) \otimes G) \text{ and}$$

$$\sup_*^{C^*(X; G)}(\mathbb{Z}(\mathcal{H}) \otimes G) = \sup_*^{C^*(Y; G)}(\mathbb{Z}(\mathcal{H}) \otimes G)$$

which means that we can replace  $X$  by any of its  $\mathcal{H}$ -subsets that contain  $\mathcal{H}$  without changing the infimum and supremum chain complexes. In particular, we have

$$H_*^{\text{emb}, X}(\mathcal{H}; G) = H_*^{\text{emb}, \Delta^X(\mathcal{H})}(\mathcal{H}; G) \tag{11}$$

where  $\Delta^X(\mathcal{H})$  is the  $\mathcal{H}$ -closure of  $\mathcal{H}$  in  $X$ , which is the minimal  $\mathcal{H}$ -subset of  $X$  containing  $\mathcal{H}$  defined in Definition 2.9.

**Definition 2.14.** A super-hypergraph  $(\mathcal{H}, X)$  is called *regular* if  $X = \Delta^X(\mathcal{H})$ .

It is straightforward to see that a super-hypergraph  $(\mathcal{H}, X)$  is regular if and only if all elements in  $X$  are obtained from the elements in  $\mathcal{H}$  together with their iterated faces in  $X$ . The following proposition may be useful for analysing for variations of parental  $\mathcal{H}$ -sets.

**Proposition 2.15.** Let  $\mathcal{H} = \{\mathcal{H}_n\}_{n \geq 0}$  be a graded set such that the cardinality of the set  $\bigsqcup_{n \geq 0} H_n$  is finite. Then there are finitely many regular super-hypergraphs  $(\mathcal{H}, X)$  up to isomorphisms.

*Proof.* Consider the collection of all possible  $\mathcal{H}$ -sets  $X$  such that  $(\mathcal{H}, X)$  is a regular super-hypergraph. If  $(\mathcal{H}, X)$  is regular, then all elements in  $X$  are given by the elements in the

graded subset  $\mathcal{H}$  together with their iterated faces in  $X$ . Therefore, as a  $\mathbb{Z}$ -set  $X$  is a finite extension of  $\mathcal{H}$  together with finitely many face operations.  $\square$

**Definition 2.16.** A super-hypergraph  $(\mathcal{H}, X)$  is *complete* if it is regular, and for any morphism of super-hypergraphs  $\phi: (\mathcal{H}, X) \rightarrow (\mathcal{H}, Y)$  with  $\phi|_{\mathcal{H}} = \text{id}_{\mathcal{H}}$ , the  $\mathbb{Z}$ -map  $\phi: X \rightarrow Y$  is injective.

Therefore, for a complete super-hypergraph  $(\mathcal{H}, X)$ , the infimum complex

$$\inf_*^{C^*(X; G)} (\mathbb{Z}(\mathcal{H}) \otimes G)$$

reaches a maximum and the gap complex (10) reaches a minimum.

By definition, a complete super-hypergraph is regular. However, the converse may not be true. For instance, the super-hypergraph  $(\mathcal{H}, X)$  in Example 2.13 is regular but not complete as the morphism to  $(\mathcal{H}, Y)$  is not an injective  $\mathbb{Z}$ -set map. Therefore completeness provides a notion of maximality in the set of super-hypergraphs related to a given hypergraph.

**Theorem 2.17** (Completeness Criterion). *A regular super-hypergraph  $(\mathcal{H}, X)$  with  $\mathcal{H} \neq \emptyset$  is complete if and only if it has the following properties:*

1. (Vertex Property) *If  $\mathcal{H}_0 \neq \emptyset$ , then  $X_0 = \mathcal{H}_0$ . If  $\mathcal{H}_0 = \emptyset$ , then  $X_0$  is a one-point set.*
2. (Matching Face Property) *Let  $z_1, z_2 \in X_n$  with  $n > 0$ . Suppose that*
  - i.  $d_i z_1 = d_i z_2, 0 \leq i < n$ , and
  - ii.  $\{z_1, z_2\} \not\subseteq \mathcal{H}_n$ .

*Then  $z_1 = z_2$ .*

*Proof.* Suppose that  $(\mathcal{H}, X)$  is complete. We first show that properties (1) and (2) hold.

1. Assume that  $\mathcal{H}_0 \neq \emptyset$ . Suppose that there exists  $z \in \Delta^X(\mathcal{H}_0) \setminus \mathcal{H}_0$ . Choose an element  $x \in \mathcal{H}_0$ . Let  $Z$  be the  $\mathbb{Z}$ -quotient of  $X$  by identifying  $z$  and  $x$ . Let  $q: X \rightarrow Z$  be the quotient map. Then  $q|_{\mathcal{H}}$  is injective, but  $q: X \rightarrow Z$  is not injective, which is a contradiction. Hence  $X_0 = \mathcal{H}_0$ .

If  $\mathcal{H}_0 = \emptyset$ , then the same argument shows that  $X_0$  must be a one-point set.

2. Suppose  $z_1 \neq z_2$ . Then, similar to the proof of (1), we can construct the  $\mathbb{Z}$ -quotient  $Z$  obtained from  $X$  by identifying  $z_1$  with  $z_2$  in dimension  $n$ . Since all faces of  $z_1$  and  $z_2$  match, the equivalence relation  $z_1 \sim z_2$  does not induce nontrivial identifications in  $X_m$  for  $m < n$ . Since  $\{z_1, z_2\} \not\subseteq \mathcal{H}_n$ , the equivalence relation  $z_1 \sim z_2$  does not effect elements in  $\mathcal{H}$ . This proves property (2).

Now let  $\phi: (\mathcal{H}, X) \rightarrow (\mathcal{H}, Y)$  be a morphism of super-hypergraphs with  $\phi|_{\mathcal{H}} = \text{id}_{\mathcal{H}}$  such that  $(\mathcal{H}, X)$  is a regular super-hypergraph satisfying properties (1) and (2).

We show that the  $\phi$ -map  $\phi: X \rightarrow Y$  is injective. By the vertex property,  $\phi: X_0 \rightarrow Y_0$  is injective. Suppose that  $\phi: X \rightarrow Y$  is not injective. Then there exists  $n > 0$  such that  $\phi: X_k \rightarrow Y_k$  is injective for  $k < n$  and  $\phi: X_n \rightarrow Y_n$  is not injective. It follows that there exists  $z_1, z_2 \in X_n$  with  $z_1 \neq z_2$  such that  $\phi(z_1) = \phi(z_2)$ . Then  $\{z_1, z_2\} \notin \mathcal{H}_n$  because  $\phi|_{\mathcal{H}}$  is injective. Since  $\phi$  is a  $\sigma$ -set map, we have

$$\phi(d_i z_1) = d_i \phi(z_1) = d_i \phi(z_2) = \phi(d_i z_2)$$

for  $0 \leq i < n$  and  $\phi: X_{n-1} \rightarrow Y_{n-1}$  is injective. Therefore  $d_i z_1 = d_i z_2$  for  $0 \leq i < n$ . However, by the matching face property,  $z_1 = z_2$ , which contradicts the assumption that  $z_1 \neq z_2$ .  $\square$

For a given regular super-hypergraph  $(\mathcal{H}, X)$ , the above proof gives a way to construct a complete super-hypergraph  $(\mathcal{H}, Y)$  with  $Y$  as a  $\sigma$ -quotient of  $X$ . The following example shows that  $(\mathcal{H}, X)$  may have non-isomorphic complete quotients.

**Example 2.18.** Let

$$\mathcal{H} = \{\{0, 1, 2\}, \{0, 1\}, \{0, 2\}, \{1, 2\}, \{1\}, \{2\}\}$$

be the graded subset of the 2-simplex  $X = \sigma^+[2]$  without the 0 vertex. Let  $Y$  be the  $\sigma$ -quotient of  $X$  by identifying vertices  $\{1\}$  and  $\{0\}$ , and let  $Z$  be the  $\sigma$ -quotient of  $X$  by identifying vertices  $\{2\}$  and  $\{0\}$ . Then

1.  $(\mathcal{H}, X)$ ,  $(\mathcal{H}, Y)$  and  $(\mathcal{H}, Z)$  are regular super-hypergraph.
2.  $(\mathcal{H}, Y)$  and  $(\mathcal{H}, Z)$  are complete, but  $(\mathcal{H}, X)$  is not complete.
3. There are non-injective  $\sigma$ -quotients  $X \twoheadrightarrow Y$  and  $X \twoheadrightarrow Z$  with  $(\mathcal{H}, Y) \not\cong (\mathcal{H}, Z)$ .

**2.3.3. Mayer-Vietoris sequence.**—The Mayer-Vietoris sequence (MV sequence) is one of the fundamental tools in topology for inductively computing homology. In general, the MV sequence fails for embedded homology of super-hypergraphs.

**Example 2.19.** Let  $X = \{f_1, f_2, e_1, e_2, v\}$  be a  $\sigma$ -set with face operations given by

1.  $d_0 f_j = d_2 f_j = e_1, d_1 f_j = e_2$  where  $j = 1, 2$ , and
2.  $d_i e_j = v$  for  $i = 0, 1$  and  $j = 1, 2$ .

Let

$$\mathcal{H} = X \setminus \{e_1, e_2\}$$

be the graded subset of  $X$  missing the edges  $e_1, e_2$ . Let  $A = \{f_1, e_1, e_2, v\}$  and  $B = \{f_2, e_1, e_2, v\}$  with face operations induced from  $X$ . Then

$$(\mathcal{H}, X) = (\mathcal{H} \cap A, A) \cup (\mathcal{H} \cap B, B)$$

with  $A \cap B = \{e_1, e_2, v\}$  and  $\mathcal{H} \cap A \cap B = \{v\}$ . Then there is no exact sequence

$$\begin{aligned} \dots &\longrightarrow H_2^{\text{emb},A}(\mathcal{H} \cap A) \oplus H_2^{\text{emb},B}(\mathcal{H} \cap B) \longrightarrow H_2^{\text{emb},X}(\mathcal{H}) \\ &\longrightarrow H_1^{\text{emb},A \cap B}(\mathcal{H} \cap A \cap B) \longrightarrow \dots \end{aligned}$$

Since  $\mathcal{H}_1 = \emptyset$  and  $\mathcal{H}_2 = X_2$ ,

$$\begin{aligned} \inf_2^{C_*(A)}(\mathbb{Z}(\mathcal{H} \cap A)) &= \mathbb{Z}(\mathcal{H}_2 \cap A) \cap \text{Ker}(\partial_2^A) = 0, \\ \inf_2^{C_*(B)}(\mathbb{Z}(\mathcal{H} \cap B)) &= \mathbb{Z}(\mathcal{H}_2 \cap B) \cap \text{Ker}(\partial_2^B) = 0, \\ \inf_2^{C_*(X)}(\mathbb{Z}(\mathcal{H})) &= \mathbb{Z}(\mathcal{H}_2) \cap \text{Ker}(\partial_2^X) = \mathbb{Z}\{f_1 - f_2\} = \mathbb{Z}, \\ \inf_1^{C_*(X)}(\mathbb{Z}(\mathcal{H})) &= \mathbb{Z}(\mathcal{H}_1) \cap (\partial_1^X)^{-1}(\mathbb{Z}(\mathcal{H}_0)) = 0, \\ \inf_1^{C_*(A \cap B)}(\mathbb{Z}(\mathcal{H} \cap A \cap B)) &= 0. \end{aligned}$$

Then  $H_2^{\text{emb},A}(\mathcal{H} \cap A) = H_2^{\text{emb},B}(\mathcal{H} \cap B) = H_1^{\text{emb},A \cap B}(\mathcal{H} \cap A \cap B) = 0$  and  $H_2^{\text{emb},X}(\mathcal{H}) = \mathbb{Z}$ .

Hence the above sequence cannot be exact.

An analogue of the classical Mayer-Vietoris sequence for super-hypergraphs is a multi-exact sequence derived from the following theorem.

**Theorem 2.20.** *Let  $(\mathcal{H}, X)$  be a super-hypergraph and let  $A$  and  $B$  be  $\mathcal{H}$ -subsets of  $X$  such that  $A \cup B = X$ . Let  $\mathcal{H}^A = \mathcal{H} \cap A$ ,  $\mathcal{H}^B = \mathcal{H} \cap B$  and  $\mathcal{H}^{A \cap B} = \mathcal{H} \cap A \cap B$ . Let  $G$  be an abelian group, and denote  $\sup_*^{C_*(X;G)}(\mathcal{H})$  and  $\inf_*^{C_*(X;G)}(\mathcal{H})$  by  $\sup_*^X(\mathcal{H})$  and  $\inf_*^X(\mathcal{H})$ , respectively.*

Then there is a commutative diagram

$$\begin{array}{ccc} & & \sup_*^X(\mathcal{H}) \\ & & \parallel \\ \sup_*^A(\mathcal{H}^A) \cap \sup_*^B(\mathcal{H}^B) & \xrightarrow{\quad} & \sup_*^A(\mathcal{H}^A) \oplus \sup_*^B(\mathcal{H}^B) \xrightarrow{j_A - j_B} \sup_*^A(\mathcal{H}^A) + \sup_*^B(\mathcal{H}^B) \\ \uparrow & & \uparrow \simeq \\ \inf_*^A(\mathcal{H}^A) \cap \inf_*^B(\mathcal{H}^B) & \xrightarrow{\quad} & \inf_*^A(\mathcal{H}^A) \oplus \inf_*^B(\mathcal{H}^B) \xrightarrow{j_A - j_B} \inf_*^A(\mathcal{H}^A) + \inf_*^B(\mathcal{H}^B) \\ \parallel & & \uparrow \\ \inf_*^{A \cap B}(\mathcal{H}^{A \cap B}) & & \end{array}$$

where the middle two rows are short exact sequences of chain complexes, the maps  $j_A$  and  $j_B$  are canonical inclusions and the vertical arrows are inclusions.

Example 2.19 shows that the left and right vertical arrows in the above diagram are not chain homotopy equivalences. However, the middle vertical arrow is always a chain homotopy equivalence.

*Proof.* We need to show that

$$\sup_*^A(\mathcal{H}^A) + \sup_*^B(\mathcal{H}^B) = \sup_*^X(\mathcal{H}) \text{ and}$$

$$\inf_*^{A \cap B}(\mathcal{H}^A \cap B) = \inf_*^A(\mathcal{H}^A) \cap \inf_*^B(\mathcal{H}^B).$$

First we prove that  $\sup_*^A(\mathcal{H}^A) + \sup_*^B(\mathcal{H}^B) = \sup_*^X(\mathcal{H})$ . Since  $\sup_*^A(\mathcal{H}^A) + \sup_*^B(\mathcal{H}^B)$  is a sub complex of  $C_*(X; G)$  containing  $Z(\mathcal{H}) \otimes G = Z(\mathcal{H}^A) \otimes G + Z(\mathcal{H}^B) \otimes G$ , we have  $\sup_*^X(\mathcal{H}) \subseteq \sup_*^A(\mathcal{H}^A) + \sup_*^B(\mathcal{H}^B)$ . On the other hand,  $\sup_*^A(\mathcal{H}^A), \sup_*^B(\mathcal{H}^B) \subseteq \sup_*^X(\mathcal{H})$ . Thus  $\sup_*^A(\mathcal{H}^A) + \sup_*^B(\mathcal{H}^B) \subseteq \sup_*^X(\mathcal{H})$ , and so  $\sup_*^X(\mathcal{H}) = \sup_*^A(\mathcal{H}^A) + \sup_*^B(\mathcal{H}^B)$

Now we show that  $\inf_*^{A \cap B}(\mathcal{H}^A \cap B) = \inf_*^A(\mathcal{H}^A) \cap \inf_*^B(\mathcal{H}^B)$ . Clearly

$$\inf_*^{A \cap B}(\mathcal{H}^A \cap B) \subseteq \inf_*^A(\mathcal{H}^A) \cap \inf_*^B(\mathcal{H}^B)$$

Conversely, note that

$$\inf_*^A(\mathcal{H}^A) \cap \inf_*^B(\mathcal{H}^B) \subseteq (Z(\mathcal{H} \cap A) \otimes G) \cap (Z(\mathcal{H} \cap B) \otimes G) = Z(\mathcal{H} \cap A \cap B) \otimes G$$

where  $(Z(\mathcal{H} \cap A) \otimes G) \cap (Z(\mathcal{H} \cap B) \otimes G) = Z(\mathcal{H} \cap A \cap B) \otimes G$  because  $C_*(X; G)$  is the direct sum of the copies of  $G$  with its coordinates labeled by the graded set  $X$ . Thus  $\inf_*^A(\mathcal{H}^A) \cap \inf_*^B(\mathcal{H}^B)$  is a subcomplex of  $C_*(X; G)$  contained in  $Z(\mathcal{H} \cap A \cap B) \otimes G$ , and so it is contained in  $\inf_*^{A \cap B}(\mathcal{H}^A \cap B)$ .  $\square$

**Corollary 2.21.** *Using the notation as in the theorem above, the inclusion*

$$\inf_*^A(\mathcal{H}^A) \cap \inf_*^B(\mathcal{H}^B) \hookrightarrow \sup_*^A(\mathcal{H}^A) \cap \sup_*^B(\mathcal{H}^B)$$

*is a chain homotopy equivalence if and only if so is the inclusion*

$$\inf_*^A(\mathcal{H}^A) + \inf_*^B(\mathcal{H}^B) \hookrightarrow \sup_*^A(\mathcal{H}^A) + \sup_*^B(\mathcal{H}^B).$$

*Proof.* The statement follows by applying the Five Lemma to the long exact sequence obtained from Theorem 2.20.  $\square$

**2.3.4. Gap complexes.**—Let  $(\mathcal{H}, X)$  be a super-hypergraph and define

$$\delta^X(\mathcal{H}) = \bigcup \{Y \subseteq \mathcal{H} \mid Y \text{ is a } \Delta - \text{subset of } X\} \tag{12}$$



to be the largest  $\delta$ -subset of  $X$  contained in  $\mathcal{H}$ . Then  $\delta^X(\mathcal{H})$  consists of the elements in  $\mathcal{H}$  whose all iterated faces lie in  $\mathcal{H}$ . The gap between  $\delta^X(\mathcal{H}) \subseteq \Delta^X(\mathcal{H})$  measures how far  $\mathcal{H}$  is from being a  $\delta$ -set. For a finite hypergraph  $\mathcal{H}$ , the differences can be expressed as

$$\#(\Delta^X(\mathcal{H}) \setminus \delta^X(\mathcal{H})) = \#(\Delta^X(\mathcal{H})) - \#(\delta^X(\mathcal{H})).$$

Topological invariants of the geometric gap complex

$$|\Delta^X(\mathcal{H})|/|\delta^X(\mathcal{H})| \tag{13}$$

such as homology groups and homotopy groups, provide different means of measuring how far  $\mathcal{H}$  is from being a  $\delta$ -set.

Algebraically, at the chain level

$$\begin{aligned} C_*(\delta^X(\mathcal{H}); G) &\hookrightarrow \inf_*^{C_*(\Delta^X; G)}(G(\mathcal{H})) \subseteq G(\mathcal{H}) \subseteq \sup_*^{C_*(\Delta^X; G)}(G(\mathcal{H})) \\ &\hookrightarrow C_*(\Delta^X(\mathcal{H}); G) \end{aligned} \tag{14}$$

where  $G(\mathcal{H}) = \mathbb{Z}(\mathcal{H}) \otimes G$ . An important consequence of Proposition 2.3 is that the inclusion

$$\inf_*^{C_*(\Delta^X; G)}(G(\mathcal{H})) \rightarrow \sup_*^{C_*(\Delta^X; G)}(G(\mathcal{H}))$$

is a chain homotopy equivalence, which implies that the algebraic gap complex (10) is acyclic. However, the geometric gap complex (13) is not contractible in general. For example, let  $X$  be any  $\delta$ -set and let  $\mathcal{H}$  be the graded subset of  $X$  by removing the vertex set  $X_0$ . Then  $\delta^X(\mathcal{H}) = \emptyset$  and so its geometric gap complex is  $|X|^+$ , the space  $|X|$  disjoint union with a one-point set. The homology of the chain complexes

$$\begin{aligned} &\inf_*^{C_*(\Delta^X; G)}(G(\mathcal{H}))/C_*(\delta^X(\mathcal{H}); G), \\ &C_*(\Delta^X(\mathcal{H}); G)/\sup_*^{C_*(\Delta^X; G)}(G(\mathcal{H})) \end{aligned}$$

could give extra information in addition to the topology of the geometric gap complex (13).

The proof of the following proposition is immediate.

**Proposition 2.22.** *Let  $(\mathcal{H}, X)$  be a super-hypergraph. Then*

$$\begin{aligned} \inf_*^{C_*(\Delta^X(\mathcal{H}); G)}(G(\mathcal{H}))/C_*(\delta^X(\mathcal{H}); G) &= \inf_*^{C_*(\Delta^X(\mathcal{H}); G)/C_*(\delta^X(\mathcal{H}); G)}(G(\mathcal{H}))/C_*(\delta^X(\mathcal{H}); G), \\ \sup_*^{C_*(\Delta^X(\mathcal{H}); G)}(G(\mathcal{H}))/C_*(\delta^X(\mathcal{H}); G) &= \sup_*^{C_*(\Delta^X(\mathcal{H}); G)/C_*(\delta^X(\mathcal{H}); G)}(G(\mathcal{H}))/C_*(\delta^X(\mathcal{H}); G). \end{aligned}$$

□

**2.3.5. Computations.**—With the potential applications in mind, it is important to consider the computability of these topological constructions. There have been various algorithms developed for computing simplicial homology that have led to the applications of topology in data analytics. The computations of embedded homology are quite similar to those of simplicial homology. Below we detail a procedure for computing embedded homology.

For a super-hypergraph  $(\mathcal{H}, X)$ , let us consider the computations for  $H_n^{\text{emb}, X}(\mathcal{H}, \mathbb{F})$  using the chain complex  $\text{inf}_*^{C_*(X; \mathbb{F})}(\mathcal{H}; \mathbb{F})$  with coefficients in a field  $\mathbb{F}$ .

By definition,

$$H_n^{\text{emb}, X}(\mathcal{H}, \mathbb{F}) = Z_n\left(\text{inf}_*^{C_*(X; \mathbb{F})}(\mathcal{H}; \mathbb{F})\right) / B_n\left(\text{inf}_*^{C_*(X; \mathbb{F})}(\mathcal{H}; \mathbb{F})\right)$$

where

$$\begin{aligned} Z_n\left(\text{inf}_*^{C_*(X; \mathbb{F})}(\mathcal{H}; \mathbb{F})\right) &= \mathbb{F}(\mathcal{H}_n) \cap \text{Ker}(\partial_n: C_n(X; \mathbb{F}) \rightarrow C_{n-1}(X; \mathbb{F})), \\ B_n\left(\text{inf}_*^{C_*(X; \mathbb{F})}(\mathcal{H}; \mathbb{F})\right) &= \mathbb{F}(\mathcal{H}_n) \cap \partial_{n+1}(C_{n+1}(X; \mathbb{F})). \end{aligned}$$

The *Betti number*  $b_n(\mathcal{H}, X)$  is defined as

$$\begin{aligned} b_n(\mathcal{H}, X) &= \dim H_n^{\text{emb}, X}(\mathcal{H}; \mathbb{F}) \\ &= \dim\left(Z_n\left(\text{inf}_*^{C_*(X; \mathbb{F})}(\mathcal{H}; \mathbb{F})\right) - \dim\left(B_n\left(\text{inf}_*^{C_*(X; \mathbb{F})}(\mathcal{H}; \mathbb{F})\right)\right)\right). \end{aligned} \tag{15}$$

To compute  $Z_n\left(\text{inf}_*^{C_*(X; \mathbb{F})}(\mathcal{H}; \mathbb{F})\right)$ , we can consider the restriction of the linear transformation  $\partial_n: C_n(X; \mathbb{F}) \rightarrow C_{n-1}(X; \mathbb{F})$  to  $\mathbb{F}(\mathcal{H}_n)$ . Namely, consider  $\mathbb{F}(\mathcal{H}_n)$  as a vector spaces over  $\mathbb{F}$  with a basis given by the elements in  $\mathcal{H}_n$ . For each element  $x$  in  $\mathcal{H}_n$ , express  $\partial_n(x)$  as an element in  $C_{n-1}(X; \mathbb{F}) = \mathbb{F}(X_{n-1})$ . This defines a linear transformation

$$\partial_n \upharpoonright : \mathbb{F}(\mathcal{H}_n) \rightarrow \mathbb{F}(X_{n-1})$$

whose kernel is  $Z_n\left(\text{inf}_*^{C_*(X; \mathbb{F})}(\mathcal{H}; \mathbb{F})\right)$ .

For computing  $B_n\left(\text{inf}_*^{C_*(X; \mathbb{F})}(\mathcal{H}; \mathbb{F})\right)$ , we can first consider  $\partial_{n+1}(C_{n+1}(X; \mathbb{F}))$  as a subspace of the vector space  $C_n(X; \mathbb{F}) = \mathbb{F}(X_n)$ , which is spanned by linear combinations of  $\partial_{n+1}(\sigma)$  for  $\sigma \in X_{n+1}$ . Then consider the decomposition

$$\mathbb{F}(X_n) = \mathbb{F}(\mathcal{X}_n) \oplus \mathbb{F}(X_n \setminus \mathcal{X}_n).$$

Let

$$p: \mathbb{F}(X_n) \rightarrow \mathbb{F}(X_n \setminus \mathcal{X}_n)$$

be the projection. Then  $B_n(\inf_*^{C_*(X; \mathbb{F})}(\mathcal{X}; \mathbb{F}))$  is the kernel of the restriction

$$p| : \partial_{n+1}(C_{n+1}(X; \mathbb{F})) \rightarrow \mathbb{F}(X_n \setminus \mathcal{X}_n).$$

If the data  $(\mathcal{X}, X)$  is large, the complexity of direct computation of  $H_*^{\text{emb}, X}(\mathcal{X}, F)$  increases. The existing computational methods for chain complexes aim at reducing this complexity.

### 3. Super-persistent homology.

The general idea of super-persistent homology is to use the geometry of  $n$ -sets and super-hypergraphs as tools to investigate collections of subgraphs in a given graph and to get topological features from for example graphic data and various networks. We will see that a  $n$ -set model performs much better than models using simplicial complexes, particularly for exploring topological features arising from the structures related to clustering.

#### 3.1. General theory.

Let  $G$  be a directed/undirected (multi-)graph. Let  $\mathcal{F}\mathcal{P}(G)$  denote the set of all finite subgraphs of  $G$ .

**Definition 3.1.** A  $n$ -set  $X$  is said to be *dominated by*  $G$  if there exists an injective map  $\phi: X \rightarrow \mathcal{F}\mathcal{P}(G)$  such that  $\phi(d_i(\sigma))$  is a subgraph of  $\phi(\sigma)$  for any  $0 \leq i \leq n$  and any element  $\sigma \in X_n$ .

A super-hypergraph  $(\mathcal{X}, X)$  is said to be *dominated by*  $G$  if its parental  $n$ -set  $X$  is dominated by  $G$ .

For a  $n$ -set  $X$  dominated by  $G$ , we identify the elements  $\sigma$  in  $X$  with its image  $\phi(\sigma)$ , a finite subgraph of  $G$ , and so we consider  $X$  as a collection of finite subgraphs of  $G$ . A super-hypergraph  $(\mathcal{X}, X)$  dominated by  $G$  can be described as a multi-layered collection of families of finite subgraphs  $X = \{X_0, X_1, \dots\}$ , where each  $X_i$  is a family of finite subgraphs, with face operations  $d_i: X_n \rightarrow X_{n-1}$ ,  $0 \leq i \leq n$ , satisfying the  $n$ -identity, and a marked graded subset  $\mathcal{X}$  of  $X$ .

To introduce persistence we need a scoring scheme on  $G$ . For graphs  $P$  and  $Q$ , denote by  $P \preceq Q$  if  $P$  is a subgraph of  $Q$ .

**Definition 3.2.** A *scoring scheme* on a directed/undirected (multi-)graph  $G$  is a function

$$\mathfrak{M}: \mathcal{F}\mathcal{P}(G) \rightarrow \mathbb{R}$$

from the set of finite subgraphs of  $G$  to the real numbers.

A scoring scheme  $\mathfrak{M}$  on  $G$  is called *regular* if for every  $P, Q \in \mathcal{F}\mathcal{P}(G)$  such that  $P \leq Q$ ,

$$\mathfrak{M}(P) \leq \mathfrak{M}(Q).$$

**Definition 3.3.** A *persistent  $\mathbb{R}$ -filtration* of a  $\mathbb{R}$ -set  $X$  over  $\mathbb{R}$  is a family of  $\mathbb{R}$ -subsets  $X(t)$  of  $X$ , indexed by  $t \in \mathbb{R}$ , such that

1.  $X(s)$  is a  $\mathbb{R}$ -subset of  $X(t)$  for  $s \leq t$ , and
2.  $X = \bigcup_{t \in \mathbb{R}} X(t)$ .

A *persistent super-hypergraph filtration* of a super-hypergraph  $(\mathcal{H}, X)$  over  $\mathbb{R}$  is a family of super-hypergraphs  $(\mathcal{H}(t), X(t))$ , indexed by  $t \in \mathbb{R}$ , such that

1. The indexed family  $X(t)$ ,  $t \in \mathbb{R}$ , is a persistent  $\mathbb{R}$ -filtration of  $X$  over  $\mathbb{R}$ ,
2.  $\mathcal{H}(t) = \mathcal{H} \cap X(t)$ .

**Proposition 3.4.** Let  $G$  be a directed/undirected (multi-)graph with a regular scoring scheme  $\mathfrak{M}: \mathcal{F}\mathcal{P}(G) \rightarrow \mathbb{R}$ . Let  $(\mathcal{H}, X)$  be a super-hypergraph dominated by  $G$ . Then

$$X(t) = \mathfrak{M}^{-1}((-\infty, t]) \cap X, \quad t \in \mathbb{R}$$

is a persistent  $\mathbb{R}$ -filtration of  $X$ , and the pair

$$(\mathcal{H}(t), X(t)) = (\mathfrak{M}^{-1}((-\infty, t]) \cap \mathcal{H}, \mathfrak{M}^{-1}((-\infty, t]) \cap X), \quad a \in \mathbb{R}$$

is a persistent super-hypergraph filtration of  $(\mathcal{H}, X)$ .

*Proof.* The proof follows from the definitions.  $\square$

Let  $\mathbb{F}$  be a fixed choice of a ground field. A (graded/ungraded) *persistence module* over  $\mathbb{F}$ , denoted by  $\mathbb{V}$ , is defined to be an indexed family of (graded/ungraded)  $\mathbb{F}$  vector spaces

$$\mathbb{V} = (V(t) \mid t \in \mathbb{R})$$

and a bi-indexed family of (graded/ungraded) linear maps

$$(v_s^t: V(s) \rightarrow V(t) \mid s \leq t)$$

which satisfy the composition law

$$v_s^t \circ v_r^s = v_r^t$$

whenever  $r \leq s \leq t$ , where  $v_t^t$  is the identity map on  $V(t)$ . For graded vector spaces  $V$  and  $W$ , a *graded linear map* of degree  $q$  is a collection of linear maps  $\phi = (\phi_n)_{n \in \mathbb{Z}}$  with  $\phi_n: V_n \rightarrow W_{n+q}$ . A *persistence morphism*  $\Phi$  of degree  $q$  between two graded persistence modules  $\mathbb{V}$  and  $\mathbb{W}$  is a collection of graded linear maps of degree  $q$ ,  $(\phi^t: V(t) \rightarrow W(t) \mid t \in \mathbb{R})$ , such that the diagram

$$\begin{array}{ccc} V(s) & \xrightarrow{v_s^t} & V(t) \\ \downarrow \phi^s & & \downarrow \phi^t \\ W(s) & \xrightarrow{w_s^t} & W(t) \end{array}$$

commutes for  $s \leq t$ . If  $q = 0$ , then  $\Phi$  is called a *persistence morphism*. A persistence morphism between ungraded persistence modules is defined in the same way.

**Definition 3.5.** Let  $(\mathcal{H}, X)$  be a super-hypergraph. Let  $A$  be an abelian group. The *relative embedded homology*  $H_*^{\text{emb}, X}(X, \mathcal{H}; A)$  with coefficients in  $A$  of  $(\mathcal{H}, X)$  is defined by

$$\begin{aligned} H_*^{\text{emb}, X}(X, \mathcal{H}; A) &= H_*\left(C_*(X; A) / \text{inf}_*^{C_*(X; A)}(\mathbb{Z}(\mathcal{H}) \otimes A)\right) \\ &\cong H_*\left(C_*(X; A) / \text{sup}_*^{C_*(X; A)}(\mathbb{Z}(\mathcal{H}) \otimes A)\right). \end{aligned}$$

Now we assume that all homology is taken with coefficients in a ground field  $\mathbb{F}$ . Therefore, we can simplify our notation of homology groups  $H_*(-; \mathbb{F})$  to  $H_*(-)$ . Let  $(\mathcal{H}, X)$  be a super-hypergraph dominated by a directed/undirected (multi-)graph  $G$  with a scoring scheme  $\mathfrak{M}$ . Let  $X(t) = \mathfrak{M}^{-1}((-\infty, t]) \cap X$  and  $\mathcal{H}(t) = \mathfrak{M}^{-1}((-\infty, t]) \cap \mathcal{H}$ . Then

$$\begin{aligned} \mathbb{H}_*(X) &= (H_*(X(t)) \mid t \in \mathbb{R}) \\ \mathbb{H}_*^{\text{emb}, X}(\mathcal{H}) &= (H_*^{\text{emb}, X}(\mathcal{H}(t)) \mid t \in \mathbb{R}) \quad \text{and} \\ \mathbb{H}_*^{\text{emb}, X}(X, \mathcal{H}) &= (H_*^{\text{emb}, X}(X, \mathcal{H}(t)) \mid t \in \mathbb{R}) \end{aligned} \tag{16}$$

are graded persistence modules. The short exact sequence of chain complexes

$$\inf_*^{C_*(X(t))}(\mathbb{F}(\mathcal{H}(t))) \xrightarrow{j^t} C_*(X(t)) \xrightarrow{p^t} C_*(X(t))/\inf_*^{C_*(X(t))}(\mathbb{F}(\mathcal{H}(t)))$$

induces a long exact sequence on homology, which can be written as an exact triangle of graded persistence modules

$$\begin{array}{ccc} \mathbb{H}_*^{\text{emb},X}(\mathcal{H}) & \xrightarrow{\mathbb{J}} & \mathbb{H}_*(X) \\ & \searrow \partial & \swarrow \mathbb{P} \\ & \mathbb{H}_*^{\text{emb},X}(X, \mathcal{H}) & \end{array}$$

(17)

where  $\mathbb{J}$  and  $\mathbb{P}$  are persistence morphisms and the boundary homomorphism  $\partial$  is persistence morphism of degree  $-1$ .

**Definition 3.6.** Let  $(\mathcal{H}, X)$  be a super-hypergraph dominated by a directed/undirected (multi-)graph  $G$  with a scoring scheme. Then the three graded persistence modules listed in (16) together with the exact triangle (17) give a *super-persistent homology* of  $(\mathcal{H}, X)$  with coefficients in a field  $\mathbb{F}$ .

Equally, we could call  $\mathbb{H}_*^{\text{emb},X}(\mathcal{H})$  super-persistent homology. However, the definition given above can carry more information than the embedded homology of the super-hypergraph, as we will now illustrate.

Let  $J \subseteq \mathbb{R}$  be an interval. The *interval persistence module*  $\mathbb{F}^J = (J(t) \mid t \in \mathbb{R})$  is defined by

$$J(t) = \begin{cases} \mathbb{F} & \text{if } t \in J, \\ 0 & \text{otherwise} \end{cases}$$

with double indexed linear maps

$$j_s^t = \begin{cases} \text{id} & \text{if } s, t \in J, \\ 0 & \text{otherwise.} \end{cases}$$

For an ungraded persistence module  $\mathbb{V}$ , the  $q$ -th *suspension*  $\Sigma^q \mathbb{V}$  is a graded persistence module with

$$\Sigma^q \mathbb{V}_n = \begin{cases} \mathbb{V} & \text{if } n = q \\ 0 & \text{otherwise.} \end{cases}$$

A *graded interval persistence module* is a  $q$ -th suspension of the ungraded interval persistence module for some  $q$ .

Recall that a  $\mathbb{N}$ -set  $X$  is called of *finite type* if  $X_n$  is finite for each  $n$ .

**Theorem 3.7** (Structure Theorem). *Let  $(\mathcal{H}, X)$  be a super-hypergraph dominated by a directed/undirected (multi-)graph  $G$  with a scoring scheme. Suppose that  $X$  is of finite type. Then the graded persistence modules  $\mathbb{H}_*(X)$ ,  $\mathbb{H}_*^{\text{emb}, X}(\mathcal{H})$  and  $\mathbb{H}_*^{\text{emb}, X}(X, \mathcal{H})$  admit direct sum decompositions in terms of graded interval persistence modules and these decompositions are unique up to the order of factors in the category of graded persistence modules.*

*Proof.* For any graded persistence module  $\mathbb{V}$ , there is a canonical decomposition

$$\mathbb{V} \cong \bigoplus_{n \in \mathbb{Z}} \Sigma^n \mathbb{V}_n$$

in the category of graded persistence modules, where  $\mathbb{V}_n$  is considered as an ungraded persistence module. It suffices to show that  $\mathbb{H}_n(X)$ ,  $\mathbb{H}_n^{\text{emb}, X}(\mathcal{H})$  and  $\mathbb{H}_n^{\text{emb}, X}(X, \mathcal{H})$  admit unique factorization as ungraded persistence modules for each  $n$ .

Since  $X$  is of finite type, the chain complex  $C_*(X)$  is of finite type and so is any subcomplex or quotient complex. The assertion follows from the structure theorem on persistence modules [46, Theorem 1.1] derived from the classical Gabriel Theorem in representation theory [65].  $\square$

We now briefly summarise persistence diagrams/persistent barcodes; for which this structure theorem is prominent in the calculations, for more details see [39].

Let  $\mathbb{V}$  be an ungraded persistence module with a unique decomposition (up to the order of factors)

$$\mathbb{V} = \bigoplus_{\alpha \in I} \mathbb{F}^{J_\alpha}$$

in terms of interval persistence modules, where  $I$  is the index set. Then the multiset given by  $\inf(J_\alpha), \sup(J_\alpha) \subset \mathbb{R}^2$ ,  $\alpha \in I$ , is the *persistence diagram* (or *barcode*) of  $\mathbb{V}$ , denoted by  $\text{dgm}(\mathbb{V})$ . In our case, under the hypothesis that  $X$  is of finite type, we have three persistence diagrams.

**Corollary 3.8.** *Let  $(\mathcal{H}, X)$  be a super-hypergraph dominated by a directed/undirected (multi-)graph  $G$  with a scoring scheme and let  $X$  be of finite type. Then there are three multi-layer persistence diagrams  $\text{dgm}(\mathbb{H}_n(X))$ ,  $\text{dgm}(\mathbb{H}_n^{\text{emb}, X}(\mathcal{H}))$  and  $\text{dgm}(\mathbb{H}_n^{\text{emb}, X}(X, \mathcal{H}))$  for  $n \geq 0$ .  $\square$*

The exact triangle (17) yields matrix data on the correlations of the multi-layer persistence diagrams as follows. Let

$$\Phi: \mathbb{V} \rightarrow \mathbb{W}$$

be a persistence morphism between ungraded persistence modules. Suppose that both  $\mathbb{V}$  and  $\mathbb{W}$  satisfy the unique factorization property with respect to decompositions in terms of interval persistence modules and let

$$\mathbb{V} = \bigoplus_{\alpha \in I_{\mathbb{V}}} \mathbb{F}^J \alpha \text{ and } \mathbb{W} = \bigoplus_{\beta \in I_{\mathbb{W}}} \mathbb{F}^J \beta.$$

Let  $\Phi_{\alpha, \beta}$  be the composite

$$\Phi_{\alpha, \beta}: \mathbb{F}^J \alpha \xrightarrow{\Phi} \mathbb{F}^J \alpha \xrightarrow{\text{proj}} \mathbb{W} \xrightarrow{\text{proj}} \mathbb{F}^J \beta.$$

According to [26, Proposition 16, Lemma 22], the set of persistence morphisms between two interval persistence modules is either a 1-dimensional vector space or 0. Thus  $\Phi_{\alpha, \beta}: \mathbb{F}^J \alpha \rightarrow \mathbb{F}^J \beta$  is either a generator for  $\text{Hom}(\mathbb{F}^J \alpha, \mathbb{F}^J \beta)$ , or zero. Let the index sets  $I_{\mathbb{V}}$  and  $I_{\mathbb{W}}$  be totally ordered and define the *correlation matrix*

$$M(\Phi) = (m_{\alpha, \beta})_{\alpha \in I_{\mathbb{V}}, \beta \in I_{\mathbb{W}}}$$

of  $\Phi$  by setting

$$m_{\alpha, \beta} = \begin{cases} 1 & \text{if } \Phi_{\alpha, \beta} \neq 0 \\ 0 & \text{otherwise.} \end{cases}$$

The correlation matrix is an analogue of adjacency matrix of graphs, which gives correlations between  $\text{dgm}(\mathbb{V})$  and  $\text{dgm}(\mathbb{W})$  by adding directed edges. In summary, we have the following information data from super-persistent homology.

**Proposition 3.9.** *Let  $(\mathcal{H}, X)$  be a super-hypergraph dominated by a directed/undirected (multi-)graph  $G$  with a scoring scheme. Suppose that  $X$  is of finite type. Then there are three multi-layer persistence diagrams  $\text{dgm}(\mathbb{H}_*(X))$ ,  $\text{dgm}(\mathbb{H}_*^{\text{emb}, X}(\mathcal{H}))$  and  $\text{dgm}(\mathbb{H}_*^{\text{emb}, X}(X, \mathcal{H}))$  together with three correlation matrices  $M(\mathbb{J})$ ,  $M(\mathbb{P})$  and  $M(\ )$  between them.  $\square$*

An important point is that we allow  $\mathcal{H}$  to be an arbitrary graded subset of  $X$ . If we fix  $X$  to be a  $\mathbb{J}$ -set dominated by a graph  $G$  with a scoring scheme and allow  $\mathcal{H}$  to be random, then  $\text{dgm}(\mathbb{H}_*(X))$  is a deterministic barcode, while  $\text{dgm}(\mathbb{H}_*^{\text{emb}, X}(\mathcal{H}))$  and  $\text{dgm}(\mathbb{H}_*^{\text{emb}, X}(X, \mathcal{H}))$  are random. The correlation matrices may help further analyse the data.

Let  $X$  be a fixed  $\mathbb{J}$ -set dominated by a directed/undirected (multi-)graph. It would be also interesting to consider the set  $\mathbb{P}(X)$  of the isomorphic classes of persistence modules (16) for all graded subsets  $\mathcal{H}$  of  $X$ . The inclusions  $\mathcal{H}' \subseteq \mathcal{H} \subseteq X$  induce a morphism of super-hypergraphs  $(\mathcal{H}', X) \rightarrow (\mathcal{H}, X)$ . By taking super-persistent homology, one would get a quiver structure on the set  $\mathbb{P}(X)$ . Moreover, there is an *interleaving distance* between persistence modules introduced in [38] corresponding to *bottleneck distance* [39], which



gives the structure of a metric space on the quiver  $\mathbb{P}(X)$ . The following example illustrates that  $\mathbb{P}(X)$  may give more robust information.

**Example 3.10.** Let  $X$  be a  $\mathbb{P}$ -set. Let  $\mathcal{H}$  be a graded subset of  $X$  consisting of *non-face* elements  $\sigma \in X$ . Here a non-face element  $\sigma$  means that there does not exist any element  $\tau \in X$  such that  $d_i \tau = \sigma$  for some  $i$ . In other words,  $\mathcal{H}$  is given by removing all face elements in  $X$ . It is straightforward to see that

$$H_*^{\text{emb}, X}(\mathcal{H}) = \mathbb{F}(\mathcal{H}) \cap \mathbb{Z}(C_*(X)).$$

Therefore, the embedded homology can detect the cycles contributed from non-face elements. From the exact triangle (17), a part of the boundaries in  $C_*(X)$  can also be detected. These detected elements would contribute to bar-codes through persistence. Hence, in addition to the persistence on the homology  $H_*(X)$ ,  $\mathbb{P}(X)$  can decode more topological features.

### 3.2. The ordinary persistent homology.

The classical persistent homology refers to the persistent homology of point cloud data using the Vietoris-Rips complex, the Čech complex or the witness complex. Persistent homology has been used as an important topological tool in data science. In this subsection, we rewrite on the classical persistent homology from the viewpoint of graphs with scoring schemes, and then give a natural generalization to persistent homology for graphs with reference maps to metric spaces.

**3.2.1. The classical persistent homology.**—A *point cloud* dataset is a finite set  $\mathcal{L}$  with a reference map that embeds  $\mathcal{L}$  into a finite dimensional Euclidean space  $\mathbb{R}^m$ , thus we can consider  $\mathcal{L}$  as a finite subset in  $\mathbb{R}^m$ . We now use a graph with a scoring scheme to describe the (persistent) Vietoris-Rips complex, Čech complex and witness complex in a unified way. The graph  $G(\mathcal{L})$  is the complete graph having  $\mathcal{L}$  as its vertex set. Intuitively, we assign one and only one edge to any two distinct points in  $\mathcal{L}$ . The main point is to show that different scoring schemes on  $G(\mathcal{L})$  can obtain different persistent complexes that are currently widely used in TDA [36]. Below we give the scoring schemes for the Vietoris-Rips complex, the Čech complex, the strong witness complex and the weak witness complex. Let  $\Lambda = \{l_0, l_1, \dots, l_n\} \subseteq \mathcal{L}$  be a subset of  $\mathcal{L}$ .

For  $x \in \mathbb{R}^m$ , denote by  $B(x, r)$  the closed ball of radius  $r$  centered at  $x$ . Define the following scoring schemes on  $\Lambda$ :

1. The *Vietoris-Rips scoring* is given by

$$\mathfrak{M}^{VR}(\Lambda) = \frac{1}{2} \sup \{ d(l_i, l_j) \mid l_i, l_j \in \Lambda \subseteq \mathbb{R}^m \}. \tag{18}$$

The balls  $B(l_0, r), \dots, B(l_n, r)$  pairwise intersect if and only if the score  $\mathfrak{M}^{VR}(\Lambda) \leq r$ .

2. The *ech scoring* is defined by

$$\mathfrak{M}^{\check{C}}(\Lambda) = \inf_{x \in \mathbb{R}^m} \max\{d(x, l) \mid l \in \Lambda \subseteq \mathbb{R}^m\}. \tag{19}$$

Note that

$$\bigcap_{i=0}^n B(l_i, r) \neq \emptyset$$

if and only if the score  $\mathfrak{m}^{\check{C}}(\Lambda) < r$ .

3. The *Strong witness scoring* is defined by

$$\mathfrak{M}^{W^s}(\Lambda) = \inf_{x \in \mathbb{R}^m} \{\sup_{y \in \Lambda} d(x, y) - \inf_{z \in \mathcal{L}} d(x, z)\}. \tag{20}$$

4. The *Vietoris-Rips Strong witness scoring* is given by

$$\mathfrak{M}^{W^s_{VR}}(\Lambda) = \sup_{0 \leq i < j \leq n} \{\inf_{x \in \mathbb{R}^m} \{\max\{d(x, l_i), d(x, l_j)\} - \inf_{z \in \mathcal{L}} d(x, z)\}\}. \tag{21}$$

5. The *Weak witness scoring* is given by

$$\mathfrak{M}^{W^w}(\Lambda) = \inf_{x \in \mathbb{R}^m} \{\sup_{y \in \Lambda} d(x, y) - \inf_{z \in \mathcal{L} \cap \Lambda} d(x, z)\}. \tag{22}$$

6. The *Vietoris-Rips weak witness scoring* is given by

$$\mathfrak{M}^{W^w_{VR}}(\Lambda) = \sup_{0 \leq i < j \leq n} \{\inf_{x \in \mathbb{R}^m} \{\max\{d(x, l_i), d(x, l_j)\} - \inf_{z \in \mathcal{L} \cap \Lambda} d(x, z)\}\}. \tag{23}$$

Let  $X$  be the clique complex of  $G(\mathcal{L})$  considered as a  $\mathbb{R}$ -set. In other words, because  $G$  is complete,  $X$  is the set of all nonempty full subgraphs of  $G(\mathcal{L})$ . Here a *full subgraph*  $H$  of  $G$  is a subgraph  $H$  such that the edge set between any two vertices  $v$  and  $w$  in  $H$  is equal to the edge set between  $v$  and  $w$  in  $G$ . Choose a linear order on the vertex set  $V(G(\mathcal{L}))$ . Then  $(X, X)$  is a super-hypergraph. It is straightforward to check that the persistent super-hypergraph filtrations on  $(X, X)$  induced by the above scoring schemes coincide with the classical persistent filtrations in [36]. Here the witness scoring schemes defined in (3)–(6) are reformulations from [36, Definition 2.7, Definition 2.8].

**3.2.2. Clique persistent homology of graphs with reference maps.**—Next, we will consider a canonical extension of ordinary persistent homology to the case of graphs with reference maps on vertices. Let  $G$  be a finite undirected (multi-)graph with a reference map that embeds the vertex set  $V(G)$  into a finite dimensional Euclidean space  $\mathbb{R}^m$  and let  $H$  be any subgraph. Then, any one of the six scoring schemes in (18)–(23) induces a persistent filtration on the clique complex  $\text{clique}(G)$ . This gives the *clique persistent homology* on  $G$ .

The clique persistent homology on  $G$  could, in general, be quite different from the ordinary persistent homology of the vertex set of  $G$  under the reference map. For instance, if the clique complex  $\text{clique}(G)$  has non-trivial reduced homology, the resulting clique persistent homology converges to  $H_*(\text{clique}(G))$  as  $t \rightarrow \infty$ , but the ordinary persistent homology converges to trivial homology as  $t \rightarrow \infty$ . The ordinary persistent homology of  $V(G)$  under the reference map is obtained by rebuilding a new graph given by the complete graph on  $V(G)$ , that is, all edges in  $G$  are forgotten and the new edges are added in depending on the scoring schemes that are obtained through the reference map. On the other hand, when we consider the clique complex of  $G$  itself, the edges in  $G$  are accounted for.

The following example illustrates how clique persistent homology can describe shapes and therefore could be useful for data analysis on protein structures or image processing on 3D objects with complicated internal structures such as the heart.

**Example 3.11.** Let  $X$  be a polyhedron in  $\mathbb{R}^m$ . Let  $K$  be a simplicial complex that is a triangulation of  $X$  and let  $G$  be the graph given by the 1-skeleton of the barycentric subdivision of  $K$ . Let the reference map on  $G$  be given by the inclusion of  $G$  in  $\mathbb{R}^m$ . Then the geometric realization of the clique complex  $\text{clique}(G)$  is homeomorphic to  $X$ , and the persistent homology of  $\text{clique}(G)$  converges to  $H_*(X)$ . In particular, the number of infinite persistence modules in the  $n^{\text{th}}$  persistent homology of  $\text{clique}(G)$  is equal to the  $n^{\text{th}}$  Betti number of  $X$ . Thus the clique persistent homology detects the topological shape of  $X$ .

More generally we can remove the embedding hypothesis of reference maps. Let  $G$  be a finite undirected (multi-)graph and let

$$f: V(G) \rightarrow \mathbb{R}^m$$

be a function (without assuming injectivity). We can use *pull-back scoring* in the following sense. Let  $H$  be any subgraph of  $G$ . Then the image  $f(V(H))$  is a finite subset located in  $\mathbb{R}^m$ . Let  $\mathfrak{M}$  be a scoring scheme on point cloud data such as one of the six aforementioned scoring schemes. Define

$$\mathfrak{M}^f(H) = \mathfrak{M}(f(V(H))) \quad (\text{Pull-back Scoring}) \quad (24)$$

which induces a persistent filtration on the clique complex  $\text{clique}(G)$  depending on the reference map  $f$ . Different choices of  $f$  would result in different persistence diagrams. For instance, a constant function does induce a trivial persistence on  $H_*(\text{clique}(G))$ . The flexibility of  $f$  could be useful. For example, if  $f$  is randomly given, it induces corresponding random persistence diagram.

The following example illustrates that a pull-back scoring on clique persistent homology may be useful for detecting higher dimensional geometric shapes.

**Example 3.12.** Let  $p: E \rightarrow B$  be a continuous map between polyhedra  $E$  and  $B$ . Assume that  $B$  is a subspace of  $\mathbb{R}^m$ . By triangulating  $B$ , we pull it back along  $p$  do define  $K^E$  and there is a simplicial map  $p': K^E \rightarrow K^B$  such that

1.  $K^E$  and  $K^B$  are simplicial complexes of triangulations of  $E$  and  $B$ , respectively.
2. Let  $G(K^E)$  and  $G(K^B)$  be the graphs given by the 1-skeleton of  $K^E$  and  $K^B$ , respectively. Then  $\text{clique}(G(K^E)) = K^E$  and  $\text{clique}(G(K^B)) = K^B$ .
3.  $p'$  is a simplicial approximation to  $p$ .

Let the reference map  $R$  on  $V(G(K^B))$  be given by the inclusion  $V(G(K^B)) \subseteq B \subseteq \mathbb{R}^m$ . Let us take the pull-back scoring  $\mathfrak{M}$  on  $G(K^E)$ . Then the persistent filtration on  $\text{clique}(G(K^E))$  induced by  $\mathfrak{M}$  is the pull-back of the persistent filtration on  $\text{Clique}(G(K^B))$  induced by  $R$ . In particular, if  $p : E \rightarrow B$  is a fibre bundle or, more generally, a fibration with fibre  $F$ , then we have a persistent Leray-Serre spectral sequence convergent to  $H_*(E)$ . It is well-known in algebraic topology that Leray-Serre spectral sequences are an important tool for computing  $H_*(E)$  starting with  $H_*(B)$  and  $H_*(F)$ .

### 3.3. Partition homology and persistent partition homology.

The methods of data science are typically aimed at finding structures and patterns within large datasets. Being able to glean information about the internal structures of graphical data would be useful in solving the typical problems given to machine learning algorithms. For example, classification problems, prediction and, in particular, partitioning data into clusters [114, Section 1.1.3].

If a collection of subgraphs forms a  $\lambda$ -set structure, then we can calculate homology. A natural question is how to introduce a  $\lambda$ -set structure on a collection of subgraphs in some natural way. More precisely, how to define face operations on subgraphs. We are going to show that any clustering on the vertex set can induce canonical face operations on subgraphs. For a dataset given by a graph, the topological features on collections of subgraphs under the face operations induced by a clustering may help for detecting correlations between the clusters.

Let  $G$  be a directed/undirected (multi-)graph. Assume that there is a disjoint clustering  $\mathbf{p}$  on the vertex set  $V(G)$ . In other words, there is a disjoint union

$$V(G) = \coprod_{i=0}^m V_i(G)$$

under the clustering  $\mathbf{p}$ , where each  $V_i(G)$  is a cluster. Let  $H$  be a subgraph of  $G$ . Then there exists a unique sequence  $(k_0, k_1, \dots, k_n)$  with  $0 \leq k_0 < k_1 < \dots < k_n \leq m$  such that  $V(H) \cap V_i(G) = \emptyset$  for  $i \in \{k_0, k_1, \dots, k_n\}$  and  $V(H) \cap V_i(G) \neq \emptyset$  if  $i \notin \{k_0, k_1, \dots, k_n\}$ .

We call  $H$  a subgraph of  $G$  linked to  $(n + 1)$  clusters. Viewing  $H$  as an abstract  $n$ -simplex, we define the  $j^{\text{th}}$ -face map,  $d_j^{\mathbf{p}}(H)$ , as the full subgraph of  $H$  formed after removing all of those vertices  $v \in V(H) \cap V_{k_j}(G)$  together with the edges incident to (or from) such  $v$ .

The resulting subgraph  $d_j^{\mathbf{p}}(H)$  is linked to  $n$  clusters with  $V(d_j^{\mathbf{p}}(H)) \cap V_{k_j}(G) = \emptyset$ . It is straightforward to show that the  $\partial$ -identity

$$d_i^{\mathbf{p}} d_j^{\mathbf{p}} = d_j^{\mathbf{p}} d_{i+1}^{\mathbf{p}}$$

for  $i < j$  holds. The face operation  $d_j^{\mathbf{p}}$  is induced by the disjoint clustering  $\mathbf{p}$ .

Now let  $\mathcal{H}$  be any collection of subgraphs of  $G$ . Define  $\mathcal{H}_n$  to be the subset of  $\mathcal{H}$  consisting of those subgraphs in  $\mathcal{H}$  linked to  $(n + 1)$  clusters. This gives a graded structure on  $\mathcal{H} = \{\mathcal{H}_n\}_{n \geq 0}$ . Let  $X(\mathcal{H})$  be the collection of subgraphs of  $G$  given by all of the elements in  $\mathcal{H}$  together with all iterated faces. Then  $X(\mathcal{H})$  is a  $\Delta$ -set and  $(\mathcal{H}, X(\mathcal{H}))$  is a super-hypergraph. The resulting homology groups

$$H_*(X), H_*^{\text{emb}, X}(\mathcal{H}), H_*^{\text{emb}, X}(X, \mathcal{H})$$

are called the *partition homology*.

In general,  $X(\mathcal{H})$  may not be a simplicial complex. In simplicial complexes, the assumption that simplices are determined by their vertices is too strict with applications in mind. For exploring topological structures arising from disjoint clusterings, a  $\Delta$ -set is a more suitable notion.

In practice, for studying possible correlations between clusters, one could start with a collection of one or more subgraphs,  $\mathcal{H}$ , linked with some or all clusters, and then produce the  $\Delta$ -set  $X(\mathcal{H})$ . Due to the nature of simplicial homology, higher dimensional homology of  $X$ , and higher dimensional embedded homology of  $(\mathcal{H}, X)$  would give topological features measuring the group correlations between more clusters. In particular, we have set up the  $\Delta$ -set  $X(\mathcal{H})$  such that the homological dimension of a given subgraph  $H$  in  $X$  is  $n - 1$ , where  $n$  is the number of clusters linked with  $H$ .

In theory, we can start with any collection of subgraphs as the initial data  $\mathcal{H}$ . For instance, we can start with  $\mathcal{H}$  given by all or some of the  $k$ -regular subgraphs of  $G$ , Eulerian subgraphs, traceable subgraphs or Hamiltonian subgraphs as initial data for constructing the  $\Delta$ -set  $X(\mathcal{H})$ . This would give different topological approaches to understanding the internal structures of  $\mathcal{H}$  in addition to what we will survey in Appendix A.

If there is a scoring scheme on  $G$ , then we can calculate the super-persistent homology of  $(\mathcal{H}, X)$  called *persistent partition homology*. A scoring scheme on  $G$  can be deterministically or randomly given. If a scoring scheme is randomly given, it may not be regular. This means that in the induced persistent filtration on  $X$ , the graded subset  $X(t)$  may not be a  $\Delta$ -subset of  $X$  for all  $t$ . In this case, the persistence system can be modified replacing chains related to the terms  $X(t)$  by infimum or supremum chains on  $X(t)$ . The resulting persistence modules and persistence diagrams can be modified accordingly.

From the perspective of data processing, a clustering may be compared against certain optimization properties. Currently, we use discrete Morse theory which works well on

simplicial complexes and cell complexes [63], and chain complexes [88] with applications in data analysis [115]. Moreover, the combinatorial Laplacian operator works well on simplicial complexes and chain complexes, where cohomology with coefficients in real numbers can be expressed as the null space of the Laplacian operator on cochains [82].

### 3.4. Other face operations.

We have shown how disjoint clusterings can induce face operations on subgraphs. This construction works well as a theory which unifies many constructions such as the clique complex and the neighborhood complex. However, if we are interested in subgraphs having some special properties, this construction has some disadvantages. For instance, if we are interested in collections of finite connected subgraphs  $H$ , then the subgraph of  $H$  given by removing some of its vertices may not be connected. In the following subsection, we will discuss some alternative ways of getting natural constructions of face operations.

**3.4.1. Link-blowup face operations.**—The idea of link-blowup face operations comes from geometric constructions on tubular neighborhoods of submanifolds or regular neighborhoods of subcomplexes. Let  $H$  be a subgraph of some graph  $G$  and let  $S$  be a subset of the vertex set  $V(H)$ . If we remove  $S$  from  $H$ , then we could add some edges from the working graph  $G$  to make a blowup for the subgraph  $H \setminus S$ . A natural way to add these edges is to consider the neighbors of vertices of  $S$  in  $H$  and to add the edges between these neighbors from the working graph  $G$ .

Let  $G$  be a directed/undirected (multi-)graph and let  $E^G(v, w)$  denote the edge set between  $v$  and  $w$  for vertices  $v, w \in V(G)$ . Let  $S \subseteq V(G)$  be a subset and define the *induced subgraph of  $S$  in  $G$* ,  $G[S]$ , to be the full subgraph of  $G$  having  $S$  as its vertex set. The *closed neighborhood set*  $N[S]$  is defined by

$$N[S] = S \cup \{u \mid u \in V(G) \text{ adjacent to a vertex } v \in S\}$$

namely,  $N[S]$  is the union of the neighborhoods of vertices  $v \in S$ . The *link set* of  $S$  in  $G$  is

$$\text{Lk}(S) = N[S] \setminus S.$$

Let  $H$  be a subgraph of  $G$  and let  $S \subseteq V(H)$  be a subset of the vertex set of  $H$ . The *link-blowup*<sup>2</sup> of  $H$  along  $S$  is defined as

$$H[V(H) \setminus S] \cup G[\text{Lk}(S) \cap V(H)].$$

Note that the graph  $H[V(H) \setminus S] \cup G[\text{Lk}(S) \cap V(H)]$  has the same vertex set of  $H[V(H) \setminus S]$ .

<sup>2</sup>This definition is taken from a geometric setting. We may consider the subgraph  $G[N[S]]$  as a regular neighborhood of  $S$ . Then  $\text{Lk}(S)$  are the vertices located in the “boundary” of the regular neighborhood. Geometrically, we add all of edges in  $G$  joining vertices in  $\text{Lk}(S) \cap V(H)$  to form a blowup on  $H[V(H) \setminus S]$ .

Now suppose that there is a disjoint clustering  $\mathbf{p}$  on the vertex set  $V(G)$  so that there is a disjoint union

$$V(G) = \coprod_{i=0}^m V_i(G)$$

under  $\mathbf{p}$  with each  $V_i(G)$  a cluster. Let  $H$  be a subgraph. Then there exists a unique sequence  $(k_0, k_1, \dots, k_n)$  with  $0 < k_0 < k_1 < \dots < k_n = m$  such that  $V(H) \cap V_i(G) = \emptyset$  for  $i \in \{k_0, k_1, \dots, k_n\}$  and  $V(H) \cap V_i(G) \neq \emptyset$  if  $i \notin \{k_0, k_1, \dots, k_n\}$ . Let  $V_j(H) = V(H) \cap V_{k_j}(G)$ . Define the link-blowup face operation as

$$d_j^{\text{lk}}(H) = H[V(H) \setminus V_j(H)] \cup G[\text{Lk}(V_j(H)) \cap V(H)] \tag{25}$$

that is, the link-blowup of  $H$  along  $V_j(H) = V(H) \cap V_{k_j}(G)$  for  $0 \leq j < n$ .

**Remark 3.13.** For helping to understand the link-blowup face operation, one can give a coloring on the vertices of  $G$  so that the vertices in each cluster has the same color under the disjoint clustering  $\mathbf{p}$ , and perform the link-blowup on the vertices having the same color.

**Proposition 3.14.** Given a disjoint clustering  $\mathbf{p}$ , let  $d_j^{\text{lk}}$  be defined as above. Then

$$d_i^{\text{lk}} d_j^{\text{lk}} = d_j^{\text{lk}} d_{i+1}^{\text{lk}} \tag{26}$$

for  $i < j$ .

*Proof.* Let  $H$  be a subgraph of  $G$ . Let  $V_j$  denote  $V_j(H)$  defined as above and let  $d_j^{\text{lk}}$  denote  $d_j^{\text{lk}}$ , the induced face operations from the disjoint clustering  $\mathbf{p}$ . We will use the fact that

$$V(d_j^{\text{lk}}(H)) = V(d_j(H)) = V(H) \setminus V_j.$$

Now

$$\begin{aligned} d_i^{\text{lk}}(d_j^{\text{lk}}(H)) &= d_j^{\text{lk}}(H)[V(d_j^{\text{lk}}(H)) \setminus V_{i+1}] \cup G[\text{Lk}(V_{i+1}) \cap V(d_j^{\text{lk}}(H))] \\ &= d_j^{\text{lk}}(H)[V(d_j^{\text{lk}}(H)) \setminus V_{i+1}] \cup G[\text{Lk}(V_{i+1}) \cap (V(H) \setminus V_j)] \\ &= d_j^{\text{lk}}(H)[V(d_j^{\text{lk}}(H)) \setminus V_{i+1}] \cup G[\text{Lk}(V_{i+1}) \cap (V(H) \setminus V_j \setminus V_{i+1})] \end{aligned}$$

because  $\text{Lk}(V_{i+1}) \cap V_{i+1} = \emptyset$ .

The term  $d_j^{\text{lk}}(H)[V(d_j^{\text{lk}}(H)) \setminus V_{i+1}]$  is the induced subgraph of  $d_j^{\text{lk}}(H)$  on its vertex subset

$$V(d_j^{\text{lk}}(H)) \setminus V_{i+1} = V(H) \setminus V_j \setminus V_{i+1}.$$

By definition,  $d_j^{\text{lk}}(H) = H[V(H) \setminus V_j(H)] \cup G[\text{Lk}(V_j(H)) \cap V(H)]$ . Restricting to the vertex subset

$$W = V(H) \setminus V_j \setminus V_{i+1}$$

we have

$$d_j^{\text{lk}}(H)[V(d_j^{\text{lk}}(H)) \setminus V_{i+1}] = H[W] \cup G[\text{Lk}(V_j) \cap W]$$

and so

$$d_i^{\text{lk}}(d_j^{\text{lk}}(H)) = H[W] \cup G[\text{Lk}(V_j) \cap W] \cup G[\text{Lk}(V_{i+1}) \cap W].$$

By the same arguments,

$$d_j^{\text{lk}}(d_{i+1}^{\text{lk}}(H)) = H[W] \cup G[\text{Lk}(V_j) \cap W] \cup G[\text{Lk}(V_{i+1}) \cap W].$$

□

**3.4.2. Face operations on subgraphs with marked starting-vertices.**—Let  $H$  be a subgraph of  $G$  and let  $S$  be a subset of the vertex set  $V(H)$ . When we remove  $S$  from  $H$ , we wish to add as few edges as possible to make a blowup for the subgraph  $H \setminus S$  with the aim of preserving particular properties of  $H$ . For this construction, consider an extension of the working graph  $G$  by adding an extra edge between any two distinct vertices  $v, w$  in  $G$  labeled as  $\infty_{vw}$ . This is an analogue to the idea of compactification in geometry. To showcase this we consider a special family of subgraphs.

Let  $G$  be a directed/undirected (multi-)graph.

**Definition 3.15.** A *subgraph with marked starting-vertices* of  $G$  is a pair  $(H, \text{SV}(H))$  satisfying the following conditions:

1.  $H$  is a subgraph of  $G$ ,
2.  $\text{SV}(H)$  is a subset of  $V(H)$ ,
3. Every vertex  $v$  of  $H$  is reachable by a directed/undirected path out from a vertex in  $\text{SV}(H)$ .

In the case of digraphs, one may require further that there are no directed edges in  $H$  incident into any vertex in  $\text{SV}(H)$ . In this case,  $\text{SV}(H)$  acts as a source set for  $H$ . We want to give a description of face operations that is consistent across graphs and digraphs, therefore we do not require such an extra condition. There could be some redundant vertices contained in  $\text{SV}(H)$ . A trivial example is to choose  $\text{SV}(H) = V(H)$ , in which case there are no face operations. If  $\text{SV}(H)$  is a proper subset of  $V(H)$ , there will be nontrivial face operations.



Let  $(H, SV(H))$  be a subgraph with marked starting-vertices of  $G$ . We will now recursively construct a partition on the vertex set  $V(H)$ , called the *neighborhood-extension partition*, in the following way. Let  $V_0(H) = SV(H)$  and suppose that  $V_j(H)$  is constructed with  $j \geq 0$ . Define  $V_{j+1}(H)$  as follows:

1. In the undirected case, let  $V_{j+1}$  be the link set of the subgraph  $H[V_j]$  in the graph  $H$ .
2. In the directed case, let  $V_{j+1}$  be the *out-link set* of the subgraph  $H[V_j]$  in the graph  $H$ . Here, for a (multi-)digraph  $\Gamma$  and a subgraph  $\Gamma'$ , the *closed out-neighborhood set* of  $\Gamma'$  in  $\Gamma$  is the union of  $\Gamma'$  and the out-neighbors of  $\Gamma'$  in  $\Gamma$ , denoted by  $N^{\text{out}}(\Gamma')$ . The *out-link set* of  $\Gamma'$  is defined as

$$Lk^{\text{out}}(\Gamma') = N^{\text{out}}(\Gamma') \setminus V(\Gamma').$$

For a subset  $S$  of  $V(\Gamma)$ , let  $N^{\text{out}}(S) = N^{\text{out}}(\Gamma[S])$  and  $Lk^{\text{out}}(S) = Lk^{\text{out}}(\Gamma[S])$ .

If  $(H, SV(H))$  is a finite subgraph with marked starting-vertices of  $G$ , then above recursive construction will stop after finitely many steps, hence this gives a finite partition on  $V(H)$ .

To define face operations taking marked starting-vertices in to account, we embed  $G$  into a larger graph  $\hat{G}$ , where  $V(\hat{G}) = V(G)$  and  $E(\hat{G})$  is the extension of  $E(G)$  by adding one edge  $\infty_{vw}$  for  $v, w \in V(G)$  in the undirected case, and by adding two directed edges  $\infty_{vw}$  from  $v$  to  $w$  and  $\infty_{wv}$  from  $w$  to  $v$  for vertices  $v, w \in V(G)$  in the directed case.

Now let  $(H, SV(H))$  be a finite subgraph with marked starting-vertices of  $\hat{G}$  with the neighborhood extension partition

$$V(H) = \bigsqcup_{i=0}^n V_i.$$

We assume that  $H \setminus \emptyset$  and so  $V_0 = SV(H) \setminus \emptyset$ . From the recursive definition,  $V_{i+1} \setminus \emptyset$  implies that  $V_i \setminus \emptyset$ . Therefore,  $V_n$  is the last nonempty set in the recursive procedure.

We define the face operation  $d_j^{\text{SV}}$  on  $H$ ,  $0 \leq j \leq n$ , with  $n > 0$  as follows:

1.  $d_0^{\text{SV}}(H) = H[V(H) \setminus V_0]$  with  $SV(d_0^{\text{SV}}(H)) = V_1$ .
2.  $d_n^{\text{SV}}(H) = H[V(H) \setminus V_n]$  with  $SV(d_n^{\text{SV}}(H)) = V_0$ .
3. For  $0 < j < n$ ,

$$d_j^{\text{SV}}(H) = H[V(H) \setminus V_j] \cup \mathcal{E}^H(V_{j-1}, V_{j+1})$$

with  $SV(d_n^{\text{SV}}(H)) = V_0$ , where  $\mathcal{E}^H(V_{j-1}, V_{j+1})$  is a subset of the edge set  $E(\hat{G})$  consisting of  $\infty_{vw}$  for  $v \in V_{j-1}$  and  $w \in V_{j+1}$  satisfying the property that there does not exist an edge from  $v$  to  $w$  in  $H^3$ .

Let  $V_{-1} = V_{n+1} = \emptyset$ . Then we can write in a unified way that

$$d_j^{SV}(H) = H[V(H) \setminus V_j] \cup \mathcal{E}^H(V_{j-1}, V_{j+1}) \text{ with } SV(d_j^{SV}(H)) = V_{\delta_j, 0} \quad (27)$$

for  $0 \leq j \leq n$ , where  $\delta_{a, b}$  is the Kronecker  $\delta$  symbol.

We need to show that  $(d_j^{SV}(H), SV(d_j^{SV}(H)))$  is also a finite subgraph with marked starting-vertices of  $\hat{G}$ . This is straightforward because we add in  $\infty_{vw}$  for possible missing edges in  $H$  from  $V_{j-1}$  to  $V_{j+1}$ . We also need to show that the  $\mathcal{F}$ -identity holds for these face operations. Let  $i \leq j$ . For  $i > j$ , we get

$$d_i^{SV}(d_j^{SV}(H)) = d_j^{SV}(d_{i+1}^{SV}(H)) = H[V(H) \setminus V_j \setminus V_{i+1}] \cup \mathcal{E}^H(V_{j-1}, V_{j+1}) \cup \mathcal{E}^H(V_i, V_{i+2}).$$

For  $i = j$ , we have

$$d_j^{SV}(d_j^{SV}(H)) = d_j^{SV}(d_{j+1}^{SV}(H)) = H[V(H) \setminus V_j \setminus V_{j+1}] \cup \mathcal{E}^H(V_{j-1}, V_{j+2}).$$

This gives the following proposition.

**Proposition 3.16.** *In the set of finite subgraphs with marked starting-vertices of  $\hat{G}$ , the operations  $d_j^{SV}$  are well-defined and satisfy the  $\mathcal{F}$ -identity for face operations.  $\square$*

Now let  $\mathcal{H}$  be a collection of finite subgraphs with marked starting-vertices of  $G$ . Under the extension  $G \leq \hat{G}$ ,  $\mathcal{H}$  is also a collection of finite subgraphs with marked starting-vertices of  $\hat{G}$ . Let  $X$  be the collection of finite subgraphs with marked starting-vertices of  $\hat{G}$  given by all elements in  $\mathcal{H}$  together with all of their iterated faces under face operations  $d_j^{SV}$ . Then  $X$  is  $\mathcal{F}$ -set dominated by  $\hat{G}$ . We call  $(\mathcal{P}(G) \cap X, X)$  the *super-hypergraph generated by  $\mathcal{H}$* .

**3.4.3. Revisiting path complexes.**—We now come back to the work of Yau’s school on path (co-)homology of graphs [76]. Following their terminology, a simple digraph  $G$  is a pair  $(V, E)$ , where  $V$  is any set and  $E \subseteq \{V \times V \setminus \text{diag}\}$ . We will show that  $d_j^{SV}$  can be used to describe the face operations given in [76]. To do this, we consider the “largest quotient simple digraph”  $\tilde{G}$  on  $\hat{G}$ . Let  $V(\tilde{G}) = V(\hat{G}) = V$ , and for any ordered pair  $(v, w) \in V \times V \setminus \text{diag}$ , identify all the directed edges from  $v$  to  $w$  to give one such directed edge. Since  $G$  is a simple graph,  $\tilde{G}$  can be chosen to be a quotient of  $\hat{G}$ . Then  $G$  is a subgraph of the simple digraph  $\tilde{G}$ . Let  $v, w \in V(G) = V(\tilde{G})$  be two distinct vertices. If there exists a directed edge  $e_{vw}$  from  $v$  to  $w$ , then  $\infty_{vw}$  is identified with  $e_{vw}$ . Otherwise  $\infty_{vw}$  is isolated. The graph  $\tilde{G}$  can be considered as the completion of  $G$  in the sense that it is the smallest complete simple digraph containing  $G$ . Here a complete simple digraph means a simple digraph with the property that for any two distinct vertices  $v$  and  $w$ , there are exactly two directed edges with one from  $v$  to  $w$  and another from  $w$  to  $v$ .

<sup>3</sup>In the directed case, we only consider directed edges from a vertex in  $V_{j-1}$  to a vertex in  $V_{j+1}$  in  $H$ . If there are no directed edges in  $H$  from  $v \in V_{j-1}$  to  $w \in V_{j+1}$ , we add  $\infty_{vw}$  to join them with direction from  $v$  to  $w$ .

Replacing  $\hat{G}$  by  $\tilde{G}$ , our face operation  $d_j^{SV}$  on paths in  $\tilde{G}$  coincides with the face operation in [76, Section 4] in the sense that it describes the  $j$ -th term in the boundary operators for chains and cochains on the path complex. Here, to match the definitions, a regular elementary path in [76, Section 4] is a path in  $\tilde{G}$  and an allowed regular elementary path in [76, Section 4] is a path in the subgraph  $G$ . Then one can obtain the same objects by going via the definition of path homology in [76] or the definition of embedded homology of hypergraphs given above. In the undirected case, the operations  $d_j^{SV}$  on paths in  $\tilde{G}$  coincide with the face operations in [76, Section 5].

For directed multi-graphs (quivers), the operations  $d_j^{SV}$  describe the face operations in [77]. Similarly to the undirected case, we need to do a certain identification on  $\hat{G}$ . Following the argument in [77, Section 3], for a complete quiver  $G$ ,  $\infty_{vw}$  is identified with the 1-chain given by the sum of all directed edges from  $v$  to  $w$  in  $G$ . This would define a chain complex on the path complex of a complete quiver. For an arbitrary quiver  $G$ , one can embed  $G$  into its completion  $\bar{G}$ , and take the infimum chain complex (in Proposition 2.2) of the path complex of  $G$  in the chain complex of the path complex of  $\bar{G}$  to define path homology for the quiver  $G$ , see [77] for details.

As notions of paths and walks are commonly used in data analytics, generalising them to higher dimensional combinatorial objects, such as path complexes, could provide new tools for various applications. From a data science point of view, a graph  $G$  is assumed as a working data. Then the path homology gives some topological information on  $G$ . Using a scoring scheme of  $G$ , one could get persistent path homology of  $G$ , this gives a persistence diagram/barcode of  $G$  as a topological feature. However, if the graphic data-set  $G$  is large, the computational complexity may be an issue. From such a perspective, it would be reasonable to consider a selected sub complex, or more generally a selected graded subset of the path complex. The general theory developed in this article gives a framework that makes it possible to explore topological features from subcomplexes or graded subsets of path complexes.

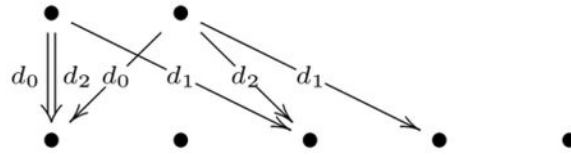
### 3.5. Descriptions of simplicial homology via $\mathbb{F}$ -sets.

In this subsection we only consider mod 2 homology, hence the coefficients are taken in  $\mathbb{F} = \mathbb{Z}/2$ .

**3.5.1.  $\mathbb{F}$ -neural network.**—We can interpret a  $\mathbb{F}$ -set as a quiver or network. A similar object is a feed-forward neural network, as defined in [141, Section 3.6.1].

*Feed-forward neural networks*, are the simplest form of artificial neural networks. The feed forward neural network was the first and arguably simplest type of artificial network devised. In this network, the information moves in only one direction, forward, from the input nodes, through the hidden nodes (if any) and to the output nodes. There are no cycles or loops in the network. In a feed-forward system, processing elements are arranged into distinct layers with each layer receiving input from the previous layer and outputting to the next layer.

Let  $X = \{X_n\}_{n=0}^n$  be a  $n$ -set. Each element in  $X$  is considered as a node (vertex), that is, the node set  $X$  is partitioned by layers labeled by  $X_0, X_1, \dots$ . For each  $x \in X_n$ , assign one and only one arrow (directed edge)  $d_i^n: x \rightarrow d_i(x)$  for  $0 \leq i \leq n$ . This forms a set of arrows whose tails lie in  $X_n$  and whose heads lie in  $X_{n-1}$ . So it forms a quiver. The following picture illustrates the arrows from  $X_2$  to  $X_1$ .



Rephrasing the definition of  $n$ -set into the terminology of network, a  $n$ -neural network is a quiver with distinct layers of nodes labeled by  $X_0, X_1, \dots$  such that for each node  $x \in X_n$  with  $n > 0$ , there are arrows  $d_i^n: x \rightarrow d_i(x)$ ,  $0 \leq i \leq n$ , tailed at  $x$  and headed at some node  $d_i^n(x) \in X_{n-1}$  such that

$$d_i^{n-1}(d_j^n(x)) = d_j^{n-1}(d_{i+1}^n(x)) \tag{28}$$

for  $0 \leq j \leq i \leq n-1$ . In a  $n$ -neural network, the information only flows in one direction, from input nodes that could be located in different layers to the output nodes.

The adjacency relationship from the  $n^{\text{th}}$  layer  $X_n$  to the  $(n-1)^{\text{th}}$  layer  $X_{n-1}$  can be described by  $(n+1)$  matrices as follows. For  $x \in X_n$  and  $y \in X_{n-1}$ , let

$$w_{x,y}^i = \begin{cases} 1 & \text{if } y = d_i(x) \\ 0 & \text{otherwise} \end{cases}$$

for  $0 \leq i \leq n$ . Let

$$W_n(i) = (w_{x,y}^i)_{x \in X_n, y \in X_{n-1}}$$

be a  $|X_n| \times |X_{n-1}|$  matrix, which is the matrix for the face operation  $d_i: X_n \rightarrow X_{n-1}$ . Equation (28) can be rewritten as the following formula

$$W_n(j)W_{n-1}^t(i) = W_n(i+1)W_{n-1}^t(j) \tag{29}$$

for  $0 \leq j \leq i \leq n-1$ , where  $A^t$  is the transpose of a matrix  $A$ .

Let  $\mathcal{H}$  be a graded subset of the  $n$ -set  $X$ . In our theory of super-hypergraphs,  $\mathcal{H}$  carries part of the  $n$ -set structure of  $X$ . More precisely, the face operation  $d_i: X_n \rightarrow X_{n-1}$  induces a partially defined face operation  $d_i: \mathcal{H}_n \rightsquigarrow \mathcal{H}_{n-1}$ . By considering this as a network,  $\mathcal{H}$  is full subnetwork of the  $n$ -neural network induced by  $X$ , where a *full subnetwork* is the induced network of the nodes of  $\mathcal{H}$  in the neural network induced by  $X$ . Therefore, one can

get homology on any full subnetwork of a  $\mathcal{H}$ -neural network using embedded homology of super-hypergraphs.

**Remark 3.17.** The product rule (28) is important to define the boundary operator on the chains. But variations are possible. For example, one could vary the product rule for weighted simplicial complexes [116], or the boundary operators on cochains could be varied to account for twisted de Rham cohomology [56].

**3.5.2. Descriptions of mod 2 homology.**—We proceed by giving the ideas behind the intuition of mod 2 homology of  $\mathcal{H}$ -sets and super-hypergraphs. Let  $(\mathcal{H}, X)$  be a super-hypergraph. Then the  $n$ -chains on  $X$  are linear combinations of the elements in  $X_n$  with coefficients in  $\mathbb{Z}/2$ . So each  $n$ -chain  $\alpha$  corresponds to a subset  $\{x_1, \dots, x_k\} \subseteq X_n$  given by  $\alpha = x_1 + x_2 + \dots + x_k$ . Since the coefficients are in  $\mathbb{Z}/2$ ,

$$\partial(\alpha) = \sum_{i=1}^k \partial(x_i) = \sum_{i=1}^k \sum_{j=0}^n d_j(x_i)$$

which is the *trace* of the multi-subset  $\{d_j(x_i) \mid 0 \leq j \leq n, 1 \leq i \leq k\}$  of  $X_{n-1}$ . Here the multiplicity of  $y = d_j(x_i) \in X_{n-1}$  is the number of pairs  $(j', i')$  such that  $d_{j'}(x_{i'}) = d_j(x_i) = y$ , that is, the in-degree of the node  $y \in X_{n-1}$  in the  $\mathcal{H}$ -neural network. Hence we have the following proposition.

**Proposition 3.18.** *An  $n$ -chain  $\alpha = x_1 + x_2 + \dots + x_k$  is a mod 2 cycle (that is,  $\partial(\alpha) = 0$ ) if and only if any node in the subset*

$$\{d_j(x_i) \mid 0 \leq j \leq n, 1 \leq i \leq k\}$$

*of  $X_{n-1}$  has even in-degree.  $\square$*

This proposition indicates that one can consider the nodes in  $X_{n-1}$  with even in-degrees in search for possible mod 2 cycles in  $n$ -chains.

The following proposition follows from the fact that  $Z_n(\text{inf}_{*}^{C_*(X)}(\mathcal{H})) = \mathbb{Z}/2(\mathcal{H}_n) \cap Z_n(C_*(X))$ .

**Proposition 3.19.** *A mod 2 cycle  $\alpha = x_1 + x_2 + \dots + x_k$  in the chains  $C_n(X)$ , with all  $x_i$  distinct represents a cycle for the mod 2 embedded homology  $H_n^{\text{emb}, X}(\mathcal{H})$  if and only if  $\{x_1, \dots, x_k\} \subseteq \mathcal{H}_n$ .  $\square$*

An  $n$ -chain  $\alpha = x_1 + \dots + x_k$  with all  $x_i$  distinct is a boundary in the chain complex  $C_*(X)$  if and only if the equation

$$\alpha = \partial(\beta)$$

where  $\beta = y_1 + \dots + y_m$  with all  $x_i$  distinct in  $X_{n+1}$  has a solution. If there is a solution  $\alpha = (\beta)$ , then  $\{x_1, \dots, x_k\}$  is the set of nodes in the multi-set  $\{d_s(y_i) \mid 0 \leq s \leq n+1, 1 \leq i \leq m\}$  which have odd in-degrees. This proves the following statement.

**Proposition 3.20.** *An  $n$ -chain  $\alpha = x_1 + \dots + x_k$  with all  $x_i$  distinct is a boundary in the mod 2 chain complex  $C_*(X)$  if and only if there exists a subset  $\{y_1, \dots, y_m\} \subseteq X_{n+1}$  with  $y_1, \dots, y_m$  distinct such that  $\{x_1, \dots, x_k\}$  is the set of nodes in the multi-set  $\{d_s(y_i) \mid 0 \leq s \leq n+1, 1 \leq i \leq m\}$  which have odd in-degrees.  $\square$*

Note that  $B_n(\inf_*^{C_*(X)}(\mathcal{H})) = \mathbb{Z}/2(\mathcal{H}_n) \cap \partial(\mathbb{Z}/2(\mathcal{H}_{n+1}))$ .

**Proposition 3.21.** *Let  $\alpha = x_1 + \dots + x_k \in \mathbb{Z}/2(\mathcal{H}_n)$  with  $x_1, \dots, x_k$  distinct. Then  $\alpha$  is a boundary in  $\inf_*^{C_*(X)}(\mathcal{H})$  if and only if there exists a subset  $\{y_1, \dots, y_m\} \subseteq \mathcal{H}_{n+1}$  with  $y_1, \dots, y_m$  distinct such that  $\{x_1, \dots, x_k\}$  is the set of nodes in the multi-set  $\{d_s(y_i) \mid 0 \leq s \leq n+1, 1 \leq i \leq m\}$  which have odd in-degrees.  $\square$*

## 4. Potential applications.

### 4.1. Potential applications in bio-molecular structures and drug design.

Applications of persistent homology to molecular biology has achieved great success in computer aided drug design [106, 131, 137]. According to [137, Paragraph 0005], theoretical models for the study of the structure-function relationships of biomolecules are conventionally based on purely geometric techniques. Mathematically, these approaches make use of local geometric information such as: coordinates, distances, angles, areas and curvatures for the physical modeling of biomolecular systems. However, conventional purely geometry based models tend to be overwhelmed by too much structural detail and are frequently computationally intractable. Topological approaches to determining the nature of structure-function relationships of biomolecules provide a dramatic simplification compared to conventional geometry based approaches [137, Paragraph 0053].

However, persistent homology neglects chemical and biological information during topological simplification and is thus not as competitive as geometry or physics-based alternatives in quantitative predictions [136]. Element-specific persistent homology, or multi-component persistent homology built on colored biomolecular networks, has been introduced to retain chemical and biological information in topological abstractions [30]. This approach encodes biological properties—such as hydrogen bonds, van der Waals interactions, hydrophilicity, and hydrophobicity—into topological invariants, rendering a potentially revolutionary representation for biomolecules, according to the SIAM news [136].

Recently, we have proposed hypergraph based persistent cohomology (HPC) for molecular representations in drug design [93]. In our HPC model, the protein-ligand interactions at the molecular level are represented as a series of element-specific hypergraphs. Figure 5 illustrates our hypergraph model for a protein-ligand complex with ID 3PB3. Its binding core region is divided into a series of element-specific atom-sets. From these

atom sets, element-specific hypergraphs can be constructed to characterize the interactions between protein atom-sets and ligand atom-sets at the level of atoms. Further, we have proposed a distance-related filtration process as illustrated in Figure 6. With the embedded homology model for hypergraphs, we have developed the hypergraph persistent homology and cohomology for molecular characterization. Molecular features and descriptors can be obtained from hypergraph persistent barcodes and hypergraph enriched barcodes, and this information can be further combined with machine learning models, in particular, the gradient boosting tree (GBT). Our HPC-GBT model has performed well for protein-ligand binding affinity predictions. Its Pearson correlation coefficients (PCCs) for the three PDBbind datasets, including PDBbind-v2007, PDBbind-v2013 and PDBbind-v2016, are consistently better than traditional machine learning models with molecular descriptors.

A molecular representation based on super-hypergraphs could give more flexibility in molecular structure and interaction characterization. Unlike simplicies and hyperedges, super-hyperedges can incorporate local topological structures, that is, subgraphs. This provides a unique way to identify and describe molecular motifs, function groups, and domains. Further, boundary operators can be defined through vertex-deletion and edge-deletion, which provide ways to define different types of homology groups and thus characterize different types of inner topological connections. Moreover, different filtration processes can be defined by considering different scoring functions, which in turn will induce different super-hypergraph based persistent homology/cohomology. Finally, molecular descriptors/fingerprints can be generated from super-hypergraph models and further combined with machine learning models for molecular data analysis in materials, chemistry and biology.

#### 4.2. Potential applications in networks with group interactions.

The abstract of a recent review article [12], citing more than 800 references, reads,

The complexity of many biological, social and technological systems stems from the richness of the interactions among their units. Over the past decades, a variety of complex systems has been successfully described as networks whose interacting pairs of nodes are connected by links. Yet, from human communications to chemical reactions and ecological systems, interactions can often occur in groups of three or more nodes and cannot be described simply in terms of dyads... We review the measures designed to characterize the structure of these systems and the models proposed to generate synthetic structures, such as random and growing bipartite graphs, hypergraphs and simplicial complexes. We introduce the rapidly growing research on higher-order dynamical systems and dynamical topology, discussing the relations between higher-order interactions and collective behavior...

Here we can see that simplicial complexes, a fundamental notion in algebraic topology, has been extensively used for providing representations of higher-order interactions [12, First paragraph of Section 2.1.3]. In some practical problems, the limitations of simplicial complexes due to completeness and vertex-determination present a problem. Hypergraphs provide a more general and unconstrained description of higher-order interactions [12, Paragraphs 2–4, page 7].

Recent progress shows that simplicial homology can be naturally extended as a homology theory on hypergraphs [23]. This provides topological invariants for geometric models using hypergraphs which have had successful applications in biomolecular structures and drug design, described in the previous subsection.

As an extension of hypergraphs, super-hypergraphs would provide more a general and unconstrained description of higher-order interactions. If we assume that the higher-order interactions take place among the nodes in a working graph, which indicates the pre-existence of the pairwise bonds or the primary pairwise links between the nodes, then the most general and unconstrained description of higher-order interactions would be a collection of finite subgraphs of the working graph, which is exactly the topic explored in this article.

## 5. Conclusion.

In this paper, we introduced a new mathematical theory which allows for topological invariants to be applied to broader range of problems, in particular enriching the methods of TDA. This new theory is suitable for both graph data and point cloud data analysis, while overcoming various limitations of the standard persistent homology theory such as the topological noise and the constraining requirements to use data with metric. Using this new theory, the upgraded pipeline of TDA becomes indeterministic in nature allowing for flexibility and adjustments. Moreover, various new topological invariants can be constructed in our flexible setting. As highlighted in Subsections 3.2–3.5, based on this topological approach, more computational tools of algebraic topology will find applications in data science. For example, in algebraic topology the computation of simplicial homology of a space can be largely simplified by homotopically deforming it into a simpler shape, see [80].

As each simplex of a simplicial complex is uniquely determined by its vertices, simplicial complexes cannot model collections of subgraphs. To explore topological structures on space of subgraphs, in this paper we use  $\mathcal{P}$ -sets. Furthermore, we introduce the notion of sup-hypergraph, as a generalization of hypergraphs, which sets a stage for the exploration of topological structures on subgraphs. The homology theory of super-hypergraph, established in Section 2, endows any collection of subgraphs with topological features.

In this work we also use the notion of scoring scheme. As highlighted in Section 3, scoring schemes are used to introduce persistence in an abstract setting without the use of any notion of metric. The classical constructions in persistent simplicial homology theory can be recovered using various scoring schemes.

We should point out that this work presents a theoretic research resulting in a framework that provides an upgraded topological approach to data science, with the aim to foster further interactions between topology and data science.

Further research on super-hypergraphs is needed as this is a new and challenging mathematical concept. From a topological perspective there are many interesting questions to consider such as, the algebraic structure of the homology of these objects as well as the homotopy aspects of super-hypergraphs that are far less-understood. Furthermore, developments in the topological study of super-hypergraphs will feed into new innovative



methods in TDA and other wide-ranging applications. Additionally, the computational complexity of the homology theory of super-hypergraphs is comparable to that of simplicial homology, therefore algorithms methods stemming from our approach will be similarly feasible as classical computation.

## Acknowledgments

The work of JW was supported in part by Natural Science Foundation of China (NSFC grant no. 11971144) and High-level Scientific Research Foundation of Hebei Province. The third author was supported by Nanyang Technological University Startup Grant M4081842 and Singapore Ministry of Education Academic Research fund Tier 1 RG109/19, Tier 2 MOE-T2EP20220-0010, and Tier 2 MOE-T2EP20120-0013. The work of GWW was supported by NIH grant GM126189, NSF grants DMS-1761320, IIS-1900473, and DMS-2052983, and NASA grant 80NSSC21M0023. We wish to thank the referees most warmly for important suggestions that have improved the exposition of this paper.

## Appendix A.: Topological structures associated to graphs.

Throughout mathematics numerous topological and categorical structures on graphs have been explored. In this appendix, we will survey various simplicial complexes associated to graphs that allow us to consider the space of subgraphs of a given graph from a topological perspective.

A *directed (multi-)graph* (or *multi-digraph* or *quiver*) is a pair  $G = (V(G), E(G))$  together with a function  $\text{end}_G : E(G) \rightarrow V(G) \times V(G)$  given by

$$\text{end}_G(e) \mapsto (i(e), t(e))$$

where  $V(G)$  is the *vertex set*,  $E(G)$  is the *edge set*,  $i(e)$  is the *initial* vertex of the edge  $e$ , and  $t(e)$  is the *terminal* vertex of  $e$ .

An *undirected (multi-)graph*<sup>4</sup> is a pair  $G = (V(G), E(G))$  together with a function  $\text{end}_G : E(G) \rightarrow (V(G) \times V(G))/\Sigma_2$  given by

$$\text{end}_G(e) \mapsto \{i(e), t(e)\}$$

where  $(V(G) \times V(G))/\Sigma_2$  is the orbit set of  $(V(G) \times V(G))$  modulo the  $\Sigma_2$ -action given by permuting the coordinates,  $V(G)$  is the *vertex set*,  $E(G)$  is the *edge set*,  $\text{end}_G$  is an *incidence relation* that associates with each edge of  $G$  an unordered pair of, possibly equal, elements of  $V(G)$ . In this definition of a directed/undirected (multi-)graph, the empty graph is allowed.

A *subgraph*  $H$  of a directed/undirected (multi-)graph  $G$  is a graph  $H = (V(H), E(H))$  with  $V(H) \subseteq V(G)$ ,  $E(H) \subseteq E(G)$  and  $\text{end}_H = \text{end}_G|_{E(H)}$ .

<sup>4</sup>We follow the definition of a multi-graph in [51, 10]. In some literature such as [113], a multi-graph is defined by requiring the edge set to be a multi-set. The difference is that the edges between two vertices are labeled by  $E(G)$  together with the incidence map  $\text{end}_G$ . Such a definition coincides with the definition on quiver (as directed multi-graph) [121].

A directed/undirected graph  $G$  is *simple* if  $\text{end}_G$  is injective and the image  $\text{end}_G(E(G))$  is disjoint from the diagonal  $\Delta(V(G))$  in  $V(G) \times V(G)$  or  $(V(G) \times V(G))/\Sigma_2$ . This means that there are no loops or multi-edges between two vertices.

From the perspective of applications, the initial data is represented by a given graph  $G$  and let  $\mathcal{K}$  be a collection of subgraphs of  $G$ . Our goal is to investigate the possible topological structures on  $\mathcal{K}$ . However, before we address this general question, we review some classical constructions of simplicial complexes associated to graphs.

### A.1. Clique complexes.

Typically, the study of collections of subgraphs has focused on measuring how strongly connected different parts of a graph are. A clique (or flag) complex and an independence complex (the clique complex of the complementary graph) are topological spaces that contain information about the connectivity of a graph. These are widely used objects in mathematics and its applications, see [4, 11, 58, 60, 87] for some recent works.

A *complete graph* is a simple graph  $G = (V(G), E(G))$  with the property that every pair of distinct vertices of  $G$  are adjacent in  $G$ .

A *clique* of a graph  $G$  is a complete subgraph of  $G$ .

The *clique complex* of a simple graph  $G$  is the abstract simplicial complex  $\text{Clique}(G)$  whose simplices consist of all cliques of  $G$ . An  $n$ -simplex  $\sigma$  in  $\text{Clique}(G)$  is a clique of  $G$  with  $(n + 1)$  vertices, and a face of a simplex  $\sigma \in \text{Clique}(G)$  is a complete subgraph obtained by deleting some vertices of  $\sigma$ .

When working with a multi-graph  $G = (V, E)$ , the set of cliques  $\text{Clique}(G)$  is generally not a simplicial complex as this requires that all simplices are uniquely determined by their vertex set. For example, let  $G$  be a multi-graph with two vertices  $v$  and  $w$  and two edges  $e_1$  and  $e_2$  joining them. Then

$$\text{Clique}(G) = \{\bar{e}_1, \bar{e}_2, v, w\}$$

has two 1-simplices  $\bar{e}_1$  and  $\bar{e}_2$  sharing the same vertices  $v$  and  $w$ , see Figure 7. Therefore, a more suitable object for describing the topological structure of  $\text{Clique}(G)$  is a  $\Delta$ -set.

For an undirected multi-graph  $G$ , the  $\Delta$ -set structure on  $\text{Clique}(G)$  is given in the following way. Assign a total ordering to  $V(G)$  and define  $\text{Clique}_\#(G)$  to be the set of cliques of  $G$  that have exactly  $n+1$  vertices. For  $\sigma \in \text{Clique}_\#(G)$  with vertices  $v_0 < v_1 < \dots < v_n$ , define  $d_i\sigma = \sigma - v_i$  the subclique of  $\sigma$  obtained by deleting the vertex  $v_i$  and the edges incident to  $v_i$  for  $0 \leq i \leq n$ . It is straightforward to check that  $\text{Clique}_\#(G)$  forms a  $\Delta$ -set<sup>5</sup>.

<sup>5</sup>The definition of the  $\Delta$ -set  $\text{Clique}_\#(G)$  depends on the given order on vertices of  $G$ , but the homology of  $\text{Clique}_\#(G)$  is independent on this choice because the geometric realization of a  $\Delta$ -set is a  $\Delta$ -complex [139, Proposition 1.39, p. 51] in the sense of Hatcher [80].

## A.2. Neighborhood complexes and Jonsson's graph complexes.

We proceed by considering a collection of simplicial complexes associated to graphs which will naturally lead to new constructions suitable for studying spaces of subgraphs. We start with a famous construction of the neighborhood complex of a graph. This was introduced by Lovász [94] in 1978 in his work on Kneser's conjecture which laid the foundations of topological combinatorial by introducing homotopy theoretical methods to combinatorics. Nowadays, the research area of topological combinatorics is very active and fruitful. The generalization by Lovász of the neighborhood complexes to the Hom complex [7, 89], which has the same homotopy type as the clique complex of an exponential graph [54, Remark 3.6], was used in a breakthrough work of Babson and Kozlov [8] to solve the Lovász conjecture which relates the chromatic number of a graph with the homology of its Hom complex. Our theory is based on the exploration of the interplay between topology and combinatorics.

The *neighborhood complex*  $\mathcal{N}(G)$  of a graph  $G$  is a simplicial complex on vertex set  $V(G)$  in which an  $n$ -simplex is a subset of  $V(G)$  with  $n + 1$  vertices such that all vertices are adjacent to an other vertex in  $G$ .

As we discussed in the previous subsection,  $\text{Clique}(G)$  may not be a simplicial complex for a multi-graph  $G$ . However, for any graph  $G$  the neighborhood complex  $\mathcal{N}(G)$  is a simplicial complex.

The topology on the geometric realization of  $\mathcal{N}(G)$  can be quite different from that of  $\text{Clique}(G)$  in general. For example, let  $G$  be a graph with three vertices  $a, b, c$  and two edges given by  $ab$  and  $bc$ . Then  $\mathcal{N}(G) = \{\{a, c\}, \{a\}, \{b\}, \{c\}\}$ , which is not connected, see Figure 8b, and  $\text{Clique}(G) = \{\{a, b\}, \{b, c\}, \{a\}, \{b\}, \{c\}\}$  which is connected, see Figure 8a. This indicates that there are various topological structures one could construct for a given working graph  $G$ .

In [84, p.26], Jonsson defines a graph complex in the following way. A *graph complex*<sup>6</sup> on a finite vertex set  $V$  is a family  $\mathcal{E}$  of simple graphs on the vertex set  $V$  such that  $\mathcal{E}$  is closed under deletion of edges; if  $H \in \mathcal{E}$  and  $e \in H$ , then  $H - e \in \mathcal{E}$ . Identifying  $H = (V, E) \in \mathcal{E}$  with the edge set  $E$ , we may interpret  $\mathcal{E}$  as a simplicial complex. There are potentially different graph complexes on a given vertex set  $V$  because the collection of simple graphs can be chosen in a different way.

With a slight modification to Jonsson's definition, namely adding a hypothesis that the simple graphs in  $\mathcal{E}$  are subgraphs of  $G$ , we retain the central ideas of Jonsson's construction but also gain control over the space of subgraphs. In contrast to clique complexes and neighborhood complexes, the face operations in Jonsson's graph complex are given by deleting edges. Also, the construction of a graph complex is not fully determined by  $G$  as there are various choices for families  $\mathcal{E}$  of simple subgraphs of  $G$  that can form graph complexes. A non-deterministic characteristic of these complexes might be useful in data science as the family  $\mathcal{E}$  can be adjusted for each iteration of the analysis.

<sup>6</sup>Kontsevich also introduced graph complexes with a different definition [86].

### A.3. Path complexes.

Considering hypergraphs as a combinatorial generalization of simplicial complexes allows the construction of the path complex of a given digraph.

The topological exploration of path complexes was first introduced by Shing-Tung Yau and his collaborators in a series of papers [73, 74, 75, 76, 68, 72, 69, 70, 71, 77]. Motivated by ideas from physical applications, A. Dimakis and F. Müller-Hoissen attempted to construct the cohomology of digraphs [53, 52]. They considered path complexes on an intuitive level without a precise definition of the corresponding cochain complex.

In this subsection, we survey the main ideas of path complexes of simple digraphs.

Let  $G$  be a simple digraph. A *directed path* in  $G$  is an alternating sequence  $\lambda = v_0 a_1 v_1 a_2 v_2 \cdots a_k v_k$ , with all vertices  $v_i$  distinct for  $0 \leq i \leq k$  and the edges,  $a_j$ , are incident out of  $v_{j-1}$  and incident into  $v_j$  for  $1 \leq j \leq k$ .

Let  $\mathcal{P}$  be the set of directed paths in  $G$ . We want to associate a combinatorial object to  $G$  built out of directed paths. Since  $G$  is simple, there is at most one edge joining two distinct vertices. So a directed path  $\lambda = v_0 a_1 v_1 a_2 v_2 \cdots a_k v_k$  is determined by its vertices  $v_0, v_1, \dots, v_k$ . Thus we consider  $\lambda = v_0 a_1 v_1 a_2 v_2 \cdots a_k v_k$  as an abstract  $k$ -simplex  $\{v_0, v_1, \dots, v_k\}$ . For  $\mathcal{P}$  to be a simplicial complex, any nonempty subset of  $\{v_0, v_1, \dots, v_k\}$  must be a simplex. In other words, any subsequence  $(v_{i_0}, v_{i_1}, v_{i_2}, \dots, v_{i_t}), 0 \leq i_0 < i_1 < \dots < i_t \leq k$  of  $\lambda$  must form a directed path in  $G$ . This is not true in general. For example, if  $v_0 a_1 v_1 a_2 v_2$  is a directed path in  $G$ , then there may not exist an edge incident out of  $v_0$  and incident into  $v_2$  in  $G$ , that is,  $(v_0, v_2)$  may not form a directed path. Therefore, a structure to consider on the set  $\mathcal{P}$  is that of a hypergraph.

The set  $\mathcal{P}$  becomes a hypergraph with its vertex  $V(G)$  and the hyperedge set given by directed paths in  $G$ . By definition, an abstract simplicial complex is a hypergraph with the additional condition that any nonempty subset of a hyperedge is a hyperedge. Therefore, a hypergraph can be viewed as a simplicial complex with some faces missing, where a hyperedge of cardinal  $k + 1$  is a  $k$ -simplex in the terminology of simplicial complexes. The approaches in [53, 52] and Yau's school lead to the embedded homology of hypergraphs as an extension of simplicial homology theory as introduced in [23].

By allowing vertex repetition in directed paths, we get directed walks. The walk complex  $\mathcal{W}(G)$  for a digraph or quiver (i.e. directed multi-graph)  $G$ , is similar to the path complex but with we replace directed paths with directed walks. Therefore,  $\mathcal{W}(G)$  is an extension of the notion of the nerve of a category in the following sense. Consider a category  $\mathcal{C}$  as a quiver with the composition operation on head-to-tail arrows. Then the nerve of category  $\mathcal{C}$  is the walk complex of quiver  $\mathcal{C}$ .

### A.4. Vertex-deletion topology.

Let  $G$  be a directed/undirected (multi-)graph. Let  $\mathcal{H}$  be a collection of finite subgraphs of  $G$ . Assign to  $\mathcal{H}$  the grading function  $f_v: \mathcal{H} \rightarrow \mathbb{N} = \{0, 1, 2, \dots\}$  given by the size, namely, for

$H \in \mathcal{H}$ , let  $f_v(H) = |V(H)| - 1$ . Let  $\mathcal{H}_n = f_v^{-1}(n)$ . Next step is to define face operations to obtain a topological structure. There are several natural approaches available.

#### A.4.1. Primary vertex-deletion topology.

Assume that the vertex set  $V(G)$  is totally ordered. A geometric way to define face operations is to delete a vertex together with all edges incident to this vertex. More precisely, let  $H \in \mathcal{H}_n$  with vertices  $v_0, v_1, \dots, v_n$ . Define  $d_i(H)$  for  $0 \leq i \leq n$  to be the subgraph of  $H$  by deleting  $v_i$  together with any edges joining with  $v_i$ . This vertex deletion does not ensure that  $d_i(H)$  lies in  $\mathcal{H}_{n-1}$ . Let

$$\Delta(\mathcal{H}) = \{d_{i_1}d_{i_2}\dots d_{i_t}(H) \mid H \in \mathcal{H}, 0 \leq i_1 < i_2 < \dots < i_t \leq |V(H)| - 1\} \tag{30}$$

be the family of subgraphs of  $G$  obtained from  $\mathcal{H}$  together with iterated faces on the subgraphs in  $\mathcal{H}$ . It is straightforward to check that  $\Delta(\mathcal{H})$  is a  $\Delta$ -set, and  $\mathcal{H} \subseteq \Delta(\mathcal{H})$  is a graded subset<sup>7</sup>. Hence  $(\mathcal{H}, \Delta(\mathcal{H}))$  is a super-hypergraph.

**Definition A.1.** Let  $G$  be a directed/undirected (multi-)graph. Let  $\mathcal{H}$  be a collection of finite subgraphs of  $G$ . The *primary vertex-deletion topological structure* on  $\mathcal{H}$  is the super-hypergraph structure defined as above.

Similarly to clique complexes on multi-graphs,  $\Delta(\mathcal{H})$  may not be a simplicial complex in general. Therefore, the notion of a super-hypergraph is the most natural and suitable topological description for  $\mathcal{H}$ .

The super-hypergraph  $(\mathcal{H}, \Delta(\mathcal{H}))$  has a structure of fibrewise topology as follows.

Let

$$V(\mathcal{H}) = \{V(H) \mid H \in \mathcal{H}\} \text{ and } V(\Delta(\mathcal{H})) = \{V(H) \mid H \in \Delta(\mathcal{H})\}$$

be a family of finite subsets of  $V(G)$ . Then  $V(\Delta(\mathcal{H}))$  is a simplicial complex, and  $V(\mathcal{H})$  is a hypergraph whose simplicial closure is  $V(\Delta(\mathcal{H}))$ . Moreover we have a  $\Delta$ -map

$$V: \Delta(\mathcal{H}) \rightarrow V(\Delta(\mathcal{H}))$$

and a morphism of super-hypergraphs

$$V: \mathcal{H} \rightarrow V(\mathcal{H}).$$

By taking geometric realization, we have a continuous map

$$|V|: |\Delta(\mathcal{H})| \rightarrow |V(\Delta(\mathcal{H}))|$$

<sup>7</sup>From the  $\Delta$ -identity (5),  $\Delta(\mathcal{H})$  contains all iterated faces on the subgraphs in  $\mathcal{H}$ , which is the smallest family of subgraphs of  $G$  containing  $\mathcal{H}$  that is closed under the face operation.

which is a fibrewise topology in the sense of James [83].

Clique complexes are typical examples of primary vertex-deletion topology, where  $\mathcal{H}$  is given by cliques in a graph  $G$ . In this case,  $\mathcal{H}$  itself is already a  $\Delta$ -set so  $\mathcal{H} = \Delta(\mathcal{H})$  and the map  $V: \mathcal{H} \rightarrow V(\mathcal{H})$  is an isomorphism.

The neighborhood complex is another good example that admits a fibrewise topological structure as follows. Let

$$\overline{\mathcal{N}(G)} = \{H \mid H \text{ is a subgraph of } G \text{ and } V(H) \in \mathcal{N}(G)\}. \tag{31}$$

Then it is straightforward to check that

$$\overline{\mathcal{N}(G)} = \Delta(\overline{\mathcal{N}(G)}) \text{ and } V(\overline{\mathcal{N}(G)}) = \mathcal{N}(G) \tag{32}$$

with a continuous map

$$|V|: |\overline{\mathcal{N}(G)}| \rightarrow |\mathcal{N}(G)| \tag{33}$$

which is called a *fibrewise neighborhood topology* of  $G$ .

#### A.4.2. Secondary vertex-deletion topology.

Consider the path complex of a simple digraph  $G$  and its face operation  $d_i$ . Let  $\lambda = v_0 a_1 v_1 a_2 v_2 \dots a_n v_n$  be a directed path. Then  $d_i(\lambda)$  is given by deleting the vertex  $v_i$ . However, we have to add back the directed edge from  $v_{i-1}$  to  $v_{i+1}$  provided that it exists to ensure that  $d_i(\lambda) \in \mathcal{P}_n$ . This gives a different type of topological structure, in which we need to redefine the edges to match the vertex removal of the face operation. This can be generalized in the following way.

Let  $G$  be a directed/undirected simple graph and let  $\mathcal{H}$  be a family of finite subgraphs of  $G$ . Let the vertex set  $V(G)$  be totally ordered. For  $H \in \mathcal{H}$ , as a finite subgraph of  $G$  with vertices  $v_0 < v_1 < \dots < v_n$ , define  $d_i H$  to be the subgraph of  $G$  by removing  $v_i$  from  $H$  and adding the edge between  $v_{i-1}$  and  $v_{i+1}$  if it exists. Then  $\mathcal{H}$  forms a super-hypergraph in a similar way as in the case of primary vertex-deletion topology. Here the notion of a super-hypergraph is necessary because there could be two subgraphs in  $\mathcal{H}$  sharing same vertices. For instance, if there is an edge joining two distinct vertices  $v$  and  $w$  in  $G$ , then the subgraphs consist of two vertices  $v$  and  $w$  with the edge joining them and without the edge joining them, respectively, are different.

**Definition A.2.** Let  $G$  be a directed/undirected simple graph. Let  $\mathcal{H}$  be a collection of finite subgraphs of  $G$ . The *secondary vertex-deletion topological structure* on  $\mathcal{H}$  is the super-hypergraph structure defined as above.

The secondary vertex-deletion topology naturally applies to subgraphs of a simple graph. However, to construct a topological structure on a space of subgraphs of a multi-graph in this a way would be more complicated.

There are other possible topological structures on special families of subgraphs. Analogously to various techniques developed in simplicial homotopy theory, for special families of subgraphs having good patterns, one could delete more than one vertex under each elementary face operation  $d_j$ .

**A.5. Edge-deletion topology.**

Let  $G$  be a directed/undirected (multi-)graph and  $\mathcal{H}$  be a collection of finite subgraphs of  $G$ . Another reasonable way to assign the grading function  $f_e: \mathcal{H} \rightarrow \mathbb{N} = \{0, 1, 2, \dots\}$  is by counting edges, that is, for  $H \in \mathcal{H}$  let  $f_e(H) = |E(H)| - 1$ . Then  $\mathcal{H}_n = f_e^{-1}(n)$ . Note that a subgraph  $H$  of  $G$  is uniquely determined by its edge set  $E(H)$ . We do not need to use the notion of a  $n$ -set for describing topological structure on  $\mathcal{H}$  from edge-deletion. If  $\mathcal{H}$  is closed under edge-deletion operation, then it forms a simplicial complex, which is exactly a path complex in the sense of Jonsson. Otherwise,  $\mathcal{H}$  is only a hypergraph.

For a fixed graph  $G$ , the edge-deletion topology could be quite different from the vertex-deletion topology because already the grading functions  $f_v$  and  $f_e$  could be quite different. The edge-deletion operation may not commute with the vertex-deletion operation, so the relationship between the edge-deletion topology and the vertex-deletion topology is not immediately clear. To better understand these structures, more exploration of the relationship between different topological structures on families of subgraphs is needed.

Finally, we should point out that there are many other ways to introduce topological structures on subgraphs, for example following ideas related to Hom complexes. The frontier of research in topological combinatorics has potential to provide new mathematical tools in data science.

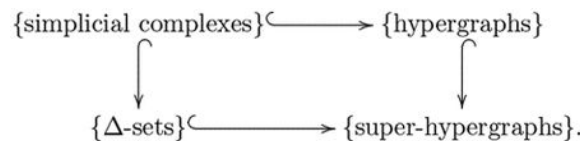
**Appendix B.: The connections of the concepts in the article.**

There are various concepts in the article, including new concepts such as super-hypergraph. We highlight the connections between them in this appendix.

The following statements are well known and important:

- A simple graph is a 1-dimensional simplicial complex.
- A multi-graph with its vertices totally ordered is a 1-dimensional  $n$ -set/  $n$ -complex.
- A quiver (multi-digraph) is a 1-dimensional  $n$ -set/  $n$ -complex.

In the following table, the arrow  $\hookrightarrow$  means an inclusion of sets.



## REFERENCES

- [1]. Dionysus: The persistent homology software, Software available at <http://www.mrzv.org/software/dionysus>.
- [2]. Adams H, Emerson T, Kirby M, Neville R, Peterson C, Shipman P, Chepushtanova S, Hanson E, Motta F and Ziegelmeier L, Persistence images: A stable vector representation of persistent homology, *J. Mach. Learn. Res.*, 18 (2017), Paper No. 8, 35 pp.
- [3]. Adcock A, Carlsson E and Carlsson G, The ring of algebraic functions on persistence bar codes, *Homology Homotopy Appl.*, 18 (2016), 381–402.
- [4]. Aharoni R, Berger E and Ziv R, Independent systems of representatives in weighted graphs, *Combinatorica*, 27 (2007), 253–267.
- [5]. Ahmed M, Fasy BT and Wenk C, Local persistent homology based distance between maps, In *Proceedings of the 22nd ACM SIGSPATIAL International Conference on Advances in Geographic Information Systems*, ACM, (2014), 43–52.
- [6]. Anirudh R, Thiagarajan JJ, Kim I and Polonik W, Autism spectrum disorder classification using graph kernels on multidimensional time series, preprint, arXiv:1611.09897.
- [7]. Babson E and Kozlov DN, Complexes of graph homomorphisms, *Israel J. Math.*, 152 (2006), 285–312.
- [8]. Babson E and Kozlov DN, Proof of the Lovász conjecture, *Ann. of Math.*, 165 (2007), 965–1007.
- [9]. Bae W, Yoo JJ and Ye JC, Beyond deep residual learning for image restoration: Persistent homology-guided manifold simplification, In *IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, (2017), 1141–1149.
- [10]. Balakrishnan R and Ranganathan K, *A Textbook of Graph Theory*, 2<sup>nd</sup> edition, Universitext, Springer, New York, 2012.
- [11]. Barmak JA, Star clusters in independence complexes of graphs, *Adv. Math.*, 241 (2013), 33–57.
- [12]. Battiston F, Cencetti G, Iacopini I, Latora V, Lucas M, Patania A, Young J-G and Petri G, Networks beyond pairwise interactions: Structure and dynamics, *Phys. Rep.*, 874 (2020), 1–92.
- [13]. Bauer U, Ripser: A lean C++ code for the computation of Vietoris-Rips persistence barcodes, Software available at <https://github.com/Ripser/ripser>.
- [14]. Bauer U, Kerber M and Reininghaus J, Distributed computation of persistent homology, In *Meeting on Algorithm Engineering and Experiments (ALENEX)*, SIAM, (2014), 31–38.
- [15]. Bauer U, Kerber M, Reininghaus J and Wagner H, PHAT–persistent homology algorithms toolbox, In *Mathematical software ICMS*, 8592 (2014), 137–143.
- [16]. Bendich P, Cohen-Steiner D, Edelsbrunner H, Harer J and Morozov D, Inferring local homology from sampled stratified spaces, In *IEEE Symposium on Foundations of Computer Science (FOCS'07)*, (2007), 536–546.
- [17]. Bendich P, Edelsbrunner H and Kerber M, Computing robustness and persistence for images, *IEEE Transactions on Visualization and Computer Graphics*, 16 (2010), 1251–1260. [PubMed: 20975165]
- [18]. Bendich P, Gasparovic E, Harer J, Izmailov R and Ness L, Multi-scale local shape analysis and feature selection in machine learning applications, In *International Joint Conference on Neural Networks (IJCNN)*, IEEE, (2015), 1–8.
- [19]. Bendich P, Wang B and Mukherjee S, Local homology transfer and stratification learning, In *Twenty-Third Annual ACM-SIAM Symposium on Discrete Algorithms*, (2012), 1355–1370.
- [20]. Bergomi MG, Ferri M, Vertechi P and Zuffi L, Beyond topological persistence: Starting from networks, *Mathematics*, 9 (2021).
- [21]. Binchi J, Merelli E, Rucco M, Petri G and Vaccarino F, jholes: A tool for understanding biological complex networks via clique weight rank persistent homology, *Electron. Notes Theor. Comput. Sci.*, 306 (2014), 5–18.
- [22]. Bonis T, Ovsjanikov M, Oudot S and Chazal F, Persistence-based pooling for shape pose recognition, In *Computational Topology in Image Context*, 9667 (2016), 19–29.
- [23]. Bressan S, Li J, Ren S and Wu J, The embedded homology of hypergraphs and applications, *Asian J. Math.*, 23 (2019), 479–500.



- [24]. Bubenik P, Statistical topological data analysis using persistence landscapes, *J. Mach. Learn. Res.*, 16 (2015), 77–102.
- [25]. Bubenik P and Kim PT, A statistical approach to persistent homology, *Homology Homotopy Appl.*, 9 (2007), 337–362.
- [26]. Bubenik P and Vergili T, Topological spaces of persistence modules and their properties, *J. Appl. Comput. Topol.*, 2 (2018), 233–269.
- [27]. Cameron PJ, Automorphisms and cohomology of switching classes, *J. Combinatorial Theory Ser. B*, 22 (1977), 297–298.
- [28]. Cameron PJ, Cohomological aspects of two-graphs, *Math. Z.*, 157 (1977), 101–119.
- [29]. Cang Z, Mu L and Wei G-W, Representability of algebraic topology for biomolecules in machine learning based scoring and virtual screening, *PLOS Computational Biology*, 14 (2018), 1–44.
- [30]. Cang Z and Wei G, Integration of element specific persistent homology and machine learning for protein-ligand binding affinity prediction, *International Journal for Numerical Methods in Biomedical Engineering*, 34 (2018), e2914.
- [31]. Cang Z and Wei G-W, Analysis and prediction of protein folding energy changes upon mutation by element specific persistent homology, *Bioinformatics*, 33 (2017), 3549–3557. [PubMed: 29036440]
- [32]. Cang Z and Wei G-W, Topologynet: Topology based deep convolutional and multi-task neural networks for biomolecular property predictions, *PLOS Computational Biology*, 13 (2017), 1–27.
- [33]. Carlsson G, Ishkhanov T, Silva V and Zomorodian A, On the local behavior of spaces of natural images, *Int. J. Comput. Vis.*, 76 (2008), 1–12.
- [34]. Carlsson G, Singh G and Zomorodian A, Computing multidimensional persistence, In *Algorithms and Computation*, 5878 (2009), 730–739.
- [35]. Carlsson G and Zomorodian A, The theory of multidimensional persistence, *Discrete Comput. Geom.*, 42 (2009), 71–93.
- [36]. Carlsson G, Topology and data, *Bull. Amer. Math. Soc. (N.S.)*, 46 (2009), 255–308.
- [37]. Cerri A and Landi C, The persistence space in multidimensional persistent homology, In *Discrete Geometry for Computer Imagery*, 7749 (2013), 180–191.
- [38]. Chazal F, Cohen-Steiner D, Glisse M, Guibas LJ and Oudot SY, Proximity of persistence modules and their diagrams, In *SCG '09: Proceedings of the twenty-fifth annual symposium on Computational Geometry*, (2009), 237–246.
- [39]. Chazal F, de Silva V, Glisse M and Oudot S, *The Structure and Stability of Persistence Modules*, SpringerBriefs in Mathematics, Springer, [Cham], 2016.
- [40]. Chazal F, Fasy B, Lecci F, Michel B, Rinaldo A and Wasserman L, Subsampling methods for persistent homology, In *Proceedings of the 32nd International Conference on Machine Learning*, (eds. Bach F and Blei D), PMLR, Lille, France, 37 (2015), 2143–2151.
- [41]. Chazal F and Michel B, An introduction to topological data analysis: Fundamental and practical aspects for data scientists, *Front. Artif. Intell.*, 2021.
- [42]. Cheng Y and Wells AL Jr., Switching classes of directed graphs, *J. Combin. Theory Ser. B*, 40 (1986), 169–186.
- [43]. Chevyrev I, Nanda V and Oberhauser H, Persistence paths and signature features in topological data analysis, *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 42 (2018), 192–202. [PubMed: 30530312]
- [44]. Chung FRK and Graham RL, Cohomological aspects of hypergraphs, *Trans. Amer. Math. Soc.*, 334 (1992), 365–388.
- [45]. Cohen-Steiner D, Edelsbrunner H and Morozov D, Vines and vineyards by updating persistence in linear time, In *Computational geometry (SCG'06)*, (2006), 119–126.
- [46]. Crawley-Boevey W, Decomposition of pointwise finite-dimensional persistence modules, *J. Algebra Appl.*, 14 (2015), 1550066, 8 pp.
- [47]. Curtis EB, Simplicial homotopy theory, *Advances in Math.*, 6 (1971), 107–209.
- [48]. de Silva V and Ghrist R, Homological sensor networks, *Notices Amer. Math. Soc.*, 54 (2007), 10–17.

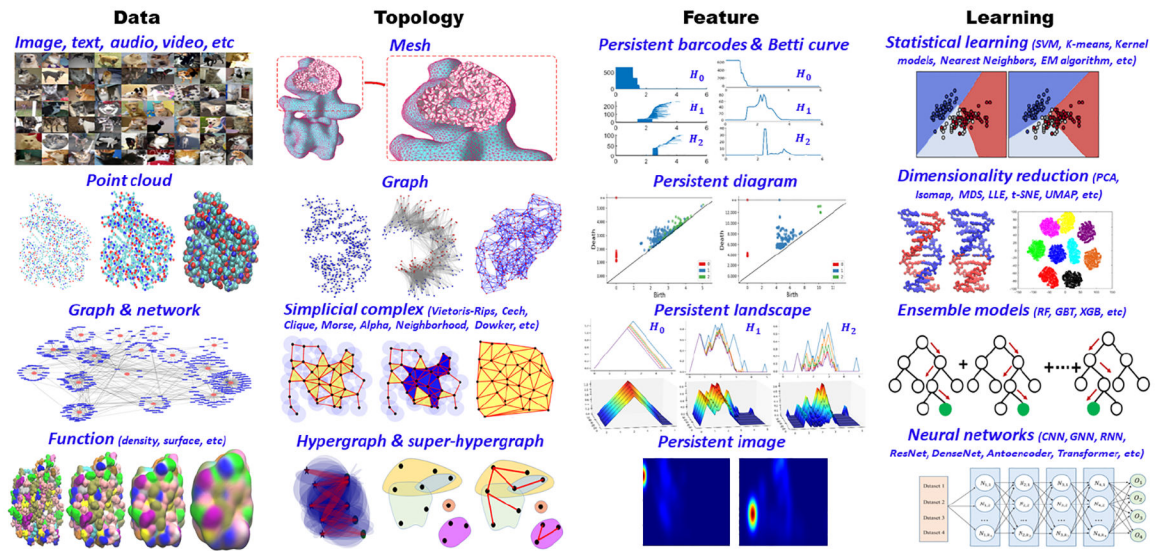
- [49]. De Silva V, Morozov D and Vejdemo-Johansson M, Persistent cohomology and circular coordinates, *Discrete Comput. Geom*, 45 (2011), 737–759.
- [50]. Dey TK and Mandal S, Protein classification with improved topological data analysis, In *LIPIcs. Leibniz Int. Proc. Inform*, 113 (2018), 13 pp.
- [51]. Diestel R, *Graph Theory*, vol. 173 of *Graduate Texts in Mathematics*, 5<sup>th</sup> edition, *Graduate Texts in Mathematics*, 173. Springer, Berlin, 2017.
- [52]. Dimakis A and Müller-Hoissen, Differential calculus and gauge theory on finite sets, *J. Phys. A*, 27 (1994), 3159–3178.
- [53]. Dimakis A and Müller-Hoissen, Discrete differential calculus: Graphs, topologies, and gauge theory, *J. Math. Phys.* 35 (1994), 6703–6735.
- [54]. Dochtermann A, Hom complexes and homotopy theory in the category of graphs, *European J. Combin.* 30 (2009), 490–509.
- [55]. Duval AM and Reiner V, Shifted simplicial complexes are Laplacian integral, *Trans. Amer. Math. Soc.* 354 (2002), 4313–4344.
- [56]. Dwork B, On the zeta function of a hypersurface, *Inst. Hautes Etudes Sci. Publ. Math.* (1962), 5–68.
- [57]. Edelsbrunner H, Letscher D and Zomorodian A, Topological persistence and simplification, *Discrete Comput. Geom*, 28 (2002), 511–533.
- [58]. Ehrenborg R and Hetyei G, The topology of the independence complex, *European J. Combin.* 27 (2006), 906–923.
- [59]. Emtander E, Betti numbers of hypergraphs, *Comm. Algebra*, 37 (2009), 1545–1571.
- [60]. Engström A, Independence complexes of claw-free graphs, *European J. Combin.* 29 (2008), 234–241.
- [61]. Fasy BT, Kim J, Lecci F and Maria C, Introduction to the r package tda, preprint, arXiv:1411.1830.
- [62]. Fasy BT and Wang B, Exploring persistent local homology in topological data analysis, In *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, (2016), 6430–6434.
- [63]. Forman R, Morse theory for cell complexes, *Adv. Math.* 134 (1998), 90–145.
- [64]. Frosini P and Landi C, Persistent Betti numbers for a noise tolerant shape-based approach to image retrieval, *Computer Analysis of Images and Patterns*, 6854 (2011), 294–301.
- [65]. Gabriel P, Unzerlegbare Darstellungen. I, *Manuscripta Math.*, 6 (1972), 71–103.
- [66]. Ghrist R, Barcodes: The persistent topology of data, *Bull. Amer. Math. Soc.* 45 (2008), 61–75.
- [67]. Giansiracusa N, Giansiracusa R and Moon C, Persistent homology machine learning for fingerprint classification, *IEEE International Conference On Machine Learning And Applications (ICMLA)*, 2019.
- [68]. Grigorian A, Lin Y, Muranov Y and Yau S-T, Homologies of path complexes and digraphs, *ArXiv*.
- [69]. Grigorian A, Muranov Y and Yau S-T, Homologies of digraphs and künnetth formulas, *Comm. Anal. Geom.* 25 (2017), 969–1018.
- [70]. Grigor'yan A, Muranov YV and Yau S-T, Graphs associated with simplicial complexes, *Homology Homotopy Appl.* 16 (2014), 295–311.
- [71]. Grigor'yan A and Muranov YV, Cohomology theories of simplicial complexes, algebras, and digraphs.
- [72]. Grigor'yan AA, Lin I, Muranov YV and Yau S, Path complexes and their homologies, *Fundam. Prikl. Mat.* 21 (2016), 79–128.
- [73]. Grigor'yan A, Jimenez R, Muranov Y and Yau S-T, On the path homology theory of digraphs and Eilenberg-Steenrod axioms, *Homology Homotopy Appl.*, 20 (2018), 179–205.
- [74]. Grigor'yan A, Jimenez R, Muranov Y and Yau S-T, Homology of path complexes and hypergraphs, *Topology Appl.*, 267 (2019), 106877, 25 pp.
- [75]. Grigor'yan A, Lin Y, Muranov Y and Yau S-T, Homotopy theory for digraphs, *Pure Appl. Math. Q.* 10 (2014), 619–674.

- [76]. Grigor'yan A, Lin Y, Muranov Y and Yau S-T, Cohomology of digraphs and (undirected) graphs, *Asian J. Math*, 19 (2015), 887–931.
- [77]. Grigor'yan A, Muranov Y, Vershinin V and Yau S-T, Path homology theory of multi-graphs and quivers, *Forum Math*, 30 (2018), 1319–1337.
- [78]. Guo W, Manohar K, Brunton SL and Banerjee AG, Sparse-tda: Sparse realization of topological data analysis for multi-way classification, *IEEE Transactions on Knowledge and Data Engineering*, 30 (2018), 1403–1408.
- [79]. Han YS, Yoo J and Ye JC, Deep residual learning for compressed sensing ct reconstruction via persistent homology analysis, preprint, arXiv:1611.06391.
- [80]. Hatcher A, *Algebraic Topology*, Cambridge University Press, Cambridge, 2002.
- [81]. Hiraoka Y, Nakamura T, Hirata A, Escolar EG, Matsue K and Nishiura Y, Hierarchical structures of amorphous solids characterized by persistent homology, *Proceedings of the National Academy of Sciences*, 113 (2016), 7035–7040.
- [82]. Horak D and Jost J, Spectra of combinatorial Laplace operators on simplicial complexes, *Adv. Math*, 244 (2013), 303–336.
- [83]. James IM, *Fibrewise Topology*, Cambridge Tracts in Mathematics, 91. Cambridge University Press, Cambridge, 1989.
- [84]. Jonsson J, *Simplicial Complexes of Graphs*, Lecture Notes in Mathematics, 1928. Springer-Verlag, Berlin, 2008.
- [85]. Kališnik S, Tropical coordinates on the space of persistence barcodes, *Foundations of Computational Mathematics*, 1–29.
- [86]. Kontsevich M, Derived Grothendieck-Teichmüller group and graph complexes [after T. Willwacher], *Séminaire Bourbaki*, 2016/2017 (2019), 1120–1135; Exposé, 1126 (2019), 183–211.
- [87]. Kozlov DN, Complexes of directed trees, *J. Combin. Theory Ser. A*, 88 (1999), 112–122.
- [88]. Kozlov DN, Discrete Morse theory for free chain complexes, *C. R. Math*, 340 (2005), 867–872.
- [89]. Kozlov DN, Simple homotopy types of Hom-complexes, neighborhood complexes, Lovász complexes, and atom crosscut complexes, *Topology Appl*, 153 (2006), 2445–2454.
- [90]. Kramár M, Levanger R, Tithof J, Suri B, Xu M, Paul M, Schatz MF and Mischaikow K, Analysis of Kolmogorov flow and Rayleigh-Bénard convection using persistent homology, *Phys. D*, 334 (2016), 82–98.
- [91]. Kramár M, Goulet A, Kondic L and Mischaikow K, Persistence of force networks in compressed granular media, *Physical Review. E, Statistical, Nonlinear, and Soft Matter Physics*, 87 (2013), 042207. [PubMed: 23679407]
- [92]. Li C, Ovsjanikov M and Chazal F, Persistence-based structural recognition, In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, (2014), 1995–2002.
- [93]. Liu X, Wang X, Wu J and Xia K, Hypergraph-based persistent cohomology (HPC) for molecular representations in drug design, *Briefings in Bioinformatics*, 22 (2021).
- [94]. Lovász L, Kneser's conjecture, chromatic number, and homotopy, *J. Combin. Theory Ser. A*, 25 (1978), 319–324.
- [95]. Makarenko N, Kalimoldayev M, Pak I and Yessenaliyeva A, Texture recognition by the methods of topological data analysis, *Open Engineering*, 6 (2016).
- [96]. Mallows CL and Sloane NJA, Two-graphs, switching classes and Euler graphs are equal in number, *SIAM J. Appl. Math*, 28 (1975), 876–880.
- [97]. Maria C, Filtered complexes, in *GUDHI User and Reference Manual*, GUDHI Editorial Board, 2015, [https://gudhi.inria.fr/doc/3.4.1/group\\_\\_simplex\\_\\_tree.html](https://gudhi.inria.fr/doc/3.4.1/group__simplex__tree.html).
- [98]. Meng Z and Xia K, Persistent spectral-based machine learning (perspect ml) for protein-ligand binding affinity prediction, *Science Advances*, 7 (2021), eabc5329. [PubMed: 33962954]
- [99]. Mielants W and Leemans H,  $Z_2$ -cohomology of projective spaces of odd order, In *Combinatorics '81 (Rome, 1981)*, *Ann. Discrete Math.*, North-Holland, Amsterdam-New York, 18 (1983), 635–651.
- [100]. Mischaikow K, Mrozek M, Reiss J and Szymczak A, Construction of symbolic dynamics from experimental time series, *Physical Review Letters*, 82 (1999), 1144–1147.

- [101]. Mischaikow K and Nanda V, Morse theory for filtrations and efficient computation of persistent homology, *Discrete Comput. Geom*, 50 (2013), 330–353.
- [102]. Munkres JR, *Elements of Algebraic Topology*, Addison-Wesley Publishing Company, Menlo Park, CA, 1984.
- [103]. Nakamura T, Hiraoka Y, Hirata A, Escolar EG and Nishiura Y, Persistent homology and many-body atomic structure for medium-range order in the glass, *Nanotechnology*, 26 (2015), 304001. [PubMed: 26150288]
- [104]. Nanda V, Perseus: The persistent homology software, Software available at <http://www.sas.upenn.edu/~vnanda/perseus>.
- [105]. Nguyen DD, Cang ZX and Wei GW, A review of mathematical representations of biomolecular data, *Physical Chemistry Chemical Physics*, 2020.
- [106]. Nguyen D, Gao K, Wang M and Wei G-W, MathDL: Mathematical deep learning for D3R grand challenge 4, *Journal of Computer-Aided Molecular Design*, 34 (2020), 131–147. [PubMed: 31734815]
- [107]. Niyogi P, Smale S and Weinberger S, A topological view of unsupervised learning from noisy data, *SIAM J. Comput.*, 40 (2011), 646–663.
- [108]. Obayashi I, Hiraoka Y and Kimura M, Persistence diagrams with linear machine learning models, *J. Appl. Comput. Topol*, 1 (2018), 421–449.
- [109]. Pachauri D, Hinrichs C, Chung MK, Johnson SC and Singh V, Topology-based kernels with application to inference problems in Alzheimer’s disease, *IEEE Transactions on Medical Imaging*, 30 (2011), 1760–1770. [PubMed: 21536520]
- [110]. Parks A, Lipscomb S and VA NSWCD., Homology and Hypergraph Acyclicity: A Combinatorial Invariant For Hypergraphs, Defense Technical Information Center, 1991, <https://apps.dtic.mil/sti/citations/ADA241584>.
- [111]. Pokorny FT, Ek CH, Kjellström H and Kragic D, Persistent homology for learning densities with bounded support, In *Proceedings of the 25th International Conference on Neural Information Processing Systems, NIPS’12*, Curran Associates Inc., Red Hook, NY, USA, 2 (2012), 1817–1825.
- [112]. Qaiser T, Tsang YW, Taniyama D, Sakamoto N, Nakane K, Epstein D and Rajpoot N, Fast and accurate tumor segmentation of histology images using persistent homology and deep convolutional features, *Medical Image Analysis*, 55 (2019), 1–14. [PubMed: 30991188]
- [113]. Rahman MS, *Basic Graph Theory, Undergraduate Topics in Computer Science*, Springer, Cham, 2017.
- [114]. Rebalá G, Ravi A and Churiwala S, *An Introduction to Machine Learning*, 2019.
- [115]. Reininghaus J, Günther D, Hotz I, Prohaska S and Hege H-C, TADD: A computational framework for data analysis using discrete Morse theory, In *Mathematical Software—ICMS 2010, Lecture Notes in Comput. Sci.*, 6327 (2010), 198–208.
- [116]. Ren S, Wu C and Wu J, Weighted persistent homology, *Rocky Mountain J. Math*, 48 (2018), 2661–2687.
- [117]. Robins V and Turner K, Principal component analysis of persistent homology rank functions with case studies of spatial point patterns, sphere packing and colloids, *Phys. D*, 334 (2016), 99–117.
- [118]. Robins V, Computational topology for point data: Betti numbers of  $\alpha$ -shapes, In *Morphology of Condensed Matter*, 600 (2002), 261–274.
- [119]. Saadatfar M, Takeuchi H, Robins V, Francois N and Hiraoka Y, Pore configuration landscape of granular crystallization, *Nature communications*, 8 (2017), 15082.
- [120]. Said A and Torra V, *Data Science in Practice, Studies in Big Data*, Springer International Publishing, 2019.
- [121]. Schiffler R, *Quiver Representations, CMS Books in Mathematics/Ouvrages de Mathématiques de la SMC*, Springer, Cham, 2014.
- [122]. Seidel JJ, A survey of two-graphs, In *Colloquio Internazionale Sulle Teorie Combinatorie (Rome, 1973), Tomo I, Atti dei Convegni Lincei*, (1976), 481–511.

- [123]. Seidel JJ and Taylor DE, Two-graphs, a second survey, In Algebraic Methods in Graph Theory, Vol. I, II (Szeged, 1978), Colloq. Math. Soc. János Bolyai, 25 (1981), 689–711.
- [124]. Seversky LM, Davis S and Berger M, On time-series topological data analysis: New data and opportunities, In IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW), (2016), 1014–1022.
- [125]. Singh G, Memoli F, Ishkhanov T, Sapiro G, Carlsson G and Ringach DL, Topological analysis of population activity in visual cortex, Journal of Vision, 8 (2008).
- [126]. Skraba P, Ovsjanikov M, Chazal F and Guibas L, Persistence-based segmentation of deformable shapes, In IEEE Computer Society Conference on Computer Vision and Pattern Recognition - Workshops, (2010), 45–52.
- [127]. Tausz A, Vejdemo-Johansson M and Adams H, Javaplex: A research software package for persistent (co)homology, Software available at <http://code.google.com/p/javaplex>, 2011.
- [128]. Tierny J, Topological Data Analysis for Scientific Visualization, Springer-Verlag, Berlin, 2017.
- [129]. Turner K, Mukherjee S and Boyer DM, Persistent homology transform for modeling shapes and surfaces, Inf. Inference, 3 (2014), 310–344.
- [130]. Umeda Y, Time series classification via topological data analysis, Transactions of the Japanese Society for Artificial Intelligence, 32 (2017), 1–12.
- [131]. Wang M, Cang Z and Wei G-W, A topology-based network tree for the prediction of protein-protein binding affinity changes following mutation, Nature Machine Intelligence, 2 (2020), 116–123.
- [132]. Wang R, Nguyen DD and Wei G-W, Persistent spectral graph, Int. J. Numer. Methods Biomed. Eng, 36 (2020), e3376, 27 pp.
- [133]. Wang Y, Ombao H and Chung MK et al. , Persistence landscape of functional signal and its application to epileptic electroencephalogram data, ENAR Distinguished Student Paper Award.
- [134]. Wee J and Xia K, Forman persistent Ricci curvature (FPRC)-based machine learning models for protein–ligand binding affinity prediction, Briefings in Bioinformatics, 22 (2021).
- [135]. Wee J and Xia K, Ollivier persistent ricci curvature-based machine learning for the protein–ligand binding affinity prediction, J. Chem. Inf. Model, 61 (2021), 1617–1626. [PubMed: 33724038]
- [136]. Wei G-W, Persistent homology analysis of biomolecular data, SIAM NEWS, 2017, <https://sinews.siam.org/Details-Page/persistent-homology-analysis-of-biomolecular-data>.
- [137]. Wei G, Nguyen D and Cang Z, System and methods for machine learning for drug design and discovery, US Patent App., 16 (2019), 239–327.
- [138]. Wells AL Jr., Even signings, signed switching classes, and  $(-1, 1)$ -matrices, J. Combin. Theory Ser. B, 36 (1984), 194–212.
- [139]. Wu J, Simplicial objects and homotopy groups, In Braids, Lect. Notes Ser. Inst. Math. Sci. Natl. Univ. Singap, World Sci. Publ., Hackensack, NJ, 19 (2010), 31–181.
- [140]. Xia K and Wei G-W, Persistent homology analysis of protein structure, flexibility, and folding, Int. J. Numer. Methods Biomed. Eng, 30 (2014), 814–844.
- [141]. Yadav N, Yadav A and Kumar M et al., An Introduction to Neural Network Methods for Differential Equations, SpringerBriefs in Applied Sciences and Technology. Springer, Dordrecht, 2015.
- [142]. Zaslavsky T, Characterizations of signed graphs, J. Graph Theory, 5 (1981), 401–406.
- [143]. Zeppelzauer M, Zieli ski B, Juda M and Seidl M, A study on topological descriptors for the analysis of 3d surface texture, Computer Vision and Image Understanding, 167 (2018), 74–88.
- [144]. Zhang ZF, Song Y, Cui HC, Wu J, Schwartz F and Qi HR, Early mastitis diagnosis through topological analysis of biosignals from low-voltage alternate current electrokinetics, In Engineering in Medicine and Biology Society (EMBC), 2015 37th Annual International Conference of the IEEE, (2015), 542–545.
- [145]. Zhou Z, Huang YZ, Wang L and Tan TN, Exploring generalized shape analysis by topological representations, Pattern Recognition Letters, 87 (2017), 177–185.
- [146]. Zhu XJ, Persistent homology: An introduction and a new text representation for natural language processing., In IJCAI, (2013), 1953–1959.

- [147]. Zielinski B, Juda M and Zeppelzauer M, Persistence codebooks for topological data analysis, *Artificial Intelligence Review* volume, 54 (2021), 1969–2009.
- [148]. Zomorodian A and Carlsson G, Computing persistent homology, *Discrete Comput. Geom*, 33 (2005), 249–274.



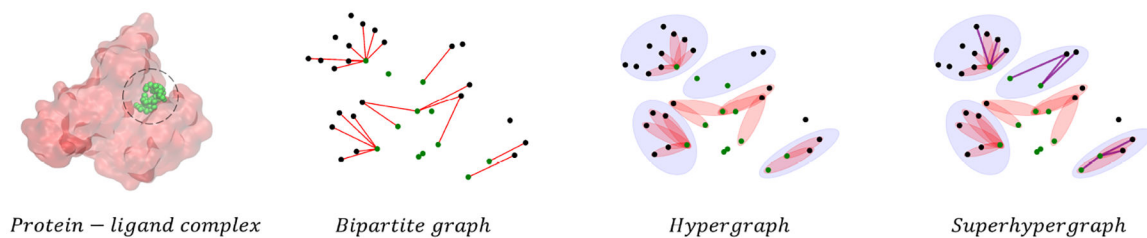
**Figure 1.** Illustration of TDA and TDA-based learning models for data analysis. Generally speaking, all TDA-based learning models have four components, including data, topology, feature and learning. More specifically, data is collected and preprocessed firstly. Second, topological representations and models are constructed to describe the inner structural and interactional information of the data. Note that efficient representations are of key importance to machine learning. Third, a series of topological features are generated by using persistent homology models. Topological-invariant-based features provide a better characterization of the most fundamental and intrinsic properties of the data, thus they have a better generalizability and transferability for machine learning models. Finally, the topological features are combined with machine learning models for various classification and regression tasks.

Author Manuscript

Author Manuscript

Author Manuscript

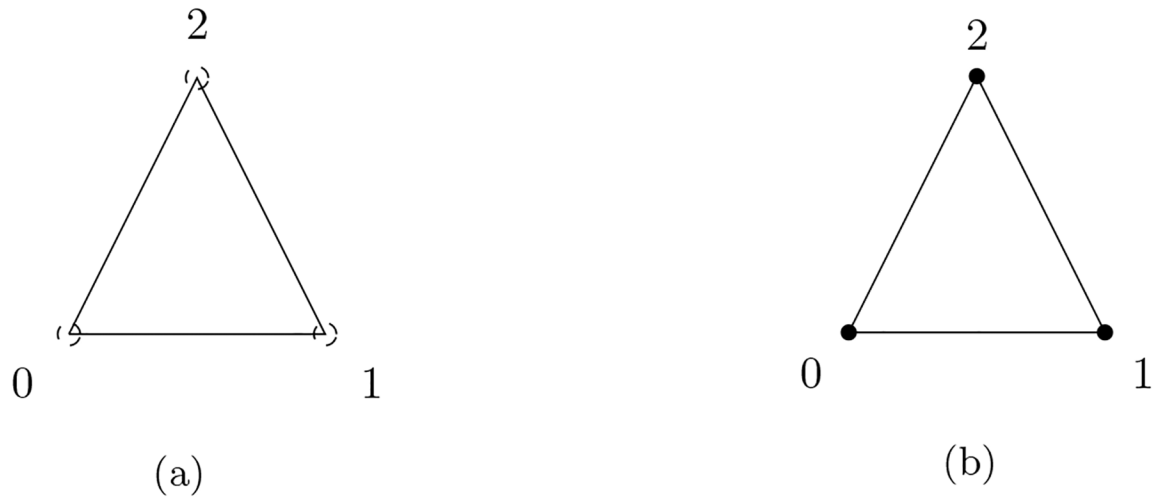
Author Manuscript



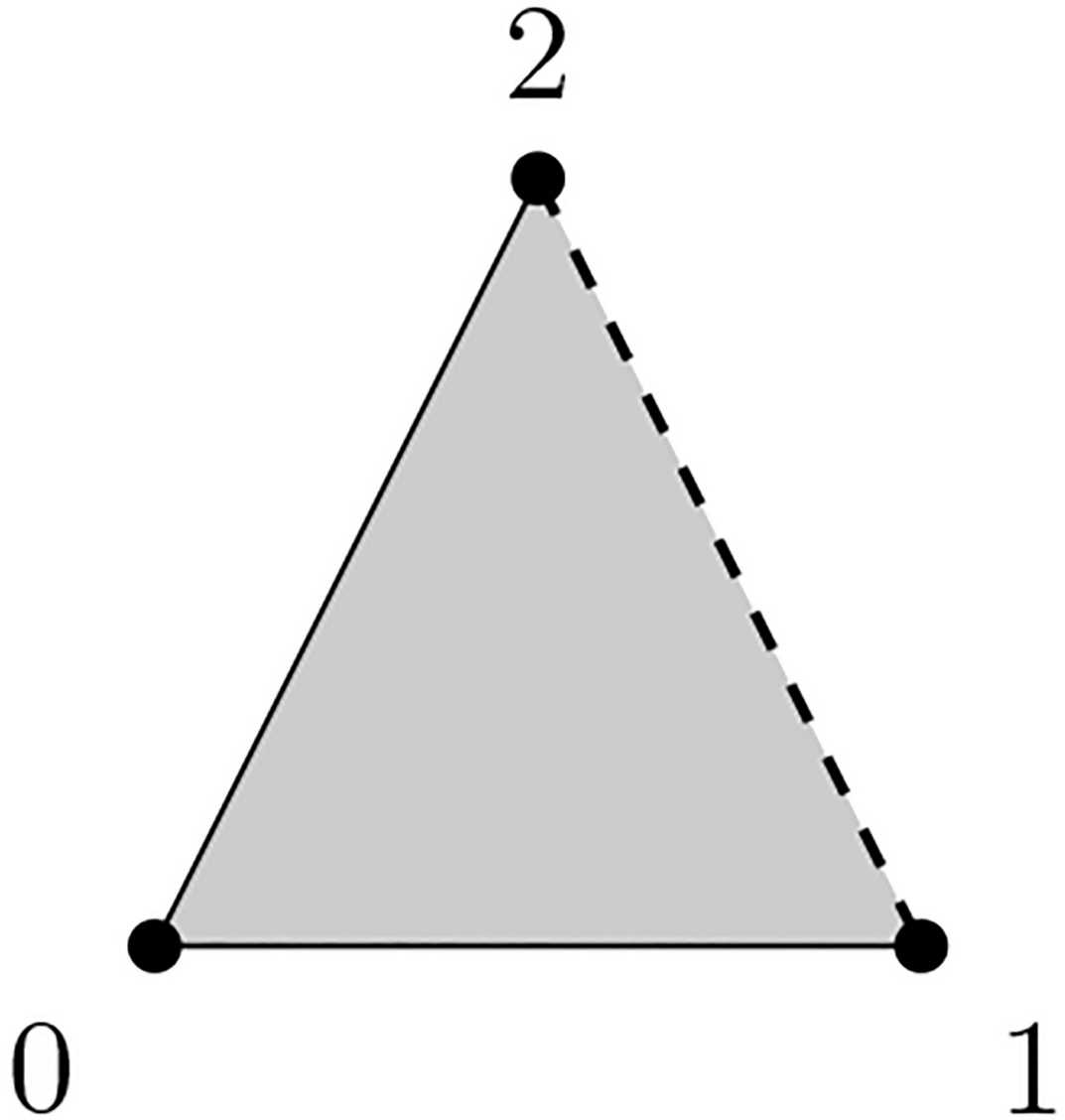
**Figure 2.**

Illustration of a super-hypergraph model constructed from the protein-ligand complex (ID: 3E6Y). The ligand (green color) is a drug that is used to cure the disease caused by the protein (red color). The potency and efficacy of the drug is directly determined by the atomic interactions between the ligand and the protein. Traditionally, atomic interactions are modeled by a graph (A). However, graphs can only characterize pair-wise interactions (by edges) and fall short for many-body interactions. Hypergraph models (C) use the hyperedge, i.e., a set of vertices, to represent many-body interactions and have demonstrated great power for biomolecular data analysis (See Section 4 for details). Mathematically, a  $n$ -hyperedge contains  $n + 1$  vertices in it. Note that 1-hyperedges are denoted by red ellipses and  $n$ -hyperedges ( $n > 1$ ) are represented by blue ellipses. The super-hypergraph (D) provides an even more flexible representation and incorporates detailed local topology within each hyperedge. Note that the hyperedge in super-hypergraph is a subgraph, i.e., a set of vertices together with edges. If we only consider vertex part of the subgraph, the super-hypergraph reduces to a hypergraph.

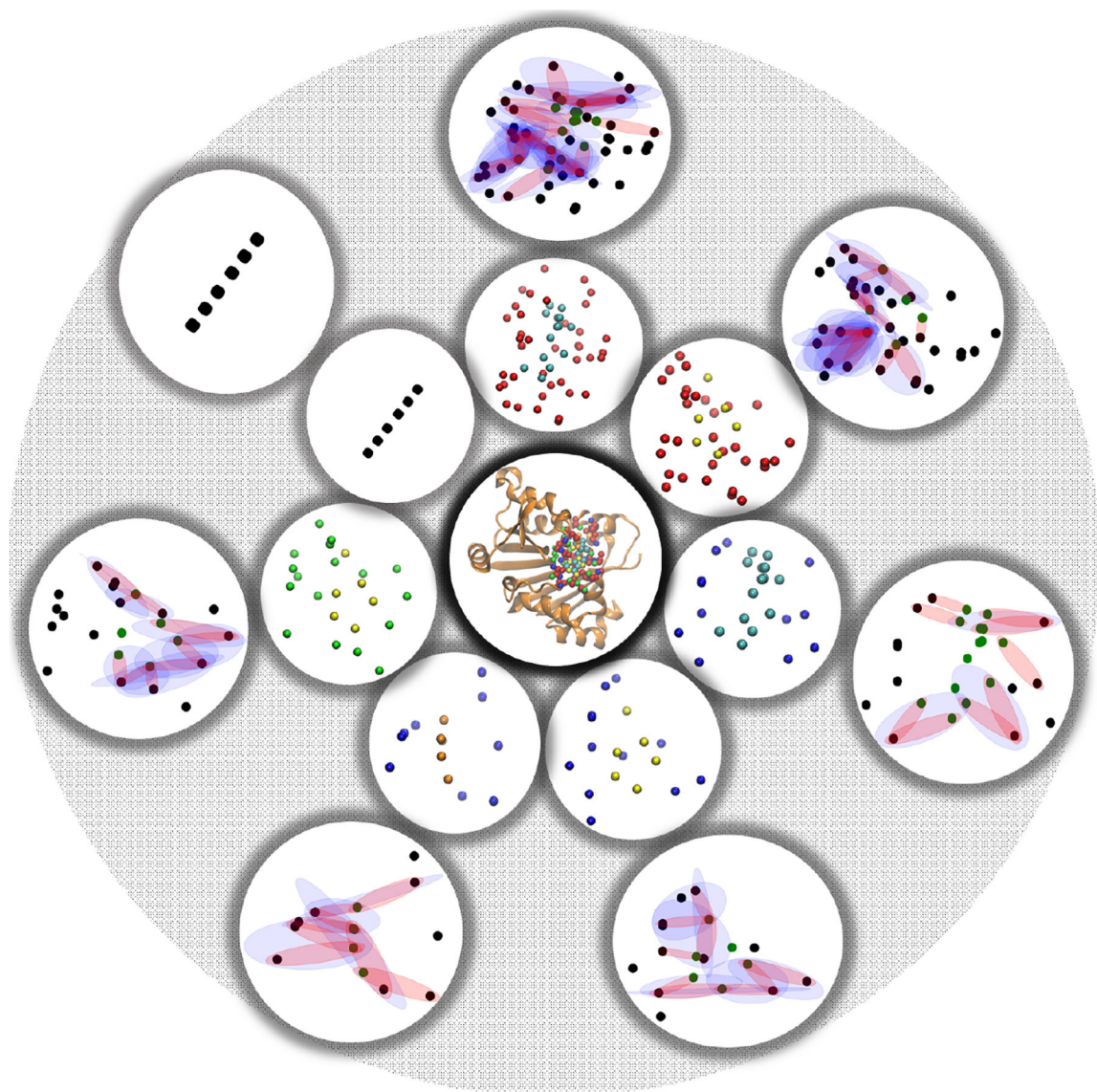




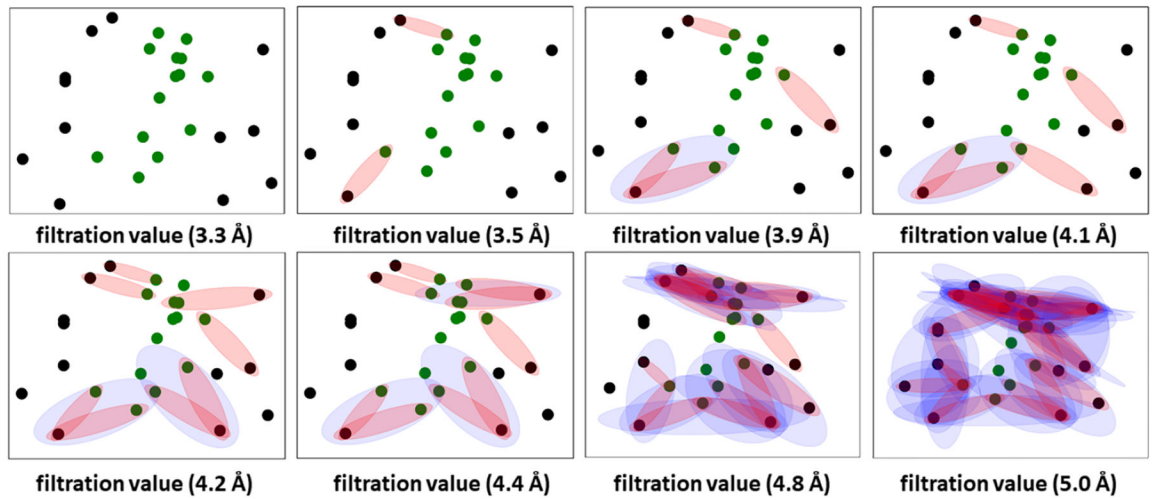
**Figure 3.** (a) The hypergraph  $\mathcal{H}$ , where the cross indicates that a vertex is missing. (b)  $\Delta(\mathcal{H})$ , the smallest  $\Delta$ -set that contains  $\mathcal{H}$ .



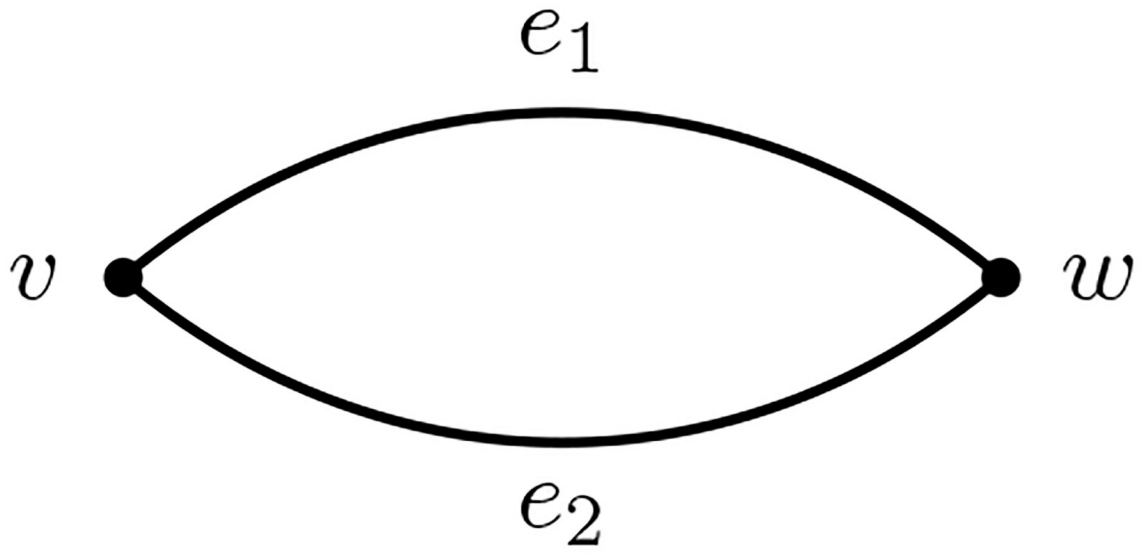
**Figure 4.**  
The hypergraph  $\mathcal{H}$  is a standard 2 simplex where the dotted edge is missing.



**Figure 5.** Illustration of an element-specific hypergraph model for a protein-ligand complex (ID 3PB3). The binding core region of the complex is decomposed into a series of element-specific atom-sets. The interactions between protein atom-sets and ligand atom-sets are modeled as a series of hypergraphs.



**Figure 6.**  
Illustration of a hypergraph-based filtration process for the protein-ligand complex with ID 3PB3.



**Figure 7.**  
The multi-graph  $G$ , which looks the same as the clique complex  $\text{Clique}(G)$ .



**Figure 8.** (a) The graph  $G$ , which is the same as  $\text{Clique}(G)$ . (b) The neighborhood complex of  $G$ ,  $\mathcal{N}(G)$ .

**Table 1.**

Topological structures associated to graphs

Constructions	Complex Type	Face Type
clique complex of a simple graph	simplicial complex	vertex-deletion
clique complex of a multi-graph	-set	vertex-deletion
neighborhood complex	simplicial complex	vertex-deletion
Jonsson's graph complex	simplicial complex	edge-deletion
path complex of a simple graph	hypergraph	vertex-deletion
path complex of a multi-graph/quiver	super-hypergraph	vertex-deletion

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript