
CAT-WALK: Inductive Hypergraph Learning via Set Walks

Ali Behrouz

Department of Computer Science
University of British Columbia
alibez@cs.ubc.ca

Farnoosh Hashemi[†]

Department of Computer Science
University of British Columbia
farsh@cs.ubc.ca

Sadaf Sadeghian[†]

Department of Computer Science
University of British Columbia
sadafsdn@cs.ubc.ca

Margo Seltzer

Department of Computer Science
University of British Columbia
mseltzer@cs.ubc.ca

Abstract

Temporal hypergraphs provide a powerful paradigm for modeling time-dependent, higher-order interactions in complex systems. Representation learning for hypergraphs is essential for extracting patterns of the higher-order interactions that are critically important in real-world problems in social network analysis, neuroscience, finance, etc. However, existing methods are typically designed only for specific tasks or static hypergraphs. We present CAT-WALK, an inductive method that learns the underlying dynamic laws that govern the temporal and structural processes underlying a temporal hypergraph. CAT-WALK introduces a temporal, higher-order walk on hypergraphs, SETWALK, that extracts higher-order causal patterns. CAT-WALK uses a novel adaptive and permutation invariant pooling strategy, SETMIXER, along with a set-based anonymization process that hides the identity of hyperedges. Finally, we present a simple yet effective neural network model to encode hyperedges. Our evaluation on 10 hypergraph benchmark datasets shows that CAT-WALK attains outstanding performance on temporal hyperedge prediction benchmarks in both inductive and transductive settings. It also shows competitive performance with state-of-the-art methods for node classification. (Code)

1 Introduction

Temporal networks have become increasingly popular for modeling interactions among entities in dynamic systems [1–5]. While most existing work focuses on pairwise interactions between entities, many real-world complex systems exhibit natural relationships among multiple entities [6–8]. Hypergraphs provide a natural extension to graphs by allowing an edge to connect any number of vertices, making them capable of representing higher-order structures in data. Representation learning on (temporal) hypergraphs has been recognized as an important machine learning problem and has become the cornerstone behind a wealth of high-impact applications in computer vision [9, 10], biology [11, 12], social networks [13, 14], and neuroscience [15, 16].

Many recent attempts to design representation learning methods for hypergraphs are equivalent to applying Graph Neural Networks (GNNs) to the clique-expansion (CE) of a hypergraph [17–21]. CE is a straightforward way to generalize graph algorithms to hypergraphs by replacing hyperedges with (weighted) cliques [18–20]. However, this decomposition of hyperedges limits expressiveness,

[†]These two authors contributed equally (ordered alphabetically) and reserve the right to swap their order.

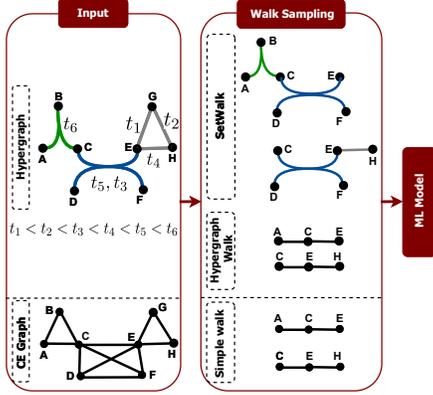


Figure 1: The advantage of SETWALKS in walk-based hypergraph learning.

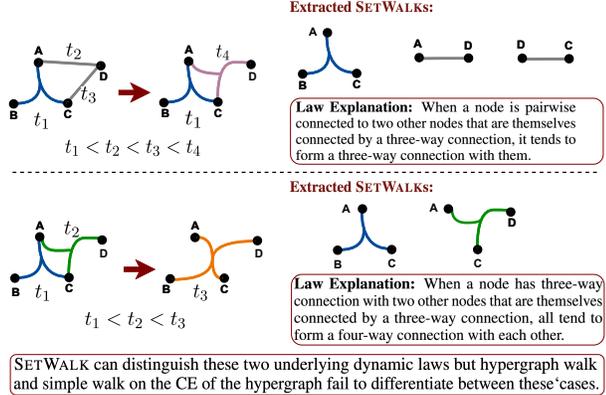


Figure 2: The advantage of SETWALKS in causality extraction and capturing complex dynamic laws.

leading to suboptimal performance [6, 22–24] (see [Theorem 1](#) and [Theorem 4](#)). New methods that encode hypergraphs directly partially address this issue [11, 25–28]. However, these methods suffer from some combination of the following three limitations: they are designed for ① learning the structural properties of *static hypergraphs* and do not consider temporal properties, ② the transductive setting, limiting their performance on unseen patterns and data, and ③ a specific downstream task (e.g., node classification [25], hyperedge prediction [26], or subgraph classification [27]) and cannot easily be extended to other downstream tasks, limiting their application.

Temporal motif-aware and neighborhood-aware methods have been developed to capture complex patterns in data [29–31]. However, counting temporal motifs in large networks is time-consuming and non-parallelizable, limiting the scalability of these methods. To this end, several recent studies suggest using temporal random walks to automatically retrieve such motifs [32–36]. One possible solution to capturing underlying temporal and higher-order structure is to extend the concept of a hypergraph random walk [37–43] to its temporal counterpart by letting the walker walk over time. However, existing definitions of random walks on hypergraphs offer limited expressivity and sometimes degenerate to simple walks on the CE of the hypergraph [40] (see [Appendix C](#)). There are two reasons for this: ① Random walks are composed of a sequence of *pair-wise* interconnected vertices, even though edges in a hypergraph connect *sets* of vertices. Decomposing them into sequences of simple pair-wise interactions loses the semantic meaning of the hyperedges (see [Theorem 4](#)). ② A sampling probability of a walk on a hypergraph must be different from its sampling probability on the CE of the hypergraph [37–43]. However, Chitra and Raphael [40] shows that each definition of the random walk with edge-independent sampling probability of nodes is equivalent to random walks on a weighted CE of the hypergraph. Existing studies on random walks on hypergraphs ignore ① and focus on ② to distinguish the walks on simple graphs and hypergraphs. However, as we show in [Table 2](#), ① is equally important, if not more so.

For example, [Figure 1](#) shows the procedure of existing walk-based machine learning methods on a temporal hypergraph. The neural networks in the model take as input only sampled walks. However, the output of the hypergraph walk [37, 38] and simple walk on the CE graph are the same. This means that the neural network cannot distinguish between pair-wise and higher-order interactions.

We present Causal Anonymous Set Walks (CAT-WALK), an inductive hyperedge learning method. We introduce a hyperedge-centric random walk on hypergraphs, called SETWALK, that automatically extracts temporal, higher-order motifs. The hyperedge-centric approach enables SETWALKS to distinguish multi-way connections from their corresponding CEs (see [Figure 1](#), [Figure 2](#), and [Theorem 1](#)). We use temporal hypergraph motifs that reflect network dynamics ([Figure 2](#)) to enable CAT-WALK to work well in the inductive setting. To make the model agnostic to the hyperedge identities of these motifs, we use two-step, set-based anonymization: ① Hide node identities by assigning them new positional encodings based on the number of times that they appear at a specific position in a set of sampled SETWALKS, and ② Hide hyperedge identities by combining the positional encodings of the nodes comprising the hyperedge using a novel permutation invariant pooling strategy, called SETMIXER. We incorporate a neural encoding method that samples a few SETWALKS starting from nodes of interest. It encodes and aggregates them via MLP-MIXER [44] and our new pooling strategy

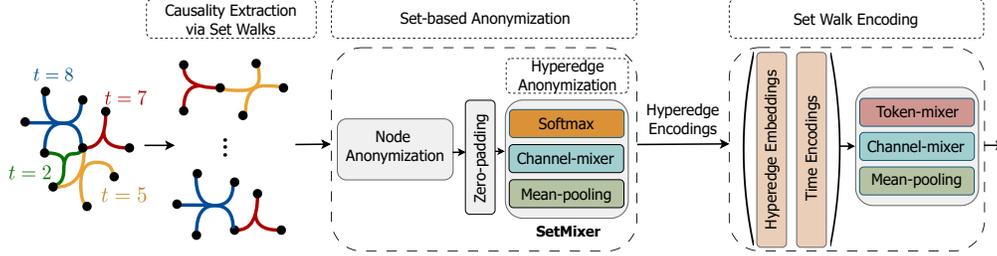


Figure 3: **Schematic of the CAT-WALK.** CAT-WALK consists of three stages: (1) Causality Extraction via Set Walks (§3.2), (2) Set-based Anonymization (§3.3), and (3) Set Walk Encoding (§3.4).

SETMIXER, respectively, to predict temporal, higher-order interactions. Finally, we discuss how to extend CAT-WALK for node classification. Figure 3 shows the schematic of the CAT-WALK.

We theoretically and experimentally discuss the effectiveness of CAT-WALK and each of its components. More specifically, we prove that SETWALKS are more expressive than existing random walk algorithms on hypergraphs. We demonstrate SETMIXER’s efficacy as a permutation invariant pooling strategy for hypergraphs and prove that using it in our anonymization process makes that process more expressive than existing anonymization processes [33, 45, 46] when applied to the CE of the hypergraphs. To the best of our knowledge, we report the most extensive experiments in the hypergraph learning literature pertaining to unsupervised hyperedge prediction with 10 datasets and eight baselines. Results show that our method produces 9% and 17% average improvement in transductive and inductive settings, outperforming all state-of-the-art baselines in the hyperedge prediction task. Also, CAT-WALK achieves the best or on-par performance on dynamic node classification tasks. All proofs appear in the Appendix.

2 Related Work

Temporal graph learning is an active research area [5, 47]. A major group of methods uses GNNs to learn node encodings and Recurrent Neural Networks (RNNs) to update these encodings over time [48–55]. More sophisticated methods based on anonymous temporal random walks [33, 34], line graphs [56], GraphMixer [57], neighborhood representation [58], and subgraph sketching [59] are designed to capture complex structures in vertex neighborhoods. Although these methods show promising results in a variety of tasks, they are fundamentally limited in that they are designed for *pair-wise* interaction among vertices and not the higher-order interactions in hypergraphs.

Representation learning on hypergraphs addresses this problem [17, 60]. We group work in this area into three overlapping categories:

① **Clique and Star Expansion:** CE-based methods replace hyperedges with (weighted) cliques and apply GNNs, sometimes with sophisticated propagation rules [21, 25], degree normalization, and nonlinear hyperedge weights [17–21, 39, 61, 62]. Although these methods are simple, it is well-known that CE causes undesired losses in learning performance, specifically when relationships within an incomplete subset of a hyperedge do not exist [6, 22–24]. Star expansion (SE) methods first use hypergraph star expansion and model the hypergraph as a bipartite graph, where one set of vertices represents nodes and the other represents hyperedges [25, 28, 63, 64]. Next, they apply modified heterogeneous GNNs, possibly with dual attention mechanisms from nodes to hyperedges and vice versa [25, 27]. Although this group does not cause as large a distortion as CE, they are neither memory nor computationally efficient. ② **Message Passing:** Most existing hypergraph learning methods, use message passing over hypergraphs [17–21, 25–27, 39, 61, 62, 65, 66]. Recently, Chien et al. [25] and Huang and Yang [28] designed universal message-passing frameworks that include propagation rules of most previous methods (e.g., [17, 19]). The main drawback of these two frameworks is that they are limited to node classification tasks and do not easily generalize to other tasks. ③ **Walk-based:** random walks are a common approach to extracting graph information for machine learning algorithms [32–34, 67, 68]. Several walk-based hypergraph learning methods are designed for a wide array of applications [43, 69–78]. However, most existing methods use simple random walks on the CE of the hypergraph (e.g., [26, 43, 78]). More complicated random walks on hypergraphs address this limitation [40–42, 79]. Although some of these studies show that their walk’s transition matrix

differs from the simple walk on the CE [40, 79], their extracted walks can still be the same, limiting their expressivity (see Figure 1, Figure 2, and Theorem 1).

Our method differs from all prior (temporal) (hyper)graph learning methods in five ways: CAT-WALK: ① Captures temporal higher-order properties in a streaming manner: In contrast to existing methods in hyperedge prediction, our method captures temporal properties in a streaming manner, avoiding the drawbacks of snapshot-based methods. ② Works in the inductive setting by extracting underlying dynamic laws of the hypergraph, making it generalizable to unseen patterns and nodes. ③ Introduces a hyperedge-centric, temporal, higher-order walk with a new perspective on the walk sampling procedure. ④ Presents a new two-step anonymization process: Anonymization of higher-order patterns (i.e., SETWALKS), requires hiding the identity of both nodes and hyperedges. We present a new permutation-invariant pooling strategy to hide hyperedges’ identity according to their vertices, making the process provably more expressive than the existing anonymization processes [33, 34, 78]. ⑤ Removes self-attention and RNNs from the walk encoding procedure, avoiding their limitations. Appendices A and B provide a more comprehensive discussion of related work.

3 Method: CAT-WALK Network

3.1 Preliminaries

Definition 1 (Temporal Hypergraphs). *A temporal hypergraph $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{X})$, can be represented as a sequence of hyperedges that arrive over time, i.e., $\mathcal{E} = \{(e_1, t_1), (e_2, t_2), \dots\}$, where \mathcal{V} is the set of nodes, $e_i \in 2^{\mathcal{V}}$ are hyperedges, t_i is the timestamp showing when e_i arrives, and $X \in \mathbb{R}^{|\mathcal{V}| \times f}$ is a matrix that encodes node attribute information for nodes in \mathcal{V} . Note that we treat each hyperedge e_i as the set of all vertices connected by e_i .*

Example 1. *The input of Figure 1 shows an example of a temporal hypergraph. Each hyperedge is a set of nodes that are connected by a higher-order connection. Each higher-order connection is specified by a color (e.g., green, blue, and gray), and also is associated with a timestamp (e.g., t_i).*

Given a hypergraph $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{X})$, we represent the set of hyperedges attached to a node u before time t as $\mathcal{E}^t(u) = \{(e, t') | t' < t, u \in e\}$. We say two hyperedges e and e' are adjacent if $e \cap e' \neq \emptyset$ and use $\mathcal{E}^t(e) = \{(e', t') | t' < t, e' \cap e \neq \emptyset\}$ to represent the set of hyperedges adjacent to e before time t . Next, we focus on the hyperedge prediction task: Given a subset of vertices $\mathcal{U} = \{u_1, \dots, u_k\}$, we want to predict whether a hyperedge among all the u_i s will appear in the next timestamp or not. In Appendix G.2 we discuss how this approach can be extended to node classification.

The (hyper)graph isomorphism problem is a decision problem that decides whether a pair of (hyper)graphs are isomorphic or not. Based on this concept, next, we discuss the measure we use to compare the expressive power of different methods.

Definition 2 (Expressive Power Measure). *Given two methods \mathcal{M}_1 and \mathcal{M}_2 , we say method \mathcal{M}_1 is more expressive than method \mathcal{M}_2 if:*

1. *for any pair of (hyper)graphs $(\mathcal{G}_1, \mathcal{G}_2)$ such that $\mathcal{G}_1 \neq \mathcal{G}_2$, if method \mathcal{M}_1 can distinguish \mathcal{G}_1 and \mathcal{G}_2 then method \mathcal{M}_2 can also distinguish them,*
2. *there is a pair of hypergraphs $\mathcal{G}'_1 \neq \mathcal{G}'_2$ such that \mathcal{M}_1 can distinguish them but \mathcal{M}_2 cannot.*

3.2 Causality Extraction via SETWALK

The collaboration network in Figure 1 shows how prior work that models hypergraph walks as sequences of vertices fails to capture complex connections in hypergraphs. Consider the two walks: $A \rightarrow C \rightarrow E$ and $H \rightarrow E \rightarrow C$. These two walks can be obtained either from a hypergraph random walk or from a simple random walk on the CE graph. Due to the symmetry of these walks with respect to (A, C, E) and (H, E, C) , they cannot distinguish A and H , although the neighborhoods of these two nodes exhibit different patterns: A, B , and C have published a paper together (connected by a hyperedge), but each pair of E, G , and H has published a paper (connected by a pairwise link). We address this limitation by defining a temporal walk on hypergraphs as a sequence of hyperedges:

Definition 3 (SETWALK). *Given a temporal hypergraph $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{X})$, a SETWALK with length ℓ on temporal hypergraph \mathcal{G} is a randomly generated sequence of hyperedges (sets):*

$$Sw : (e_1, t_{e_1}) \rightarrow (e_2, t_{e_2}) \rightarrow \dots \rightarrow (e_\ell, t_{e_\ell}),$$

where $e_i \in \mathcal{E}$, $t_{e_{i+1}} < t_{e_i}$, and the intersection of e_i and e_{i+1} is not empty, $e_i \cap e_{i+1} \neq \emptyset$. In other words, for each $1 \leq i \leq \ell - 1$: $e_{i+1} \in \mathcal{E}^h(e_i)$. We use $\text{Sw}[i]$ to denote the i -th hyperedge-time pair in the SETWALK. That is, $\text{Sw}[i][0] = e_i$ and $\text{Sw}[i][1] = t_{e_i}$.

Example 2. Figure 1 illustrates a temporal hypergraph with two sampled SETWALKS, hypergraph random walks, and simple walks. As an example, $(\{A, B, C\}, t_6) \rightarrow (\{C, D, E, F\}, t_5)$ is a SETWALK that starts from hyperedge e including $\{A, B, C\}$ in time t_6 , backtracks overtime and then samples hyperedge e' including $\{C, D, E, F\}$ in time $t_5 < t_6$. While this higher-order random walk with length two provides information about all $\{A, B, C, D, E, F\}$, its simple hypergraph walk counterpart, i.e. $A \rightarrow C \rightarrow E$, provides information about only three nodes.

Next, we formally discuss the power of SETWALKS. The proof of the theorem can be found in Appendix E.1.

Theorem 1. A random SETWALK is equivalent to neither the hypergraph random walk, the random walk on the CE graph, nor the random walk on the SE graph. Also, for a finite number of samples of each, SETWALK is more expressive than existing walks.

In Figure 1, SETWALKS capture higher-order interactions and distinguish the two nodes A and H , which are indistinguishable via hypergraph random walks and graph random walks in the CE graph. We present a more detailed discussion and comparison with previous definitions of random walks on hypergraphs in Appendix C.

Causality Extraction. We introduce a sampling method to allow SETWALKS to extract temporal higher-order motifs that capture causal relationships by backtracking over time and sampling adjacent hyperedges. As discussed in previous studies [33, 34], more recent connections are usually more informative than older connections. Inspired by Wang et al. [33], we use a hyperparameter $\alpha \geq 0$ to sample a hyperedge e with probability proportional to $\exp(\alpha(t - t_p))$, where t and t_p are the timestamps of e and the previously sampled hyperedge in the SETWALK, respectively. Additionally, we want to bias sampling towards pairs of adjacent hyperedges that have a greater number of common nodes to capture higher-order motifs. However, as discussed in previous studies, the importance of each node for each hyperedge can be different [25, 27, 40, 65]. Accordingly, the transferring probability from hyperedge e_i to its adjacent hyperedge e_j depends on the importance of the nodes that they share. We address this via a temporal SETWALK sampling process with *hyperedge-dependent node weights*. Given a temporal hypergraph $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{X})$, a hyperedge-dependent node-weight function $\Gamma : \mathcal{V} \times \mathcal{E} \rightarrow \mathbb{R}^{\geq 0}$, and a previously sampled hyperedge (e_p, t_p) , we sample a hyperedge (e, t) with probability:

$$\mathbb{P}[(e, t)|(e_p, t_p)] \propto \frac{\exp(\alpha(t - t_p))}{\sum_{(e', t') \in \mathcal{E}^p(e_p)} \exp(\alpha(t' - t_p))} \times \frac{\exp(\varphi(e, e_p))}{\sum_{(e', t') \in \mathcal{E}^p(e_p)} \exp(\varphi(e', e_p))}, \quad (1)$$

where $\varphi(e, e') = \sum_{u \in e \cap e'} \Gamma(u, e) \Gamma(u, e')$, representing the assigned weight to $e \cap e'$. We refer to the first and second terms as *temporal bias* and *structural bias*, respectively.

The pseudocode of our SETWALK sampling algorithm and its complexity analysis are in Appendix D. We also discuss this hyperedge-dependent sampling procedure and how it is provably more expressive than existing hypergraph random walks in Appendix C.

Given a (potential) hyperedge $e_0 = \{u_1, u_2, \dots, u_k\}$ and a time t_0 , we say a SETWALK, Sw, starts from u_i if $u_i \in \text{Sw}[1][0]$. We use the above procedure to generate M SETWALKS with length m starting from each $u_i \in e_0$. We use $\mathcal{S}(u_i)$ to store SETWALKS that start from u_i .

3.3 Set-based Anonymization of Hyperedges

In the anonymization process, we replace hyperedge identities with position encodings, capturing structural information while maintaining the inductiveness of the method. Micali and Zhu [45] studied Anonymous Walks (AWs), which replace a *node's identity* by the order of its appearance in each walk. The main limitation of AWs is that the position encoding of each node depends only on its specific walk, missing the dependency and correlation of different sampled walks [33]. To mitigate this drawback, Wang et al. [33] suggest replacing node identities with the hitting counts of the nodes based on a set of sampled walks. In addition to the fact that this method is designed for walks on simple graphs, there are two main challenges to adopting it for SETWALKS: ① SETWALKS

are a sequence of hyperedges, so we need an encoding for the position of hyperedges. Natural attempts to replace hyperedges’ identity with the hitting counts of the hyperedges based on a set of sampled SETWALKS, misses the similarity of hyperedges with many of the same nodes. ② Each hyperedge is a set of vertices and natural attempts to encode its nodes’ positions and aggregate them to obtain a position encoding of the hyperedge requires a permutation invariant pooling strategy. This pooling strategy also requires consideration of the higher-order dependencies between obtained nodes’ position encodings to take advantage of higher-order interactions (see [Theorem 2](#)). To address these challenges we present a set-based anonymization process for SETWALKS. Given a hyperedge $e_0 = \{u_1, \dots, u_k\}$, let w_0 be any node in e_0 . For each node w that appears on at least one SETWALK in $\bigcup_{i=1}^k \mathcal{S}(u_i)$, we assign a relative, node-dependent node identity, $\mathcal{R}(w, \mathcal{S}(w_0)) \in \mathbb{Z}^m$, as follows:

$$\mathcal{R}(w, \mathcal{S}(w_0))[i] = |\{\text{Sw} | \text{Sw} \in \mathcal{S}(w_0), w \in \text{Sw}[i][0]\}| \quad \forall i \in \{1, 2, \dots, m\}. \quad (2)$$

For each node w we further define $\text{Id}(w, e_0) = \{\mathcal{R}(w, \mathcal{S}(u_i))\}_{i=1}^k$. Let $\Psi(\cdot) : \mathbb{R}^{d_1} \rightarrow \mathbb{R}^{d_2}$ be a pooling function that gets a set of d_1 -dimensional vectors and aggregates them to a d_2 -dimensional vector. Given two instances of this pooling function, $\Psi_1(\cdot)$ and $\Psi_2(\cdot)$, for each hyperedge $e = \{w_1, w_2, \dots, w_k\}$ that appears on at least one SETWALK in $\bigcup_{i=1}^k \mathcal{S}(u_i)$, we assign a relative hyperedge identity as:

$$\text{Id}(e, e_0) = \Psi_1 \left(\{\Psi_2(\text{Id}(w_i, e_0))\}_{i=1}^k \right). \quad (3)$$

That is, for each node $w_i \in e$ we first aggregate its relative node-dependent identities (i.e., $\mathcal{R}(w_i, \mathcal{S}(u_j))$) to obtain the relative hyperedge-dependent identity. Then we aggregate these hyperedge-dependent identities for all nodes in e . Since the size of hyperedges can vary, we zero-pad to a fixed length. Note that this zero padding is important to capture the size of the hyperedge. The hyperedge with more zero-padded dimensions has fewer nodes.

This process addresses the first challenge and encodes the position of hyperedges. Unfortunately, many simple and known pooling strategies (e.g., $\text{SUM}(\cdot)$, $\text{ATTN-SUM}(\cdot)$, $\text{MEAN}(\cdot)$, etc.) can cause missing information when applied to hypergraphs. We formalize this in the following theorem:

Theorem 2. *Given an arbitrary positive integer $k \in \mathbb{Z}^+$, let $\Psi(\cdot)$ be a pooling function such that for any set $S = \{w_1, \dots, w_d\}$:*

$$\Psi(S) = \sum_{\substack{S' \subseteq S \\ |S'|=k}} f(S'), \quad (4)$$

where f is some function. Then the pooling function can cause missing information, meaning that it limits the expressiveness of the method to applying to the projected graph of the hypergraph.

While simple concatenation does not suffer from this undesirable property, it is not permutation invariant. To overcome these challenges, we design an all-MLP permutation invariant pooling function, SETMIXER, that not only captures higher-order dependencies of set elements but also captures dependencies across the number of times that a node appears at a certain position in SETWALKS.

SETMIXER. MLP-MIXER [44] is a family of models based on multi-layer perceptrons, widely used in the computer vision community, that are simple, amenable to efficient implementation, and robust to over-squashing and long-term dependencies (unlike RNNs and attention mechanisms) [44, 57]. However, the token-mixer phase of these methods is sensitive to the order of the input (see [Appendix A](#)). To address this limitation, inspired by MLP-MIXER [44], we design SETMIXER as follows: Let $S = \{\mathbf{v}_1, \dots, \mathbf{v}_d\}$, where $\mathbf{v}_i \in \mathbb{R}^{d_1}$, be the input set and $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_d] \in \mathbb{R}^{d_1 \times d}$ be its matrix representation:

$$\Psi(S) = \text{MEAN} \left(\mathbf{H}_{\text{token}} + \sigma \left(\text{LayerNorm} \left(\mathbf{H}_{\text{token}} \mathbf{W}_s^{(1)} \right) \mathbf{W}_s^{(2)} \right) \right), \quad (5)$$

where

$$\mathbf{H}_{\text{token}} = \mathbf{V} + \sigma \left(\text{Softmax} \left(\text{LayerNorm} \left(\mathbf{V}^T \right) \right) \right)^T. \quad (6)$$

Here, $\mathbf{W}_s^{(1)}$ and $\mathbf{W}_s^{(2)}$ are learnable parameters, $\sigma(\cdot)$ is an activation function (we use GeLU [80] in our experiments), and LayerNorm is layer normalization [81]. [Equation 5](#) is the channel mixer and [Equation 6](#) is the token mixer. The main intuition of SETMIXER is to use the Softmax(.) function to bind token-wise information in a non-parametric manner, avoiding permutation variant operations in the token mixer. We formally prove the following theorem in [Appendix E.3](#).

Theorem 3. *SETMIXER is permutation invariant and is a universal approximator of invariant multi-set functions. That is, SETMIXER can approximate any invariant multi-set function.*

Based on the above theorem, SETMIXER can overcome the challenges we discussed earlier as it is permutation invariant. Also, it is a universal approximator of multi-set functions, which shows its power to learn any arbitrary function. Accordingly, in our anonymization process, we use $\Psi(\cdot) = \text{SETMIXER}(\cdot)$ in Equation 3 to hide hyperedge identities. Next, we guarantee that our anonymization process does not depend on hyperedges or nodes identities, which justifies the claim of inductiveness of our model:

Proposition 1. *Given two (potential) hyperedges $e_0 = \{u_1, \dots, u_k\}$ and $e'_0 = \{u'_1, \dots, u'_k\}$, if there exists a bijective mapping π between node identities such that for each SETWALK like $\text{Sw} \in \bigcup_{i=1}^k \mathcal{S}(u_i)$ can be mapped to one SETWALK like $\text{Sw}' \in \bigcup_{i=1}^k \mathcal{S}(u'_i)$, then for each hyperedge $e = \{w_1, \dots, w_{k'}\}$ that appears in at least one SETWALK in $\bigcup_{i=1}^k \mathcal{S}(u_i)$, we have $\text{Id}(e, e_0) = \text{Id}(\pi(e), e'_0)$, where $\pi(e) = \{\pi(w_1), \dots, \pi(w_{k'})\}$.*

Finally, we guarantee that our anonymization approach is more expressive than existing anonymization process [33, 45] when applied to the CE of the hypergraphs:

Theorem 4. *The set-based anonymization method is more expressive than any existing anonymization strategies on the CE of the hypergraph. More precisely, there exists a pair of hypergraphs $\mathcal{G}_1 = (\mathcal{V}_1, \mathcal{E}_1)$ and $\mathcal{G}_2 = (\mathcal{V}_2, \mathcal{E}_2)$ with different structures (i.e., $\mathcal{G}_1 \not\cong \mathcal{G}_2$) that are distinguishable by our anonymization process and are not distinguishable by the CE-based methods.*

3.4 SETWALK Encoding

Previous walk-based methods [33, 34, 78] view a walk as a sequence of nodes. Accordingly, they plug nodes' positional encodings in a RNN [82] or Transformer [83] to obtain the encoding of each walk. However, in addition to the computational cost of RNN and Transformers, they suffer from over-squashing and fail to capture long-term dependencies. To this end, we design a simple and low-cost SETWALK encoding procedure that uses two steps: ① A time encoding module to distinguish different timestamps, and ② A mixer module to summarize temporal and structural information extracted by SETWALKS.

Time Encoding. We follow previous studies [33, 84] and adopt random Fourier features [85, 86] to encode time. However, these features are periodic, so they capture only periodicity in the data. We add a learnable linear term to the feature representation of the time encoding. We encode a given time t as follows:

$$\mathbf{T}(t) = (\boldsymbol{\omega}_l t + \mathbf{b}_l) \parallel \cos(t\boldsymbol{\omega}_p), \quad (7)$$

where $\boldsymbol{\omega}_l, \mathbf{b}_l \in \mathbb{R}$ and $\boldsymbol{\omega}_p \in \mathbb{R}^{d_2-1}$ are learnable parameters and \parallel shows concatenation.

MIXER Module. To summarize the information in each SETWALK, we use a MLP-MIXER [44] on the sequence of hyperedges in a SETWALK as well as their corresponding encoded timestamps. Contrary to the anonymization process, where we need a permutation invariant procedure, here, we need a permutation variant procedure since the order of hyperedges in a SETWALK is important. Given a (potential) hyperedge $e_0 = \{u_1, \dots, u_k\}$, we first assign $\text{Id}(e, e_0)$ to each hyperedge e that appears on at least one sampled SETWALK starting from e_0 (Equation 3). Given a SETWALK, $\text{Sw} : (e_1, t_{e_1}) \rightarrow \dots \rightarrow (e_m, t_{e_m})$, we let \mathbf{E} be a matrix that $\mathbf{E}_i = \text{Id}(e_i, e_0) \parallel \mathbf{T}(t_{e_i})$:

$$\text{ENC}(\text{Sw}) = \text{MEAN} \left(\mathbf{H}_{\text{token}} + \sigma \left(\text{LayerNorm}(\mathbf{H}_{\text{token}}) \mathbf{W}_{\text{channel}}^{(1)} \right) \mathbf{W}_{\text{channel}}^{(2)} \right), \quad (8)$$

where

$$\mathbf{H}_{\text{token}} = \mathbf{E} + \mathbf{W}_{\text{token}}^{(2)} \sigma \left(\mathbf{W}_{\text{token}}^{(1)} \text{LayerNorm}(\mathbf{E}) \right). \quad (9)$$

3.5 Training

In the training phase, for each hyperedge in the training set, we adopt the commonly used negative sample generation method [60] to generate a negative sample. Next, for each hyperedge in the training set such as $e_0 = \{u_1, u_2, \dots, u_k\}$, including both positive and negative samples, we sample M SETWALKS with length m starting from each $u_i \in e_0$ to construct $\mathcal{S}(u_i)$. Next, we anonymize each hyperedge that appears in at least one SETWALK in $\bigcup_{i=1}^k \mathcal{S}(u_i)$ by Equation 3 and then use the MIXER module to encode each $\text{Sw} \in \bigcup_{i=1}^k \mathcal{S}(u_i)$. To encode each node $u_i \in e_0$, we use $\text{MEAN}(\cdot)$ pooling over SETWALKS in $\mathcal{S}(u_i)$. Finally, to encode e_0 we use SETMIXER to mix obtained node encodings. For

hyperedge prediction, we use a 2-layer perceptron over the hyperedge encodings to make the final prediction. We discuss node classification in [Appendix G.2](#).

4 Experiments

We evaluate the performance of our model on two important tasks: hyperedge prediction and node classification (see [Appendix G.2](#)) in both inductive and transductive settings. We then discuss the importance of our model design and the significance of each component in CAT-WALK.

4.1 Experimental Setup

Baselines. We compare our method to eight previous state-of-the-art baselines on the hyperedge prediction task. These methods can be grouped into three categories: ① Deep hypergraph learning methods including HYPER-SAGCN [26], NHP [87], and CHESHIRE [11]. ② Shallow methods including HPRA [88] and HPLSF [89]. ③ CE methods: CE-CAW [33], CE-EVOLVEGCN [90] and CE-GCN [91]. Details on these models and hyperparameters used appear in [Appendix F.2](#).

Datasets. We use 10 available benchmark datasets [6] from the existing hypergraph neural networks literature. These datasets’ domains include drug networks (i.e., NDC [6]), contact networks (i.e., High School [92] and Primary School [93]), the US. Congress bills network [94, 95], email networks (i.e., Email Enron [6] and Email Eu [96]), and online social networks (i.e., Question Tags and Users-Threads [6]). Detailed descriptions of these datasets appear in [Appendix F.1](#).

Evaluation Tasks. We focus on Hyperedge Prediction: In the transductive setting, we train on the temporal hyperedges with timestamps less than or equal to T_{train} and test on those with timestamps greater than T_{train} . Inspired by Wang et al. [33], we consider two inductive settings. In the **Strongly Inductive** setting, we predict hyperedges consisting of some unseen nodes. In the **Weakly Inductive** setting, we predict hyperedges with *at least* one seen and some unseen nodes. We first follow the procedure used in the transductive setting, and then we randomly select 10% of the nodes and remove all hyperedges that include them from the training set. We then remove all hyperedges with seen nodes from the validation and testing sets. For dynamic node classification, see [Appendix G.2](#). For all datasets, we fix $T_{\text{train}} = 0.7 T$, where T is the last timestamp. To evaluate the models’ performance we follow the literature and use Area Under the ROC curve (AUC) and Average Precision (AP).

4.2 Results and Discussion

Hyperedge Prediction. We report the results of CAT-WALK and baselines in [Table 1](#) and [Appendix G](#). The results show that CAT-WALK achieves the best overall performance compared to the baselines in both transductive and inductive settings. In the transductive setting, not only does our method outperform baselines in all but one dataset, but it achieves near perfect results (i.e., ≈ 98.0) on the NDC and Primary School datasets. In the Weakly Inductive setting, our model achieves high scores (i.e., > 91.5) in all but one dataset, while most baselines perform poorly as they are not designed for the inductive setting and do not generalize well to unseen nodes or patterns. In the Strongly Inductive setting, CAT-WALK still achieves high AUC (i.e., > 90.0) on most datasets and outperforms baselines on *all* datasets. There are three main reasons for CAT-WALK’s superior performance: ① Our SETWALKS capture higher-order patterns. ② CAT-WALK incorporates temporal properties (both from SETWALKS and our time encoding module), thus learning underlying dynamic laws of the network. The other temporal methods (CE-CAW and CE-EVOLVEGCN) are CE-based methods, limiting their ability to capture higher-order patterns. ③ CAT-WALK’s set-based anonymization process that avoids using node and hyperedge identities allows it to generalize to unseen patterns and nodes.

Ablation Studies. We next conduct ablation studies on the High School, Primary School, and Users-Threads datasets to validate the effectiveness of CAT-WALK’s critical components. [Table 2](#) shows AUC results for inductive hyperedge prediction. The first row reports the performance of the complete CAT-WALK implementation. Each subsequent row shows results for CAT-WALK with one module modification: row 2 replace SETWALK by edge-dependent hypergraph walk [40], row 3 removes the time encoding module, row 4 replaces SETMIXER with MEAN(.) pooling, row 5 replaces the SETMIXER with sum-based universal approximator for sets [97], row 6 replaces the MLP-MIXER

Table 1: Performance on hyperedge prediction: Mean AUC (%) \pm standard deviation. Boldfaced letters shaded blue indicate the best result, while gray shaded boxes indicate results within one standard deviation of the best result. N/A: the method has numerical precision or computational issues.

Methods	NDC Class	High School	Primary School	Congress Bill	Email Enron	Email Eu	Question Tags	Users-Threads
Strongly Inductive								
CE-GCN	52.31 \pm 2.99	60.54 \pm 2.06	52.34 \pm 2.75	49.18 \pm 3.61	63.04 \pm 1.80	52.76 \pm 2.41	56.10 \pm 1.88	57.91 \pm 1.56
CE-EvolveGCN	49.78 \pm 3.13	46.12 \pm 3.83	58.01 \pm 2.56	54.00 \pm 1.84	57.31 \pm 4.19	44.16 \pm 1.27	64.08 \pm 2.75	52.00 \pm 2.32
CE-CAW	76.45 \pm 0.29	83.73 \pm 1.42	80.31 \pm 1.46	75.38 \pm 1.25	70.81 \pm 1.13	72.99 \pm 0.20	70.14 \pm 1.89	73.12 \pm 1.06
NHP	70.53 \pm 4.95	65.29 \pm 3.80	70.86 \pm 3.42	69.82 \pm 2.19	49.71 \pm 6.09	65.35 \pm 2.07	68.23 \pm 3.34	71.83 \pm 2.64
HYPER-SAGCN	79.05 \pm 2.48	88.12 \pm 3.01	80.13 \pm 1.38	79.51 \pm 1.27	73.09 \pm 2.60	78.01 \pm 1.24	73.66 \pm 1.95	73.94 \pm 2.57
CHESHIRE	72.24 \pm 2.63	82.54 \pm 0.88	77.26 \pm 1.01	79.43 \pm 1.58	70.03 \pm 2.55	69.98 \pm 2.71	N/A	76.99 \pm 2.82
CAT-WALK	98.89 \pm 1.82	96.03 \pm 1.50	95.32 \pm 0.89	93.54 \pm 0.56	73.45 \pm 2.92	91.68 \pm 2.78	88.03 \pm 3.38	89.84 \pm 6.02
Weakly Inductive								
CE-GCN	51.80 \pm 3.29	50.33 \pm 3.40	52.19 \pm 2.54	52.38 \pm 2.75	50.81 \pm 2.87	49.60 \pm 3.96	55.13 \pm 2.76	57.06 \pm 3.16
CE-EvolveGCN	55.39 \pm 5.16	57.85 \pm 3.51	51.50 \pm 4.07	55.63 \pm 3.41	45.66 \pm 2.10	52.44 \pm 2.38	61.79 \pm 1.63	55.81 \pm 2.54
CE-CAW	77.61 \pm 1.05	83.77 \pm 1.41	82.98 \pm 1.06	79.51 \pm 0.94	80.54 \pm 1.02	73.54 \pm 1.19	77.29 \pm 0.86	80.79 \pm 0.82
NHP	75.17 \pm 2.02	67.25 \pm 5.19	71.92 \pm 1.83	69.58 \pm 4.07	60.38 \pm 4.45	67.19 \pm 4.33	70.46 \pm 3.52	76.44 \pm 1.90
HYPER-SAGCN	79.45 \pm 2.18	88.53 \pm 1.26	85.08 \pm 1.45	80.12 \pm 2.00	78.86 \pm 0.63	77.26 \pm 2.09	78.15 \pm 1.41	75.38 \pm 1.43
CHESHIRE	79.03 \pm 1.24	88.40 \pm 1.06	83.55 \pm 1.27	79.67 \pm 0.83	74.53 \pm 0.91	77.31 \pm 0.95	N/A	81.27 \pm 0.85
CAT-WALK	99.16 \pm 1.08	94.68 \pm 2.37	96.53 \pm 1.39	98.38 \pm 0.21	64.11 \pm 7.96	91.98 \pm 2.41	90.28 \pm 2.81	97.15 \pm 1.81
Transductive								
HPRA	70.83 \pm 0.01	94.91 \pm 0.00	89.86 \pm 0.06	79.48 \pm 0.03	78.62 \pm 0.00	72.51 \pm 0.00	83.18 \pm 0.00	70.49 \pm 0.02
HPLSF	76.19 \pm 0.82	92.14 \pm 0.29	88.57 \pm 1.09	79.31 \pm 0.52	75.73 \pm 0.05	75.27 \pm 0.31	83.45 \pm 0.93	74.38 \pm 1.11
CE-GCN	66.83 \pm 3.74	62.99 \pm 3.02	59.14 \pm 3.87	64.42 \pm 3.11	58.06 \pm 3.80	64.19 \pm 2.79	55.18 \pm 5.12	62.78 \pm 2.69
CE-EvolveGCN	67.08 \pm 3.51	65.19 \pm 2.26	63.15 \pm 1.32	69.30 \pm 2.27	69.98 \pm 5.38	64.36 \pm 4.17	72.56 \pm 1.72	68.55 \pm 2.26
CE-CAW	76.30 \pm 0.84	81.63 \pm 0.97	86.53 \pm 0.84	76.99 \pm 1.02	79.57 \pm 0.14	78.19 \pm 1.10	81.73 \pm 2.48	80.86 \pm 0.45
NHP	82.39 \pm 2.81	76.85 \pm 3.08	80.04 \pm 3.42	80.27 \pm 2.53	63.17 \pm 3.79	78.90 \pm 4.39	79.14 \pm 3.52	82.33 \pm 1.02
HYPER-SAGCN	80.76 \pm 2.64	94.98 \pm 1.30	90.77 \pm 2.05	82.84 \pm 1.61	83.59 \pm 0.98	79.61 \pm 2.35	84.07 \pm 2.50	79.62 \pm 2.04
CHESHIRE	84.91 \pm 1.05	95.11 \pm 0.94	91.62 \pm 1.18	86.81 \pm 1.24	82.27 \pm 0.86	86.38 \pm 1.23	N/A	82.75 \pm 1.99
CAT-WALK	98.72 \pm 1.38	95.30 \pm 0.43	97.91 \pm 3.30	88.15 \pm 1.46	80.45 \pm 5.30	96.74 \pm 1.28	91.63 \pm 1.41	93.51 \pm 1.27

Table 2: Ablation study on CAT-WALK. AUC scores on inductive hyperedge prediction.

Methods	High School	Primary School	Users in Threads	Congress bill	Question Tags U
1 CAT-WALK	96.03 \pm 1.50	95.32 \pm 0.89	89.84 \pm 6.02	93.54 \pm 0.56	97.59 \pm 2.21
2 Replace SETWALK by Random Walk	92.10 \pm 2.18	51.56 \pm 5.63	53.24 \pm 1.73	80.27 \pm 0.02	67.74 \pm 2.92
3 Remove Time Encoding	95.94 \pm 0.19	86.80 \pm 6.33	70.58 \pm 9.32	92.56 \pm 0.49	96.91 \pm 1.89
4 Replace SETMIXER by MEAN(.)	94.58 \pm 1.22	95.14 \pm 4.36	63.59 \pm 5.26	91.06 \pm 0.24	68.62 \pm 1.25
5 Replace SETMIXER by Sum-based	94.77 \pm 0.67	90.86 \pm 0.57	60.03 \pm 1.16	91.07 \pm 0.70	89.76 \pm 0.45
6 Universal Approximator for Sets					
7 Replace MLP-MIXER by RNN	92.85 \pm 1.53	50.29 \pm 4.07	58.11 \pm 1.60	54.90 \pm 0.50	65.18 \pm 1.99
8 Replace MLP-MIXER by Transformer	55.98 \pm 0.83	86.64 \pm 3.55	60.65 \pm 1.56	89.38 \pm 1.66	56.16 \pm 4.03
9 Fix $\alpha = 0$	74.06 \pm 14.9	58.3 \pm 18.62	74.41 \pm 10.69	93.31 \pm 0.13	62.41 \pm 4.34

module with a RNN (see Appendix G for more experiments on the significance of using MLP-MIXER in walk encoding), row 7 replaces the MLP-MIXER module with a Transformer [83], and row 8 replaces the hyperparameter α with uniform sampling of hyperedges over all time periods. These results show that each component is critical for achieving CAT-WALK’s superior performance. The greatest contribution comes from SETWALK, MLP-MIXER in walk encoding, α in temporal hyperedge sampling, and SETMIXER pooling, respectively.

Hyperparameter Sensitivity. We analyze the effect of hyperparameters used in CAT-WALK, including temporal bias coefficient α , SETWALK length m , and sampling number M . The mean AUC performance on all inductive test hyperedges is reported in Figure 4. As expected, the left figure shows that

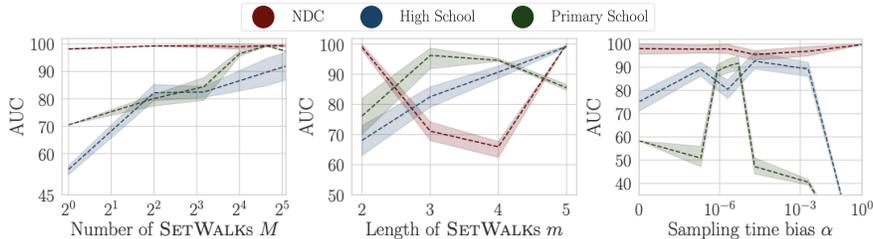


Figure 4: Hyperparameter sensitivity in CAT-WALK. AUC on inductive test hyperedges are reported.

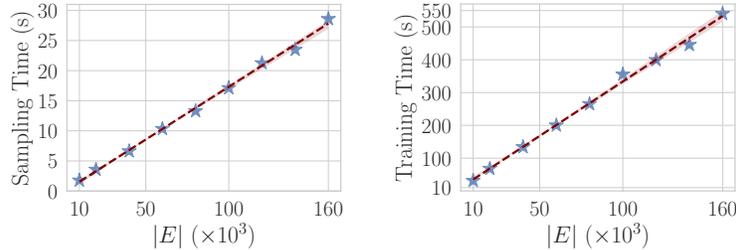


Figure 5: Scalability evaluation: The runtime of (left) SETWALK extraction and (right) the training time of CAT-WALK over one epoch on High School (using different $|E|$ for training).

increasing the number of sampled SETWALKS produces better performance. The main reason is that the model has more extracted structural and temporal information from the network. Also, notably, we observe that only a small number of sampled SETWALKS are needed to achieve competitive performance: in the best case 1 and in the worst case 16 sampled SETWALKS per each hyperedge are needed. The middle figure reports the effect of the length of SETWALKS on performance. These results show that performance peaks at certain SETWALK lengths and the exact value varies with the dataset. That is, longer SETWALKS are required for the networks that evolve according to more complicated laws encoded in temporal higher-order motifs. The right figure shows the effect of the temporal bias coefficient α . Results suggest that α has a dataset-dependent optimal interval. That is, a small α suggests an almost uniform sampling of interaction history, which results in poor performance when the short-term dependencies (interactions) are more important in the dataset. Also, very large α might damage performance as it makes the model focus on the most recent few interactions, missing long-term patterns.

Scalability Analysis. In this part, we investigate the scalability of CAT-WALK. To this end, we use different versions of the High School dataset with different numbers of hyperedges from 10^4 to 1.6×10^5 . Figure 5 (left) reports the runtimes of SETWALK sampling and Figure 5 (right) reports the runtimes of CAT-WALK for training one epoch using $M = 8$, $m = 3$ with batch-size = 64. Interestingly, our method scales linearly with the number of hyperedges, which enables it to be used on long hyperedge streams and large hypergraphs.

5 Conclusion, Limitation, and Future Work

We present CAT-WALK, an inductive hypergraph representation learning that learns both higher-order patterns and the underlying dynamic laws of temporal hypergraphs. CAT-WALK uses SETWALKS, a new temporal, higher-order random walk on hypergraphs that are provably more expressive than existing walks on hypergraphs, to extract temporal higher-order motifs from hypergraphs. CAT-WALK then uses a two-step, set-based anonymization process to establish the correlation between the extracted motifs. We further design a permutation invariant pooling strategy, SETMIXER, for aggregating nodes’ positional encodings in a hyperedge to obtain hyperedge level positional encodings. Consequently, the experimental results show that CAT-WALK (i) produces superior performance compared to the state-of-the-art in temporal hyperedge prediction tasks, and (ii) competitive performance in temporal node classification tasks. These results suggest many interesting directions for future studies: Using CAT-WALK as a positional encoder in existing anomaly detection frameworks to design an inductive anomaly detection method on hypergraphs. There are, however, a few limitations: Currently, CAT-WALK uses *fixed-length* SETWALKS, which might cause suboptimal performance. Developing a procedure to learn from SetWALKS of varying lengths might produce better results.

Acknowledgments and Disclosure of Funding

We acknowledge the support of the Natural Sciences and Engineering Research Council of Canada (NSERC).

Nous remercions le Conseil de recherches en sciences naturelles et en génie du Canada (CRSNG) de son soutien.

References

- [1] Serina Chang, Emma Pierson, Pang Wei Koh, Jaline Gerardin, Beth Redbird, David Grusky, and Jure Leskovec. Mobility network models of covid-19 explain inequities and inform reopening. *Nature*, 589(7840):82–87, 2021.
- [2] Michael Simpson, Farnoosh Hashemi, and Laks V. S. Lakshmanan. Misinformation mitigation under differential propagation rates and temporal penalties. *Proc. VLDB Endow.*, 15(10): 2216–2229, jun 2022. ISSN 2150-8097. doi: 10.14778/3547305.3547324. URL <https://doi.org/10.14778/3547305.3547324>.
- [3] Martino Ciaperoni, Edoardo Galimberti, Francesco Bonchi, Ciro Cattuto, Francesco Gullo, and Alain Barrat. Relevance of temporal cores for epidemic spread in temporal networks. *Scientific reports*, 10(1):1–15, 2020.
- [4] Farnoosh Hashemi, Ali Behrouz, and Laks V.S. Lakshmanan. Firmcore decomposition of multi-layer networks. In *Proceedings of the ACM Web Conference 2022, WWW '22*, page 1589–1600, New York, NY, USA, 2022. Association for Computing Machinery. ISBN 9781450390965. doi: 10.1145/3485447.3512205. URL <https://doi.org/10.1145/3485447.3512205>.
- [5] Antonio Longa, Veronica Lachi, Gabriele Santin, Monica Bianchini, Bruno Lepri, Pietro Lio, Franco Scarselli, and Andrea Passerini. Graph neural networks for temporal graphs: State of the art, open challenges, and opportunities. *arXiv preprint arXiv:2302.01018*, 2023.
- [6] Austin R. Benson, Rediet Abebe, Michael T. Schaub, Ali Jadbabaie, and Jon Kleinberg. Simplicial closure and higher-order link prediction. *Proceedings of the National Academy of Sciences*, 2018. ISSN 0027-8424. doi: 10.1073/pnas.1800683115.
- [7] Federico Battiston, Enrico Amico, Alain Barrat, Ginestra Bianconi, Guilherme Ferraz de Arruda, Benedetta Franceschiello, Iacopo Iacopini, Sonia Kéfi, Vito Latora, Yamir Moreno, Micah M. Murray, Tiago P. Peixoto, Francesco Vaccarino, and Giovanni Petri. The physics of higher-order interactions in complex systems. *Nature Physics*, 17(10):1093–1098, Oct 2021. ISSN 1745-2481. doi: 10.1038/s41567-021-01371-4. URL <https://doi.org/10.1038/s41567-021-01371-4>.
- [8] Yuanzhao Zhang, Maxime Lucas, and Federico Battiston. Higher-order interactions shape collective dynamics differently in hypergraphs and simplicial complexes. *Nature Communications*, 14(1):1605, Mar 2023. ISSN 2041-1723. doi: 10.1038/s41467-023-37190-9. URL <https://doi.org/10.1038/s41467-023-37190-9>.
- [9] Eun-Sol Kim, Woo Young Kang, Kyoung-Woon On, Yu-Jung Heo, and Byoung-Tak Zhang. Hypergraph attention networks for multimodal learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 14581–14590, 2020.
- [10] Yichao Yan, Jie Qin, Jiabin Chen, Li Liu, Fan Zhu, Ying Tai, and Ling Shao. Learning multi-granular hypergraphs for video-based person re-identification. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2899–2908, 2020.
- [11] Can Chen, Chen Liao, and Yang-Yu Liu. Teasing out missing reactions in genome-scale metabolic networks through hypergraph learning. *Nature Communications*, 14(1):2375, 2023.
- [12] Guobo Xie, Yinting Zhu, Zhiyi Lin, Yuping Sun, Guosheng Gu, Jianming Li, and Weiming Wang. Hbrwrlda: predicting potential lncrna–disease associations based on hypergraph bi-random walk with restart. *Molecular Genetics and Genomics*, 297(5):1215–1228, 2022.
- [13] Junliang Yu, Hongzhi Yin, Jundong Li, Qinyong Wang, Nguyen Quoc Viet Hung, and Xiangliang Zhang. Self-supervised multi-channel hypergraph convolutional network for social recommendation. In *Proceedings of the web conference 2021*, pages 413–424, 2021.
- [14] Dingqi Yang, Bingqing Qu, Jie Yang, and Philippe Cudre-Mauroux. Revisiting user mobility and social relationships in lbsns: a hypergraph embedding approach. In *The world wide web conference*, pages 2147–2157, 2019.

- [15] Junren Pan, Baiying Lei, Yanyan Shen, Yong Liu, Zhiguang Feng, and Shuqiang Wang. Characterization multimodal connectivity of brain network by hypergraph gan for alzheimer’s disease analysis. In *Pattern Recognition and Computer Vision: 4th Chinese Conference, PRCV 2021, Beijing, China, October 29–November 1, 2021, Proceedings, Part III 4*, pages 467–478. Springer, 2021.
- [16] Li Xiao, Junqi Wang, Peyman H Kassani, Yipu Zhang, Yuntong Bai, Julia M Stephen, Tony W Wilson, Vince D Calhoun, and Yu-Ping Wang. Multi-hypergraph learning-based brain functional connectivity analysis in fmri data. *IEEE transactions on medical imaging*, 39(5): 1746–1758, 2019.
- [17] Yifan Feng, Haoxuan You, Zizhao Zhang, Rongrong Ji, and Yue Gao. Hypergraph neural networks. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 3558–3565, 2019.
- [18] Dengyong Zhou, Jiayuan Huang, and Bernhard Schölkopf. Learning with hypergraphs: Clustering, classification, and embedding. In B. Schölkopf, J. Platt, and T. Hoffman, editors, *Advances in Neural Information Processing Systems*, volume 19. MIT Press, 2006. URL https://proceedings.neurips.cc/paper_files/paper/2006/file/dff8e9c2ac33381546d96deea9922999-Paper.pdf.
- [19] Naganand Yadati, Madhav Nimishakavi, Prateek Yadav, Vikram Nitin, Anand Louis, and Partha Talukdar. Hypergcn: A new method for training graph convolutional networks on hypergraphs. *Advances in neural information processing systems*, 32, 2019.
- [20] Sameer Agarwal, Jongwoo Lim, Lihi Zelnik-Manor, Pietro Perona, David Kriegman, and Serge Belongie. Beyond pairwise clustering. In *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR’05)*, volume 2, pages 838–845. IEEE, 2005.
- [21] Song Bai, Feihu Zhang, and Philip HS Torr. Hypergraph convolution and hypergraph attention. *Pattern Recognition*, 110:107637, 2021.
- [22] Matthias Hein, Simon Setzer, Leonardo Jost, and Syama Sundar Rangapuram. The total variation on hypergraphs - learning on hypergraphs revisited. In C.J. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 26. Curran Associates, Inc., 2013. URL https://proceedings.neurips.cc/paper_files/paper/2013/file/8a3363abe792db2d8761d6403605aeb7-Paper.pdf.
- [23] Pan Li and Olgica Milenkovic. Submodular hypergraphs: p-laplacians, Cheeger inequalities and spectral clustering. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 3014–3023. PMLR, 10–15 Jul 2018. URL <https://proceedings.mlr.press/v80/li18e.html>.
- [24] I (Eli) Chien, Huozhi Zhou, and Pan Li. hs^2 : Active learning over hypergraphs with pointwise and pairwise queries. In Kamalika Chaudhuri and Masashi Sugiyama, editors, *Proceedings of the Twenty-Second International Conference on Artificial Intelligence and Statistics*, volume 89 of *Proceedings of Machine Learning Research*, pages 2466–2475. PMLR, 16–18 Apr 2019. URL <https://proceedings.mlr.press/v89/chien19a.html>.
- [25] Eli Chien, Chao Pan, Jianhao Peng, and Olgica Milenkovic. You are allset: A multiset function framework for hypergraph neural networks. In *International Conference on Learning Representations*, 2022. URL https://openreview.net/forum?id=hpBTIv2uy_E.
- [26] Ruochi Zhang, Yuesong Zou, and Jian Ma. Hyper-sagmn: a self-attention based graph neural network for hypergraphs. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=ryeHuJBtPH>.
- [27] Yuan Luo. SHINE: Subhypergraph inductive neural network. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022. URL <https://openreview.net/forum?id=IsHRUzXPqhI>.

- [28] Jing Huang and Jie Yang. Unignn: a unified framework for graph and hypergraph neural networks. In Zhi-Hua Zhou, editor, *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence, IJCAI-21*, pages 2563–2569. International Joint Conferences on Artificial Intelligence Organization, 8 2021. doi: 10.24963/ijcai.2021/353. URL <https://doi.org/10.24963/ijcai.2021/353>. Main Track.
- [29] Ghadeer AbuOda, Gianmarco De Francisci Morales, and Ashraf Aboulnaga. Link prediction via higher-order motif features. In *Machine Learning and Knowledge Discovery in Databases: European Conference, ECML PKDD 2019, Würzburg, Germany, September 16–20, 2019, Proceedings, Part I*, pages 412–429. Springer, 2020.
- [30] Mayank Lahiri and Tanya Y Berger-Wolf. Structure prediction in temporal networks using frequent subgraphs. In *2007 IEEE Symposium on computational intelligence and data mining*, pages 35–42. IEEE, 2007.
- [31] Mahmudur Rahman and Mohammad Al Hasan. Link prediction in dynamic networks using graphlet. In *Machine Learning and Knowledge Discovery in Databases: European Conference, ECML PKDD 2016, Riva del Garda, Italy, September 19-23, 2016, Proceedings, Part I 16*, pages 394–409. Springer, 2016.
- [32] Giang H Nguyen, John Boaz Lee, Ryan A Rossi, Nesreen K Ahmed, Eunye Koh, and Sungchul Kim. Dynamic network embeddings: From random walks to temporal random walks. In *2018 IEEE International Conference on Big Data (Big Data)*, pages 1085–1092. IEEE, 2018.
- [33] Yanbang Wang, Yen-Yu Chang, Yunyu Liu, Jure Leskovec, and Pan Li. Inductive representation learning in temporal networks via causal anonymous walks. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=KYPz4YsCPj>.
- [34] Ming Jin, Yuan-Fang Li, and Shirui Pan. Neural temporal walks: Motif-aware representation learning on continuous-time dynamic graphs. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022. URL <https://openreview.net/forum?id=NqbktPUkZf7>.
- [35] Zhining Liu, Dawei Zhou, Yada Zhu, Jinjie Gu, and Jingrui He. Towards fine-grained temporal network representation via time-reinforced random walk. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 4973–4980, 2020.
- [36] Giang Hoang Nguyen, John Boaz Lee, Ryan A Rossi, Nesreen K Ahmed, Eunye Koh, and Sungchul Kim. Continuous-time dynamic network embeddings. In *Companion proceedings of the the web conference 2018*, pages 969–976, 2018.
- [37] Timoteo Carletti, Federico Battiston, Giulia Cencetti, and Duccio Fanelli. Random walks on hypergraphs. *Physical review E*, 101(2):022308, 2020.
- [38] Koby Hayashi, Sinan G Aksoy, Cheong Hee Park, and Haesun Park. Hypergraph random walks, laplacians, and clustering. In *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*, pages 495–504, 2020.
- [39] Josh Payne. Deep hyperedges: a framework for transductive and inductive learning on hypergraphs, 2019.
- [40] Uthsav Chitra and Benjamin Raphael. Random walks on hypergraphs with edge-dependent vertex weights. In *International conference on machine learning*, pages 1172–1181. PMLR, 2019.
- [41] Sinan G. Aksoy, Cliff Joslyn, Carlos Ortiz Marrero, Brenda Praggastis, and Emilie Purvine. Hypernetwork science via high-order hypergraph walks. *EPJ Data Science*, 9(1):16, Jun 2020. ISSN 2193-1127. doi: 10.1140/epjds/s13688-020-00231-0. URL <https://doi.org/10.1140/epjds/s13688-020-00231-0>.

- [42] Austin R. Benson, David F. Gleich, and Lek-Heng Lim. The spacey random walk: A stochastic process for higher-order data. *SIAM Review*, 59(2):321–345, 2017. doi: 10.1137/16M1074023. URL <https://doi.org/10.1137/16M1074023>.
- [43] Ankit Sharma, Shafiq Joty, Himanshu Kharkwal, and Jaideep Srivastava. Hyperedge2vec: Distributed representations for hyperedges, 2018. URL <https://openreview.net/forum?id=rJ5C67-C->.
- [44] Ilya Tolstikhin, Neil Houlsby, Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Thomas Unterthiner, Jessica Yung, Andreas Peter Steiner, Daniel Keysers, Jakob Uszkoreit, Mario Lucic, and Alexey Dosovitskiy. MLP-mixer: An all-MLP architecture for vision. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, 2021. URL <https://openreview.net/forum?id=EI2K0XKdnP>.
- [45] Silvio Micali and Zeyuan Allen Zhu. Reconstructing markov processes from independent and anonymous experiments. *Discrete Applied Mathematics*, 200:108–122, 2016. ISSN 0166-218X. doi: <https://doi.org/10.1016/j.dam.2015.06.035>. URL <https://www.sciencedirect.com/science/article/pii/S0166218X15003212>.
- [46] Ali Behrouz and Margo Seltzer. ADMIRE++: Explainable anomaly detection in the human brain via inductive learning on temporal multiplex networks. In *ICML 3rd Workshop on Interpretable Machine Learning in Healthcare (IMLH)*, 2023. URL <https://openreview.net/forum?id=t4H8acYudJ>.
- [47] Joakim Skarding, Bogdan Gabrys, and Katarzyna Musial. Foundations and modeling of dynamic networks using dynamic graph neural networks: A survey. *IEEE Access*, 9:79143–79168, 2021.
- [48] Youngjoo Seo, Michaël Defferrard, Pierre Vandergheynst, and Xavier Bresson. Structured sequence modeling with graph convolutional recurrent networks. In *International conference on neural information processing*, pages 362–373. Springer, 2018.
- [49] Ling Zhao, Yujiao Song, Chao Zhang, Yu Liu, Pu Wang, Tao Lin, Min Deng, and Haifeng Li. T-gcn: A temporal graph convolutional network for traffic prediction. *IEEE Transactions on Intelligent Transportation Systems*, 21(9):3848–3858, 2019.
- [50] Hao Peng, Hongfei Wang, Bowen Du, Md Zakirul Alam Bhuiyan, Hongyuan Ma, Jianwei Liu, Lihong Wang, Zeyu Yang, Linfeng Du, Senzhang Wang, et al. Spatial temporal incidence dynamic graph neural networks for traffic flow forecasting. *Information Sciences*, 521:277–290, 2020.
- [51] Xiaoyang Wang, Yao Ma, Yiqi Wang, Wei Jin, Xin Wang, Jiliang Tang, Caiyan Jia, and Jian Yu. Traffic flow prediction via spatial temporal graph neural network. In *Proceedings of The Web Conference 2020*, pages 1082–1092, 2020.
- [52] Jiaxuan You, Tianyu Du, and Jure Leskovec. Roland: Graph learning framework for dynamic graphs. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, KDD ’22, page 2358–2366, New York, NY, USA, 2022. Association for Computing Machinery. ISBN 9781450393850. doi: 10.1145/3534678.3539300. URL <https://doi.org/10.1145/3534678.3539300>.
- [53] Farnoosh Hashemi, Ali Behrouz, and Milad Rezaei Hajidehi. Cs-tgn: Community search via temporal graph neural networks. In *Companion Proceedings of the Web Conference 2023*, WWW ’23, New York, NY, USA, 2023. Association for Computing Machinery. doi: 10.1145/3543873.3587654. URL <https://doi.org/10.1145/3543873.3587654>.
- [54] Srijan Kumar, Xikun Zhang, and Jure Leskovec. Predicting dynamic embedding trajectory in temporal interaction networks. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, KDD ’19, page 1269–1278, New York, NY, USA, 2019. Association for Computing Machinery. ISBN 9781450362016. doi: 10.1145/3292500.3330895. URL <https://doi.org/10.1145/3292500.3330895>.

- [55] Ali Behrouz and Margo Seltzer. Anomaly detection in multiplex dynamic networks: from blockchain security to brain disease prediction. In *NeurIPS 2022 Temporal Graph Learning Workshop*, 2022. URL <https://openreview.net/forum?id=UDGZDfwmay>.
- [56] Sudhanshu Chanpuriya, Ryan A. Rossi, Sungchul Kim, Tong Yu, Jane Hoffswell, Nedim Lipka, Shunan Guo, and Cameron N Musco. Direct embedding of temporal network edges via time-decayed line graphs. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=Qamz7Q_Talk.
- [57] Weilin Cong, Si Zhang, Jian Kang, Baichuan Yuan, Hao Wu, Xin Zhou, Hanghang Tong, and Mehrdad Mahdavi. Do we really need complicated model architectures for temporal networks? In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=ayPPc0SyLv1>.
- [58] Yuhong Luo and Pan Li. Neighborhood-aware scalable temporal network representation learning. In *The First Learning on Graphs Conference*, 2022. URL <https://openreview.net/forum?id=EPUtNe7a9ta>.
- [59] Benjamin Paul Chamberlain, Sergey Shirobokov, Emanuele Rossi, Fabrizio Frasca, Thomas Markovich, Nils Yannick Hammerla, Michael M. Bronstein, and Max Hansmire. Graph neural networks for link prediction with subgraph sketching. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=m1oqEOAozQU>.
- [60] Can Chen and Yang-Yu Liu. A survey on hyperlink prediction. *arXiv preprint arXiv:2207.02911*, 2022.
- [61] Devanshu Arya, Deepak Gupta, Stevan Rudinac, and Marcel Worring. Hyper{sage}: Generalizing inductive representation learning on hypergraphs, 2021. URL <https://openreview.net/forum?id=cKnKJcTPRcV>.
- [62] Devanshu Arya, Deepak K Gupta, Stevan Rudinac, and Marcel Worring. Adaptive neural message passing for inductive learning on hypergraphs. *arXiv preprint arXiv:2109.10683*, 2021.
- [63] Sameer Agarwal, Kristin Branson, and Serge Belongie. Higher order learning with graphs. In *Proceedings of the 23rd international conference on Machine learning*, pages 17–24, 2006.
- [64] Chaoqi Yang, Ruijie Wang, Shuochao Yao, and Tarek Abdelzaher. Hypergraph learning with line expansion. *arXiv preprint arXiv:2005.04843*, 2020.
- [65] Kaize Ding, Jianling Wang, Jundong Li, Dingcheng Li, and Huan Liu. Be more with less: Hypergraph attention networks for inductive text classification. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 4927–4936, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.399. URL <https://aclanthology.org/2020.emnlp-main.399>.
- [66] Boxin Du, Changhe Yuan, Robert Barton, Tal Neiman, and Hanghang Tong. Hypergraph pre-training with graph neural networks. *arXiv preprint arXiv:2105.10862*, 2021.
- [67] Aditya Grover and Jure Leskovec. node2vec: Scalable feature learning for networks. In *Proceedings of the 22nd ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 855–864, 2016.
- [68] Bryan Perozzi, Rami Al-Rfou, and Steven Skiena. Deepwalk: Online learning of social representations. In *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 701–710, 2014.
- [69] Xin-Jian Xu, Chong Deng, and Li-Jie Zhang. Hyperlink prediction via local random walks and jensen–shannon divergence. *Journal of Statistical Mechanics: Theory and Experiment*, 2023(3):033402, mar 2023. doi: 10.1088/1742-5468/acc31e. URL <https://dx.doi.org/10.1088/1742-5468/acc31e>.

- [70] Valerio La Gatta, Vincenzo Moscato, Mirko Pennone, Marco Postiglione, and Giancarlo Sperli. Music recommendation via hypergraph embedding. *IEEE Transactions on Neural Networks and Learning Systems*, pages 1–13, 2022. doi: 10.1109/TNNLS.2022.3146968.
- [71] Dingqi Yang, Bingqing Qu, Jie Yang, and Philippe Cudre-Mauroux. Revisiting user mobility and social relationships in lbsns: A hypergraph embedding approach. In *The World Wide Web Conference, WWW '19*, page 2147–2157, New York, NY, USA, 2019. Association for Computing Machinery. ISBN 9781450366748. doi: 10.1145/3308558.3313635. URL <https://doi.org/10.1145/3308558.3313635>.
- [72] Aurélien Ducournau and Alain Bretto. Random walks in directed hypergraphs and application to semi-supervised image segmentation. *Computer Vision and Image Understanding*, 120: 91–102, 2014. ISSN 1077-3142. doi: <https://doi.org/10.1016/j.cviu.2013.10.012>. URL <https://www.sciencedirect.com/science/article/pii/S1077314213002038>.
- [73] Lei Ding and Alper Yilmaz. Interactive image segmentation using probabilistic hypergraphs. *Pattern Recognition*, 43(5):1863–1873, 2010. ISSN 0031-3203. doi: <https://doi.org/10.1016/j.patcog.2009.11.025>. URL <https://www.sciencedirect.com/science/article/pii/S0031320309004440>.
- [74] Sai Nageswar Satchidanand, Harini Ananthapadmanaban, and Balaraman Ravindran. Extended discriminative random walk: A hypergraph approach to multi-view multi-relational transductive learning. In *Proceedings of the 24th International Conference on Artificial Intelligence, IJCAI'15*, page 3791–3797. AAAI Press, 2015. ISBN 9781577357384.
- [75] Jie Huang, Chuan Chen, Fanghua Ye, Jiajing Wu, Zibin Zheng, and Guohui Ling. Hyper2vec: Biased random walk for hyper-network embedding. In Guoliang Li, Jun Yang, Joao Gama, Jugapong Natwichai, and Yongxin Tong, editors, *Database Systems for Advanced Applications*, pages 273–277, Cham, 2019. Springer International Publishing. ISBN 978-3-030-18590-9.
- [76] Jie Huang, Xin Liu, and Yangqiu Song. Hyper-path-based representation learning for hyper-networks. In *Proceedings of the 28th ACM International Conference on Information and Knowledge Management, CIKM '19*, page 449–458, New York, NY, USA, 2019. Association for Computing Machinery. ISBN 9781450369763. doi: 10.1145/3357384.3357871. URL <https://doi.org/10.1145/3357384.3357871>.
- [77] Jie Huang, Chuan Chen, Fanghua Ye, Weibo Hu, and Zibin Zheng. Nonuniform hyper-network embedding with dual mechanism. *ACM Trans. Inf. Syst.*, 38(3), may 2020. ISSN 1046-8188. doi: 10.1145/3388924. URL <https://doi.org/10.1145/3388924>.
- [78] Yunyu Liu, Jianzhu Ma, and Pan Li. Neural predicting higher-order patterns in temporal networks. In *Proceedings of the ACM Web Conference 2022, WWW '22*, page 1340–1351, New York, NY, USA, 2022. Association for Computing Machinery. ISBN 9781450390965. doi: 10.1145/3485447.3512181. URL <https://doi.org/10.1145/3485447.3512181>.
- [79] Jiying Zhang, Fuyang Li, Xi Xiao, Tingyang Xu, Yu Rong, Junzhou Huang, and Yatao Bian. Hypergraph convolutional networks via equivalency between hypergraphs and undirected graphs, 2022. URL <https://openreview.net/forum?id=zFyCvjXof60>.
- [80] Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus), 2020.
- [81] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.
- [82] David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. Learning representations by back-propagating errors. *nature*, 323(6088):533–536, 1986.
- [83] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [84] da Xu, chuanwei ruan, evren korpeoglu, sushant kumar, and kannan achan. Inductive representation learning on temporal graphs. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=rJeWlyHYwH>.

- [85] Da Xu, Chuanwei Ruan, Evren Korpeoglu, Sushant Kumar, and Kannan Achan. Self-attention with functional time representation learning. *Advances in neural information processing systems*, 32, 2019.
- [86] Seyed Mehran Kazemi, Rishab Goel, Sepehr Eghbali, Janahan Ramanan, Jaspreet Sahota, Sanjay Thakur, Stella Wu, Cathal Smyth, Pascal Poupart, and Marcus Brubaker. Time2vec: Learning a vector representation of time. *arXiv preprint arXiv:1907.05321*, 2019.
- [87] Naganand Yadati, Vikram Nitin, Madhav Nimishakavi, Prateek Yadav, Anand Louis, and Partha Talukdar. Nhp: Neural hypergraph link prediction. In *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*, pages 1705–1714, 2020.
- [88] Tarun Kumar, K Darwin, Srinivasan Parthasarathy, and Balaraman Ravindran. Hpra: Hyper-edge prediction using resource allocation. In *12th ACM conference on web science*, pages 135–143, 2020.
- [89] Ye Xu, Dan Rockmore, and Adam M Kleinbaum. Hyperlink prediction in hypernetworks using latent social features. In *Discovery Science: 16th International Conference, DS 2013, Singapore, October 6-9, 2013. Proceedings 16*, pages 324–339. Springer, 2013.
- [90] Aldo Pareja, Giacomo Domeniconi, Jie Chen, Tengfei Ma, Toyotaro Suzumura, Hiroki Kanezashi, Tim Kaler, Tao Schardl, and Charles Leiserson. Evolvegcn: Evolving graph convolutional networks for dynamic graphs. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 5363–5370, 2020.
- [91] Felix Wu, Amauri Souza, Tianyi Zhang, Christopher Fifty, Tao Yu, and Kilian Weinberger. Simplifying graph convolutional networks. In *International conference on machine learning*, pages 6861–6871. PMLR, 2019.
- [92] Rossana Mastrandrea, Julie Fournet, and Alain Barrat. Contact patterns in a high school: A comparison between data collected using wearable sensors, contact diaries and friendship surveys. *PLOS ONE*, 10(9):e0136497, 2015. doi: 10.1371/journal.pone.0136497. URL <https://doi.org/10.1371/journal.pone.0136497>.
- [93] Juliette Stehlé, Nicolas Voirin, Alain Barrat, Ciro Cattuto, Lorenzo Isella, Jean-François Pinton, Marco Quaggiotto, Wouter Van den Broeck, Corinne Régis, Bruno Lina, and Philippe Vanhems. High-resolution measurements of face-to-face contact patterns in a primary school. *PLoS ONE*, 6(8):e23176, 2011. doi: 10.1371/journal.pone.0023176. URL <https://doi.org/10.1371/journal.pone.0023176>.
- [94] James H. Fowler. Connecting the congress: A study of cosponsorship networks. *Political Analysis*, 14(04):456–487, 2006. doi: 10.1093/pan/mpi002. URL <https://doi.org/10.1093/pan/mpi002>.
- [95] James H. Fowler. Legislative cosponsorship networks in the US house and senate. *Social Networks*, 28(4):454–465, oct 2006. doi: 10.1016/j.socnet.2005.11.003. URL <https://doi.org/10.1016/j.socnet.2005.11.003>.
- [96] Hao Yin, Austin R. Benson, Jure Leskovec, and David F. Gleich. Local higher-order graph clustering. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM Press, 2017. doi: 10.1145/3097983.3098069. URL <https://doi.org/10.1145/3097983.3098069>.
- [97] Manzil Zaheer, Satwik Kottur, Siamak Ravanbakhsh, Barnabas Poczos, Russ R Salakhutdinov, and Alexander J Smola. Deep sets. *Advances in neural information processing systems*, 30, 2017.
- [98] Fan RK Chung. The laplacian of a hypergraph. In "", 1993.
- [99] Timoteo Carletti, Duccio Fanelli, and Renaud Lambiotte. Random walks and community detection in hypergraphs. *Journal of Physics: Complexity*, 2(1):015011, 2021.
- [100] Linyuan Lu and Xing Peng. High-order random walks and generalized laplacians on hypergraphs. *Internet Mathematics*, 9(1):3–32, 2013.

- [101] Chenzi Zhang, Shuguang Hu, Zhihao Gavin Tang, and TH Hubert Chan. Re-revisiting learning on hypergraphs: confidence interval and subgradient method. In *International Conference on Machine Learning*, pages 4026–4034. PMLR, 2017.
- [102] T-H Hubert Chan, Zhihao Gavin Tang, Xiaowei Wu, and Chenzi Zhang. Diffusion operator and spectral analysis for directed hypergraph laplacian. *Theoretical Computer Science*, 784: 46–64, 2019.
- [103] Pan Li and Olgica Milenkovic. Inhomogeneous hypergraph clustering with applications. *Advances in neural information processing systems*, 30, 2017.
- [104] Ziyu Wang, Wenhao Jiang, Yiming M Zhu, Li Yuan, Yibing Song, and Wei Liu. Dynamixer: a vision mlp architecture with dynamic mixing. In *International Conference on Machine Learning*, pages 22691–22701. PMLR, 2022.
- [105] Ali Behrouz, Parsa Delavari, and Farnoosh Hashemi. Unsupervised representation learning of brain activity via bridging voxel activity and functional connectivity. In *NeurIPS 2023 AI for Science Workshop*, 2023. URL <https://openreview.net/forum?id=HSvg7qFFd2>.
- [106] Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. *Advances in neural information processing systems*, 30, 2017.
- [107] Marta Garnelo, Dan Rosenbaum, Christopher Maddison, Tiago Ramalho, David Saxton, Murray Shanahan, Yee Whye Teh, Danilo Rezende, and SM Ali Eslami. Conditional neural processes. In *International conference on machine learning*, pages 1704–1713. PMLR, 2018.
- [108] David Lopez-Paz, Robert Nishihara, Soumith Chintala, Bernhard Scholkopf, and Léon Bottou. Discovering causal signals in images. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6979–6987, 2017.
- [109] Juho Lee, Yoonho Lee, Jungtaek Kim, Adam Kosiorek, Seungjin Choi, and Yee Whye Teh. Set transformer: A framework for attention-based permutation-invariant neural networks. In *International conference on machine learning*, pages 3744–3753. PMLR, 2019.
- [110] Jinheon Baek, Minki Kang, and Sung Ju Hwang. Accurate learning of graph representations with graph multiset pooling. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=JHcqXGaqiGn>.
- [111] Ronald Harry Atkin. *Mathematical structure in human affairs*. Heinemann Educational London, 1974.
- [112] David I Spivak. Higher-dimensional models of networks. *arXiv preprint arXiv:0909.4314*, 2009.
- [113] Jacob Charles Wright Billings, Mirko Hu, Giulia Lerda, Alexey N Medvedev, Francesco Mottes, Adrian Onicas, Andrea Santoro, and Giovanni Petri. Simplex2vec embeddings for community detection in simplicial complexes. *arXiv preprint arXiv:1906.09068*, 2019.
- [114] Mustafa Hajij, Kyle Istvan, and Ghada Zamzmi. Cell complex neural networks. In *TDA & Beyond*, 2020. URL <https://openreview.net/forum?id=6Tq18ySFpGU>.
- [115] Celia Hacker. \mathbb{S}^k -simplex2vec: a simplicial extension of node2vec. In *TDA & Beyond*, 2020. URL <https://openreview.net/forum?id=Aw9DUXPjq55>.
- [116] Stefania Ebli, Michaël Defferrard, and Gard Spreemann. Simplicial neural networks. In *TDA & Beyond*, 2020. URL <https://openreview.net/forum?id=nPct39DVIfk>.
- [117] Maosheng Yang, Elvin Isufi, and Geert Leus. Simplicial convolutional neural networks. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 8847–8851. IEEE, 2022.
- [118] Eric Bunch, Qian You, Glenn Fung, and Vikas Singh. Simplicial 2-complex convolutional neural networks. In *TDA & Beyond*, 2020. URL <https://openreview.net/forum?id=TLbnsKrt6J->.

- [119] Christopher Wei Jin Goh, Cristian Bodnar, and Pietro Lio. Simplicial attention networks. In *ICLR 2022 Workshop on Geometrical and Topological Representation Learning*, 2022. URL <https://openreview.net/forum?id=ScfRNWkpec>.
- [120] Mustafa Hajij, Karthikeyan Natesan Ramamurthy, Aldo Guzmán-Sáenz, and Ghada Za. High skip networks: A higher order generalization of skip connections. In *ICLR 2022 Workshop on Geometrical and Topological Representation Learning*, 2022.
- [121] Matthias Hein, Simon Setzer, Leonardo Jost, and Syama Sundar Rangapuram. The total variation on hypergraphs - learning on hypergraphs revisited. In C.J. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 26. Curran Associates, Inc., 2013. URL https://proceedings.neurips.cc/paper_files/paper/2013/file/8a3363abe792db2d8761d6403605aeb7-Paper.pdf.
- [122] Sameer Agarwal, Kristin Branson, and Serge Belongie. Higher order learning with graphs. In *Proceedings of the 23rd International Conference on Machine Learning, ICML '06*, page 17–24, New York, NY, USA, 2006. Association for Computing Machinery. ISBN 1595933832. doi: 10.1145/1143844.1143847. URL <https://doi.org/10.1145/1143844.1143847>.
- [123] Edmund Ihler, Dorothea Wagner, and Frank Wagner. Modeling hypergraphs by graphs with the same mincut properties. *Inf. Process. Lett.*, 45(4):171–175, mar 1993. ISSN 0020-0190. doi: 10.1016/0020-0190(93)90115-P. URL [https://doi.org/10.1016/0020-0190\(93\)90115-P](https://doi.org/10.1016/0020-0190(93)90115-P).
- [124] Peihao Wang, Shenghao Yang, Yunyu Liu, Zhangyang Wang, and Pan Li. Equivariant hypergraph diffusion neural operators. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=RiTjKoscnNd>.
- [125] Patrick Kidger, James Morrill, James Foster, and Terry Lyons. Neural controlled differential equations for irregular time series. *Advances in Neural Information Processing Systems*, 33: 6696–6707, 2020.
- [126] Giang Hoang Nguyen, John Boaz Lee, Ryan A. Rossi, Nesreen K. Ahmed, Eunye Koh, and Sungchul Kim. Continuous-time dynamic network embeddings. In *Companion Proceedings of the The Web Conference 2018, WWW '18*, page 969–976, Republic and Canton of Geneva, CHE, 2018. International World Wide Web Conferences Steering Committee. ISBN 9781450356404. doi: 10.1145/3184558.3191526. URL <https://doi.org/10.1145/3184558.3191526>.
- [127] Weilin Cong, Si Zhang, Jian Kang, Baichuan Yuan, Hao Wu, Xin Zhou, Hanghang Tong, and Mehrdad Mahdavi. Do we really need complicated model architectures for temporal networks? In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=ayPPc0SyLv1>.

Appendix

Table of Contents

A Preliminaries, Backgrounds, and Motivations	21
A.1 Anonymous Random Walks	21
A.2 Random Walk on Hypergraphs	21
A.3 MLP-Mixer	22
B Additional Related Work	22
B.1 Learning (Multi)Set Functions	22
B.2 Simplicial Complexes Representation Learning	23
B.3 How Does CAT-WALK Differ from Existing Works? (Contributions)	23
C SETWALK and Random Walk on Hypergraphs	23
C.1 Extension of SETWALKS	24
D Efficient Hyperedge Sampling	25
E Theoretical Results	26
E.1 Proof of Theorem 1	26
E.2 Proof of Theorem 2	28
E.3 Proof of Theorem 3	28
E.4 Proof of Theorem 4	29
E.5 Proof of Proposition 2	30
F Experimental Setup Details	31
F.1 Datasets	31
F.2 Baselines	32
F.3 Implementation and Training Details	33
G Additional Experimental Results	34
G.1 Results on More Datasets	34
G.2 Node Classification	34
G.3 Performance in Average Precision	34
G.4 More Results on RNN v.s. MLP-MIXER in Walk Encoding	35
H Broader Impacts	36

A Preliminaries, Backgrounds, and Motivations

We begin by reviewing the preliminaries and background concepts that we refer to in the main paper. Next, we discuss the fundamental differences between our method and techniques from prior work.

A.1 Anonymous Random Walks

Micali and Zhu [45] studied Anonymous Walks (AWs), which replace a *node’s identity* by the order of its appearance in each walk. Given a simple network, an AW starts from a node, performs random walks over the graph to collect a sequence of nodes, $W : (u_1, u_2, \dots, u_k)$, and then replaces the node identities by their order of appearance in each walk. That is:

$$ID_{AW}(w_0, W) = \{|u_1, u_2, \dots, u_{k^*}\} \text{ where } k^* \text{ is the smallest index such that } u_{k^*} = w_0. \quad (10)$$

While this method is a simple anonymization process, it misses the correlation between different walks and assigns new node identities based on only one single walk. The correlation between different walks is more important in temporal networks to assign new node identities, as a single walk cannot capture the frequency of a pattern over time [33]. To this end, Wang et al. [33] design a set-based anonymization process that assigns new node identities based on a set of sampled walks. Given a vertex u , they sample M walks with length m starting from u and store them in S_u . Next, for each node w_0 that appears on at least one walk in S_u , they assign a vector to each node as its hidden identity [33]:

$$g(w_0, S_u)[i] = |\{W|W \in S_u, W[i] = w_0\}| \quad \forall i \in \{0, \dots, m\}, \quad (11)$$

where $W[i]$ shows the i -th node in the walk W . This anonymization process not only hides the identity of vertices but it also can establish such hidden identity based on different sampled walks, capturing the correlation between several walks starting from a vertex.

Both of these anonymization processes are designed for graphs with pair-wise interactions, and there are three main challenges in adopting them for hypergraphs: ① To capture higher-order patterns, we use SETWALKS, which are a sequence of hyperedges. Accordingly, we need an encoding for the position of hyperedges. A natural attempt to encode the position of hyperedges is to count the position of hyperedges across sampled SETWALKS, as CAW [33] does for nodes. However, this approach misses the similarity of hyperedges with many nodes in common. That is, given two hyperedges $e_1 = \{u_1, u_2, \dots, u_k\}$ and $e_2 = \{u_1, u_2, \dots, u_k, u_{k+1}\}$. Although we want to encode the position of these two hyperedges, we also want these two hyperedges to have almost the same encoding as they share many vertices. Accordingly, we suggest viewing a hyperedge as a set of vertices. We first encode the position of vertices, and then we aggregate the position encodings of nodes that are connected by a hyperedge to compute the positional encoding of the hyperedge. ② However, since we focus on undirected hypergraphs, the order of a hyperedge’s vertices in the aggregation process should not affect the hyperedge positional encodings. Therefore, we need a permutation invariant pooling strategy. ③ While several existing studies used simple pooling functions such as MEAN(.) or SUM(.) [26], these pooling functions do not capture the higher-order dependencies between obtained nodes’ position encodings, missing the advantage of higher-order interactions. That is, a pooling function such as MEAN(.) is a non-parametric method that sees the positional encoding of each node in a hyperedge separately. Therefore, it is unable to aggregate them in a non-linear manner, which, depending on the data, can miss information. To address challenges ② and ③, we design SETMIXER, a permutation invariant pooling strategy that uses MLPs to learn how to aggregate positional encodings of vertices in a hyperedge to compute the hyperedge positional encoding.

A.2 Random Walk on Hypergraphs

Chung [98] presents some of the earliest research on the hypergraph Laplacian, defining the Laplacian of the k -uniform hypergraph. Following this line of work, Zhou et al. [18] defined a two-step CE-based random walk-based Laplacian for general hypergraphs. Given a node u , in the first step, we uniformly sample a hyperedge e including node u , and in the second step, we uniformly sample a node in e . Following this idea, several studies developed more sophisticated (weighted) CE-based random walks on hypergraphs [20]. However, Chitra and Raphael [40] shows that random walks on hypergraphs with edge-independent node weights are limited to capturing pair-wise interactions, making them unable to capture higher-order information. To address this limitation, they designed an edge-dependent sampling procedure of random walks on hypergraphs. Carletti et al. [37] and Carletti

et al. [99] argued that to sample more informative walks from a hypergraph, we must consider the degree of hyperedges in measuring the importance of vertices in the first step. Concurrently, some studies discuss the dependencies among hyperedges and define the s -th Laplacian based on simple walks on the dual hypergraphs [41, 100]. Finally, more sophisticated random walks with non-linear Laplacians have been designed [23, 101–103].

SETWALKS addresses three main drawbacks from existing methods: ① None of these methods are designed for temporal hypergraphs so they cannot capture temporal properties of the network. Also, natural attempts to extend them to temporal hypergraphs and let the walker uniformly walk over time ignore the fact that recent hyperedges are more informative than older ones (see Table 2). To address this issue, SETWALK uses a temporal bias factor in its sampling procedure (Equation 1). ② Existing hypergraph random walks are unable to capture either higher-order interactions of vertices or higher-order dependencies of hyperedges. That is, random walks with edge-independent weights [37] are not able to capture higher-order interactions and are equivalent to simple random walks on the CE of the hypergraph [40]. The expressivity of random walks on hypergraphs with edge-dependent walks is also limited when we have a limited number of sampled walks (see Theorem 1). Finally, defining a hypergraph random walk as a random walk on the dual hypergraph also cannot capture the higher-order dependencies of hyperedges (see Appendix C and Appendix D). SETWALK by its nature is able to walk over hyperedges (instead of vertices) and time and can capture higher-order interactions. Also, with a structural bias factor in its sampling procedure, which is based on hyperedge-dependent node weights, it is more informative than a simple random walk on the dual hypergraph, capturing higher-order dependencies of hyperedges. See Appendix C for further discussion.

A.3 MLP-Mixer

MLP-MIXER [44] is a family of models, based on multi-layer perceptions (MLPs), that are simple, amenable to efficient implementation, and robust to long-term dependencies (unlike RNNs, attention mechanisms, and Transformers [83]) with a wide array of applications from computer vision [104] to neuroscience [105]. The original architecture is designed for image data, where it takes image tokens as inputs. It then encodes them with a linear layer, which is equivalent to a convolutional layer over the image tokens, and updates their representations with a sequence of feed-forward layers applied to image tokens and features. Accordingly, we can divide the architecture of MLP-MIXER into two main parts: ① Token Mixer: The main intuition of the token mixer is to clearly separate the cross-location operations and learn the cross-feature (cross-location) dependencies. ② Channel Mixer: The intuition behind the channel mixer is to clearly separate the per-location operations and provide positional invariance, a prominent feature of convolutions. In both MIXER and SETMIXER we use the channel mixer as designed in MLP-MIXER. Next, we discuss the token mixer and its limitation in mixing features in a permutation variant manner:

Token Mixer. Let \mathbf{E} be the input of the MLP-MIXER, then the token mixer phase is defined as:

$$\mathbf{H}_{\text{token}} = \mathbf{E} + \mathbf{W}_{\text{token}}^{(2)} \sigma \left(\mathbf{W}_{\text{token}}^{(1)} \text{LayerNorm}(\mathbf{E})^T \right)^T, \quad (12)$$

where $\sigma(\cdot)$ is nonlinear activation function (usually GeLU [80]). Since it feeds the input’s columns to an MLP, it mixes the cross-feature information, which results in the MLP-MIXER being sensitive to permutation. Although natural attempts to remove the token mixer or its linear layer can produce a permutation invariant method, it misses cross-feature dependencies, which are the main motivation for using the MLP-MIXER architecture. To address this issue, SETMIXER uses the $\text{Softmax}(\cdot)$ function over features. Using Softmax over features can be seen as cross-feature normalization, which can capture their dependencies. While $\text{Softmax}(\cdot)$ is a non-parametric method that can bind token-wise information, it is also permutation equivariant, and as we prove in Appendix E.3, makes the SETMIXER permutation invariant.

B Additional Related Work

B.1 Learning (Multi)Set Functions

(Multi)set functions are pooling architectures for (multi)sets with a wide array of applications in many real-world problems including few-shot image classification [106], conditional regression [107], and causality discovery [108]. Zaheer et al. [97] develop DEEPSSETS, a universal approach to parameterize

the (multi)set functions. Following this direction, some works design attention mechanisms to learn multiset functions [109], which also inspired Baek et al. [110] to adopt attention mechanisms designed for (multi)set functions in graph representation learning. Finally, Chien et al. [25] build the connection between learning (multi)set functions with propagations on hypergraphs. To the best of our knowledge, SETMIXER is the first adaptive permutation invariant pooling strategy for hypergraphs, which views each hyperedge as a set of vertices and aggregates node encodings by considering their higher-order dependencies.

B.2 Simplicial Complexes Representation Learning

Simplicial complexes can be considered a special case of hypergraphs and are defined as a collection of polytopes such as triangles and tetrahedra, which are called simplices [111]. While these frameworks can be used to represent higher-order relations, simplicial complexes require the downward closure property [112]. That is, every substructure or face of a simplex contained in a complex \mathcal{K} is also in \mathcal{K} . Recently, to encode higher-order interactions, representation learning on simplicial complexes has attracted much attention [6, 113–119]. The first group of methods extend node2vec [67] to simplicial complexes with random walks on interactions through Hasse diagrams and simplex connections inside p -chains [113, 115]. With the recent advances in message-passing-based methods, several studies focus on designing neural networks on simplicial complexes [116–119]. Ebli et al. [116] introduced Simplicial neural networks (SNN), a generalization of spectral graph convolution to simplicial complexes with higher-order Laplacian matrices. Following this direction, some works propose simplicial convolutional neural networks with different simplicial filters to exploit the relationships in upper- and lower-neighborhoods [117, 118]. Finally, the last group of studies use an encoder-decoder architecture as well as message-passing to learn the representation of simplicial complexes [114, 120].

CAT-WALK is different from all these methods in three main aspects: ① Contrary to these methods, CAT-WALK is designed for temporal hypergraphs and is capable of capturing higher-order temporal properties in a streaming manner, avoiding the drawbacks of snapshot-based methods. ② CAT-WALK works in the inductive setting by extracting underlying dynamic laws of the hypergraph, making it generalizable to unseen patterns and nodes. ③ All these methods are designed for simplicial complexes, which are special cases of hypergraphs, while CAT-WALK is designed for general hypergraphs and does not require any assumption of the downward closure property.

B.3 How Does CAT-WALK Differ from Existing Works? (Contributions)

As we discussed in Appendix A.2, existing random walks on hypergraphs are unable to capture either ① higher-order interactions between nodes or ② higher-order dependencies of hyperedges. Moreover, all these walks are for static hypergraphs and are not able to capture temporal properties. To this end, we design SETWALK a higher-order temporal walk on hypergraphs. Naturally, SETWALKS are capable of capturing higher-order patterns as a SETWALK is defined as a sequence of hyperedges. We further design a new sampling procedure with temporal and structural biases, making SETWALKS capable of capturing higher-order dependencies of hyperedges. To take advantage of complex information provided by SETWALKS as well as training the model in an inductive manner, we design a two-step anonymization process with a novel pooling strategy, called SETMIXER. The anonymization process starts with encoding the position of vertices with respect to a set of sampled SETWALKS and then aggregates node positional encodings via a non-linear permutation invariant pooling function, SETMIXER, to compute their corresponding hyperedge positional encodings. This two-step process lets us capture structural properties while we also care about the similarity of hyperedges. Finally, to take advantage of continuous-time dynamics in data and avoid the limitations of sequential encoding, we design a neural network for temporal walk encoding that leverages a time encoding module to encode time as well as a MIXER module to encode the structure of the walk.

C SETWALK and Random Walk on Hypergraphs

We reviewed existing random walks in Appendix A.2. Here, we discuss how these concepts are different from SETWALKS and investigate whether SETWALKS are more expressive than these methods.

As we discussed in Sections 1 and 3.2, there are two main challenges for designing random walks on hypergraphs: ① Random walks are a sequence of *pair-wise* interconnected vertices, even though

edges in a hypergraph connect *sets* of vertices. ② A sampling probability of a walk on a hypergraph must be different from its sampling probability on the CE of the hypergraph [37–43]. To address these challenges, most existing works on random walks on hypergraphs ignore ① and focus on ② to distinguish the walks on simple graphs and hypergraphs, and ① is relatively unexplored. To this end, we answer the following questions:

Q1: Can ② alone be sufficient to take advantage of higher-order interactions? First, semantically, decomposing hyperedges into sequences of simple pair-wise interactions (CE) loses the semantic meaning of the hyperedges. Consider the collaboration network in Figure 1. When decomposing the hyperedges into pair-wise interactions, both (A, B, C) and (H, G, E) have the same structure (a triangle), while the semantics of these two structures in the data are completely different. That is, (A, B, C) have *all* published a paper together, while each pair of (H, G, E) separately have published a paper. One might argue that although the output of hypergraph random walks and simple random walks on the CE might be the same, the sampling probability of each walk is different and with a large number of samples, our model can distinguish these two structures. In Theorem 1 (proof in Appendix E) we theoretically show that when we have a finite number of hypergraph walk samples, M , there is a hypergraph \mathcal{G} such that with M hypergraph walks, the \mathcal{G} and its CE are not distinguishable. Note that in reality, the bottleneck for the number of sampled walks in machine learning-based methods is memory. Accordingly, even with tuning the number of samples for each dataset, the size of samples is bounded by a small number. This theorem shows that with a limited budget for walk sampling, ② alone is not enough to capture higher-order patterns.

Q2: Can addressing ① alone be sufficient to take advantage of higher-order interactions? To answer this question, we use the extended version of the edge-to-vertex dual graph concept for hypergraphs:

Definition 4 (Dual Hypergraph). Given a hypergraph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, the dual hypergraph of \mathcal{G} is defined as $\tilde{\mathcal{G}} = (\tilde{\mathcal{V}}, \tilde{\mathcal{E}})$, where $\tilde{\mathcal{V}} = \mathcal{E}$ and a hyperedge $\tilde{e} = \{e_1, e_2, \dots, e_k\} \in \tilde{\mathcal{E}}$ shows that $\bigcap_{i=1}^k e_i \neq \emptyset$.

To address ①, we need to see walks on hypergraphs as a sequence of hyperedges (instead of a sequence of pair-wise connected nodes). One can interpret this as a hypergraph walk on the dual hypergraph. That is, each hypergraph walk on the dual graph is a sequence of \mathcal{G} 's hyperedges: $e_1 \rightarrow e_2 \rightarrow \dots \rightarrow e_k$. However, as shown by Chitra and Raphael [40], each walk on hypergraphs with edge-independent weights for sampling vertices is equivalent to a simple walk on the (weighted) CE graph. To this end, addressing ② alone can be equivalent to sample walks on the CE of the dual hypergraph, which misses the higher-order interdependencies of hyperedges and their intersections.

Based on the above discussion, both ① and ② are required to capture higher-order interaction between nodes as well as higher-order interdependencies of hyperedges. The definition of SETWALKS (Definition 3) with *structural bias*, introduced in Equation 1, satisfies both ① and ②. In the next section, we discuss how a simple extension of SETWALKS can not only be more expressive than all existing walks on hypergraphs and their CEs, but its definition is universal, and all these methods are special cases of extended SETWALK.

C.1 Extension of SETWALKS

Random walks on hypergraphs are simple but less expressive methods for extracting network motifs, while SETWALKS are more complex patterns that provide more expressive motif extraction approaches. One can model the trade-off of simplicity and expressivity to connect all these concepts in a single notion of walks. To establish a connection between SETWALKS and existing walks on hypergraphs, as well as a universal random walk model on hypergraphs, we extend SETWALKS to r -SETWALKS, where parameter r controls the size of hyperedges that appear in the walk:

Definition 5 (r -SETWALK). Given a temporal hypergraph $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{X})$, and a threshold $r \in \mathbb{Z}^+$, a r -SETWALK with length ℓ on temporal hypergraph \mathcal{G} is a randomly generated sequence of hyperedges (sets):

$$\text{Sw} : (e_1, t_{e_1}) \rightarrow (e_2, t_{e_2}) \rightarrow \dots \rightarrow (e_\ell, t_{e_\ell}),$$

where $e_i \in \mathcal{E}$, $|e_i| \leq r$, $t_{e_{i+1}} < t_{e_i}$, and the intersection of e_i and e_{i+1} is not empty, $e_i \cap e_{i+1} \neq \emptyset$. In other words, for each $1 \leq i \leq \ell - 1$: $e_{i+1} \in \mathcal{E}^i(e_i)$. We use $\text{Sw}[i]$ to denote the i -th hyperedge-time pair in the SETWALK. That is, $\text{Sw}[i][0] = e_i$ and $\text{Sw}[i][1] = t_{e_i}$.

The only difference between this definition and [Definition 3](#) is that r -SETWALK limits hyperedges in the walk to hyperedges with size at most r . The sampling process of r -SETWALKS is the same as that of SETWALK (introduced in [Section 3.2](#) and [Appendix D](#)), while we only sample hyperedges with size at most r . Now to establish the connection of r -SETWALKS and existing walks on hypergraphs, we define the extended version of the clique expansion technique:

Definition 6 (r -Projected Hypergraph). *Given a hypergraph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ and an integer $r \geq 2$, we construct the (weighted) r -projected hypergraph of \mathcal{G} as a hypergraph $\hat{\mathcal{G}}_r = (\mathcal{V}, \hat{\mathcal{E}}_r)$, where for each $e = \{u_1, u_2, \dots, u_k\} \in \mathcal{E}$:*

1. if $k \leq r$: add e to the $\hat{\mathcal{E}}_r$,
2. if $k \geq r + 1$: add $e_i = \{u_{i_1}, u_{i_2}, \dots, u_{i_r}\}$ to $\hat{\mathcal{E}}_r$, for every possible $\{i_1, i_2, \dots, i_r\} \subseteq \{1, 2, \dots, k\}$.

Each of steps 1 or 2 can be done in a weighted manner. In other words, we approximate each hyperedge with a size of more than r with $\binom{k}{r}$ (weighted) hyperedges with size r . For example, when $r = 2$, the 2-projected graph of \mathcal{G} is equivalent to its clique expansion, and $r = \infty$ is the hypergraph itself. Furthermore, we define the Union Projected Hypergraph (UP hypergraph) as the union of all r -projected hypergraphs, i.e., $\mathcal{G}^* = (\mathcal{V}, \bigcup_{r=2}^{\infty} \hat{\mathcal{E}}_r)$. Note that the UP hypergraph has the downward closure property and is equivalent to the simplicial complex representation of the hypergraph \mathcal{G} . The next proposition establishes the universality of the r -SETWALK concept.

Proposition 2. *Edge-independent random walks on hypergraphs [37], edge-dependent random walks on hypergraphs [40], and simple random walks on the CE of hypergraphs are all special cases of r -SETWALK, when applied to the 2-projected graph, UP hypergraph, and 2-projected graph, respectively. Furthermore, all the above methods are less expressive than r -SETWALKS.*

The proof of this proposition is in [Appendix E.5](#).

D Efficient Hyperedge Sampling

For sampling SETWALKS, inspired by Wang et al. [33], we use two steps: ① Online score computation: we assign a set of scores to each incoming hyperedge. ② Iterative sampling: we use assigned scores in the previous step to sample hyperedges in a SETWALK.

Algorithm 1 Online Score Computation

Input: Given a hypergraph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ and $\alpha \in [0, 1]$

Output: A probability score for each vertex

```

1:  $P \leftarrow \emptyset$ ;
2: for  $(e, t) \in \mathcal{E}$  with an increasing order of  $t$  do
3:    $P_{e,t}[0] \leftarrow \exp(\alpha t)$ ;
4:    $P_{e,t}[1] \leftarrow 0, P_{e,t}[2] \leftarrow 0, P_{e,t}[3] \leftarrow \emptyset$ ;
5:   for  $u \in e$  do
6:     for  $e_n \in \mathcal{E}^t(u)$  do
7:       if  $e_n$  is not visited then
8:          $P_{e,t}[2] \leftarrow P_{e,t}[2] + \exp(\varphi(e_n, e))$ ;
9:          $P_{e,t}[3] \leftarrow P_{e,t}[3] \cup \{\exp(\varphi(e_n, e))\}$ ;
10:         $P_{e,t}[1] \leftarrow P_{e,t}[1] + \exp(\alpha \times t_n)$ ;
return  $P$ ;
```

Online Score Computation. The first part essentially works in an online manner and assigns each new incoming hyperedge e a four-tuple of scores:

$$\begin{aligned}
P_{e,t}[0] &= \exp(\alpha \times t), & P_{e,t}[1] &= \sum_{(e', t') \in \mathcal{E}^t(e)} \exp(\alpha \times t') \\
P_{e,t}[2] &= \sum_{(e', t') \in \mathcal{E}^t(e)} \exp(\varphi(e, e')), & P_{e,t}[3] &= \{\exp(\varphi(e, e'))\}_{(e', t') \in \mathcal{E}^t(e)}
\end{aligned}$$

Iterative Sampling. In the iterative sampling algorithm, we use pre-computed scores by [Algorithm 1](#) and sample a hyperedge (e, t) given a previously sampled hyperedge (e_p, t_p) . In the next proposition,

Algorithm 2 Iterative SETWALK Sampling

Input: Given a hypergraph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, $\alpha \in [0, 1]$, and previously sampled hyperedge (e_p, t_p)

Output: Next sampled hyperedge (e, t)

- 1: **for** $(e, t) \in \mathcal{E}^p(e_p)$ with an decreasing order of t **do**
 - 2: Sample $b \sim \text{UNIFORM}(0, 1)$;
 - 3: Get $P_{e,t}[0], P_{e_p,t_p}[1], P_{e_p,t_p}[2]$ and $\varphi(e, e_p)$ from the output of [Algorithm 1](#);
 - 4: $\mathcal{P} \leftarrow \text{Normalize } \frac{P_{e,t}[0]}{P_{e_p,t_p}[1]} \times \frac{\exp(\varphi(e, e_p))}{P_{e_p,t_p}[2]}$;
 - 5: **if** $b < \mathcal{P}$ **then return** (e, t) ;
- return** (e_X, t_X) ; $\triangleright (e_X, t_X)$ is a dummy empty hyperedge signaling the end of algorithm.
-

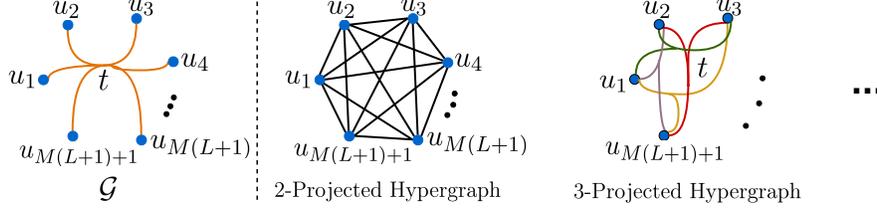


Figure 6: The example of a hypergraph \mathcal{G} and its 2- and 3-projected hypergraphs.

we show that this sampling algorithm samples each hyperedge with the probability mentioned in [Section 3.2](#).

Proposition 3. *Algorithm 2* sample a hyperedge (e, t) after (e_p, t_p) with a probability proportional to $\mathbb{P}[(e, t)|(e_p, t_p)]$ ([Equation 1](#)).

How can this sampling procedure capture higher-order patterns? As discussed in [Appendix C](#),

SETWALKS on \mathcal{G} can be interpreted as a random walk on the dual hypergraph of \mathcal{G} , $\tilde{\mathcal{G}}$. However, a simple (or hyperedge-independent) random walk on the dual hypergraph is equivalent to the walk on the CE of the dual hypergraph [\[40, 41\]](#), missing the higher-order dependencies of hyperedges. Inspired by Chitra and Raphael [\[40\]](#), we use hyperedge-dependent weights $\Gamma : \mathcal{V} \times \mathcal{E} \rightarrow \mathbb{R}^{\geq 0}$ and sample hyperedges with a probability proportional to $\exp\left(\sum_{u \in e \cap e_p} \Gamma(u, e)\Gamma(u, e')\right)$, where e_p is the previously sampled hyperedge. In the dual hypergraph $\tilde{\mathcal{G}} = (\mathcal{E}, \mathcal{V})$, we assign a score $\tilde{\Gamma} : \mathcal{E} \times \mathcal{V} \rightarrow \mathbb{R}^{\geq 0}$ to each pair of (e, u) as $\tilde{\Gamma}(e, u) = \Gamma(u, e)$. Now, a SETWALK with this sampling procedure is equivalent to the edge-dependent hypergraph walk on the dual hypergraph of \mathcal{G} with edge-dependent weight $\tilde{\Gamma}(\cdot)$. Chitra and Raphael [\[40\]](#) show that an edge-dependent hypergraph random walk can capture some information about higher-order interactions and is not equivalent to a simple walk on the weighted CE of the hypergraph. Accordingly, even on the dual hypergraph, SETWALK with this sampling procedure can capture higher-order dependencies of hyperedges and is not equivalent to a simple walk on the CE of the dual hypergraph $\tilde{\mathcal{G}}$. We conclude that, unlike existing random walks on hypergraphs [\[37, 38, 41, 79\]](#), SETWALK can capture both higher-order interactions of nodes, and, based on its sampling procedure, higher-order dependencies of hyperedges.

E Theoretical Results

E.1 Proof of [Theorem 1](#)

Theorem 1. *A random SETWALK is equivalent to neither the hypergraph random walk, the random walk on the CE graph, nor the random walk on the SE graph. Also, for a finite number of samples of each, SETWALK is more expressive than existing walks.*

Proof. In this proof, we focus on the hypergraph random walk and simple random walk on the CE. The proof for the SE graph is the same and also it has been proven that the SE graph and the CE of a hypergraph have close (or equal in uniform hypergraphs) Laplacian and have the same expressiveness power in the representation of hypergraphs [\[121–123\]](#).

First, note that each SETWALK can be approximately decomposed to a set of either hypergraph walks, simple random walks, or walk on the SE. Moreover, each of these walks can be mapped to

a corresponding SETWALK (but not a bijective mapping), by sampling hyperedges corresponding to each consecutive pair of nodes in these walks. Accordingly, SETWALKS includes the information provided by these walks and so its expressiveness is not less than these methods. To this end, next, we discuss two examples in two different tasks for which SETWALKS are successful while other walks fail.

① In the first task, we want to see if there exists a pair of hypergraphs with different semantics that SETWALKS can distinguish, but other walks cannot. We construct such hypergraphs. Let $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ be a hypergraph with $\mathcal{V} = \{u_1, u_2, \dots, u_N\}$ and $\mathcal{E} = \{(e, t_i)\}_{i=1}^T$, where $e = \{u_1, u_2, \dots, u_N\}$ and $t_1 < t_2 < \dots < t_T$. Also, let \mathcal{A} be an edge-independent hypergraph random walk (or random walk on the CE) sampling algorithm. Chitra and Raphael [40] show that each of these walks is equivalent to a random walk on the weighted CE. Assume that $\xi(\cdot)$ is a function that assigns weights to edges in $\mathcal{G}^* = (\mathcal{V}, \mathcal{E}^*)$, the weighted CE of the \mathcal{G} , such that a hypergraph random walk on \mathcal{G} is equivalent to a walk on this weighted CE graph. Next, we construct a weighted hypergraph $\mathcal{G}' = (\mathcal{V}, \mathcal{E}')$ with the same set of vertices but with $\mathcal{E}' = \bigcup_{k=1}^T \{(u_i, u_j), t_k\}_{u_i, u_j \in \mathcal{V}}$, such that each edge $e_{i,j} = (u_i, u_j)$ is associated with a weight $\xi(e_{i,j})$. Clearly, sampling procedure \mathcal{A} on \mathcal{G} and \mathcal{G}' are the same, while they have different semantics. For example, assume that both are collaboration networks. In \mathcal{G} , all vertices have published a single paper together, while in \mathcal{G}' , each pair of vertices have published a separate paper together. The proof for the hypergraph random walk with hyperedge-dependent weights is the same, while we construct weights of the hypergraph \mathcal{G}' based on the sampling probability of hyperedges in the hypergraph random walk procedure.

② Next, in the second task, we investigate the expressiveness of these walks for reconstructing hyperedges. That is, we want to see that given a perfect classifier, can these walks provide enough information to detect higher-order patterns in the network. To this end, we show that for a finite number of samples of each walk, SETWALK is more expressive than all of these walks in detecting higher-order patterns. Let M be the maximum number of samples and L be the maximum length of walks, we show that for any $M \geq 2$ and $L \geq 2$ there exists a pair of hypergraphs \mathcal{G} , with higher-order interactions, and \mathcal{G}' , with pairwise interactions, such that SETWALKS can distinguish them, while they are indistinguishable by any of these walks. We construct a temporal hypergraph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ as a hypergraph with $\mathcal{V} = \{u_1, u_2, \dots, u_{M(L+1)+1}\}$ and $\mathcal{E} = \{(e, t_i)\}_{i=1}^L$, where $e = \{u_1, u_2, \dots, u_{M(L+1)+1}\}$ and $t_1 < t_2 < \dots < t_L$. We further construct $\mathcal{G}' = (\mathcal{V}, \mathcal{E}')$ with the same set of vertices but with $\mathcal{E}' = \bigcup_{k=1}^L \{(u_i, u_j), t_k\}_{u_i, u_j \in \mathcal{V}}$. Figure 6 illustrates \mathcal{G} and its projected graphs at a given timestamp $t \in \{t_1, t_2, \dots, t_L\}$.

SETWALK with only one sample, Sw, can distinguish interactions in these two hypergraphs. That is, let $\text{Sw} : (e, t_L) \rightarrow (e, t_{L-1}) \rightarrow \dots \rightarrow (e, t_1)$ be the sample SETWALK from \mathcal{G} (note that masking the time, this is the only SETWALK on \mathcal{G} , so in any case the sampled SETWALK is Sw). Since all interactions in \mathcal{G}' are pairwise, any sampled SETWALK on \mathcal{G}' , Sw', includes only pairwise interactions, so $\text{Sw} \neq \text{Sw}'$, in any case. Accordingly, SETWALK can distinguish interactions in these two hypergraphs.

Since the output of hypergraph random walks, simple walks on the CE, and walks on the SE include only pairwise interactions, it seems that they are unable to detect higher-order patterns, so are unable to distinguish these two hypergraphs. However, one might argue that by having a large number of sampled walks and using a perfect classifier, which learned the distribution of sampled random walks and can detect whether a set of sampled walks is from a higher-order interaction, we might be able to detect higher-order interactions. To this end, we next assume that we have a perfect classifier $C(\cdot)$ that can detect whether a set of sampled hypergraph walks, simple walks on the CE, or walks on the SE are sampled from a higher-order structure or pair-wise patterns. Next, we show that hypergraph random walks cannot provide enough information about every vertex for $C(\cdot)$ to detect whether all vertices in \mathcal{V} shape a hyperedge. To this end, assume that we sample $S = \{W_1, W_2, \dots, W_M\}$ walks from hypergraph \mathcal{G} and $S' = \{W'_1, W'_2, \dots, W'_M\}$ walks from hypergraph \mathcal{G}' . In the best case scenario, since $C(\cdot)$ is a perfect classifier, it can detect that \mathcal{G}' includes only pair-wise interactions based on sampled walk S' . To distinguish these two hypergraphs, we need $C(\cdot)$ to detect sampled walks from \mathcal{G} (i.e., S) that come from a higher-order pattern. For any M sampled walks with length L from \mathcal{G} , we observe at most $M \times (L + 1)$ vertices, so we have information about at most $M \times (L + 1)$ vertices, unable to capture any information about the neighborhood of at least one vertex. Due to the symmetry of vertices, without loss of generality, we can assume that this vertex is u_1 . This means that with these M sampled hypergraph random walks with length L , we are not able to provide any information about node u_1 at any timestamp for $C(\cdot)$. Therefore, even a perfect classifier $C(\cdot)$ cannot verify whether u_1

is a part of higher-order interaction or pair-wise interaction, which completes the proof. Note that the proof for the simple random walk is completely the same. \square

Remark 1. Note that while the first task investigates the expressiveness of these methods with respect to their sampling procedure, the second tasks discuss the limitation and difference in their outputs.

Remark 2. Note that in reality, we can have neither an unlimited number of samples nor an unlimited walk length. Also, the upper bound for the number of samples or walk length depends on the RAM of the machine on which the model is being trained. In our experiments, we observe that usually, we cannot sample more than 125 walks with a batch size of 32.

E.2 Proof of Theorem 2

Before discussing the proof of Theorem 2 we first formally define what missing information means in this context.

Definition 7 (Missing Information). We say a pooling strategy like $\Psi(\cdot)$ misses information if there is a model \mathcal{M} such that using $\Psi(\cdot)$ on top of the \mathcal{M} (call it $\hat{\mathcal{M}}$) decreases the expressive power of \mathcal{M} . That is, $\hat{\mathcal{M}}$ has less expressive power than \mathcal{M} .

Theorem 2. Given an arbitrary positive integer $k \in \mathbb{Z}^+$, let $\Psi(\cdot)$ be a pooling function such that for any set $S = \{w_1, \dots, w_d\}$:

$$\Psi(S) = \sum_{\substack{S' \subseteq S \\ |S'|=k}} f(S'), \quad (13)$$

where f is some function. Then the pooling function can cause missing information, limiting the expressiveness of the method to applying to the projected (hyper)graph of the hypergraph.

Proof. The main intuition of this theorem is that a pooling function needs to capture higher-order dependencies of its input's elements and if it can be decomposed to a summation of functions that capture lower-order dependencies, it misses information. We show that, in the general case for a given $k \in \mathbb{Z}^+$, the pooling function $\Psi(\cdot)$ when applied to a hypergraph \mathcal{G} is at most as expressive as $\Psi(\cdot)$ when applied to the k -projected hypergraph of \mathcal{G} (Definition 6). Let $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ be a hypergraph with $\mathcal{V} = \{u_1, u_2, \dots, u_{k+1}\}$ and $\mathcal{E} = \{(\mathcal{V}, t)\} = \{(\{u_1, u_2, \dots, u_{k+1}\}, t)\}$ for a given time t , and $\hat{\mathcal{G}} = (\mathcal{V}, \hat{\mathcal{E}})$ be its k -projected graph, i.e., $\hat{\mathcal{E}} = \{(e_1, t), \dots, (e_{\binom{k+1}{k}}, t)\}$, where $e_i \subset \{u_1, u_2, \dots, u_{k+1}\}$ such that $|e_i| = k$. Applying pooling function $\Psi(\cdot)$ on the hypergraph \mathcal{G} is equivalent to applying $\Psi(\cdot)$ to the hyperedge $(\mathcal{V}, t) \in \mathcal{E}$, which provides $\Psi(\mathcal{V}) = \sum_{i=1}^{k+1} f(e_i)$. On the other hand, applying $\Psi(\cdot)$ on projected graph $\hat{\mathcal{G}}$ means applying it on each hyperedge $e_i \in \hat{\mathcal{E}}$. Accordingly, since for each hyperedge $e_i \in \hat{\mathcal{E}}$ we have $\Psi(e_i) = f(e_i)$, all captured information by pooling function $\Psi(\cdot)$ on $\hat{\mathcal{G}}$ is the set of $S = \{f(e_i)\}_{i=1}^{k+1}$. It is clear that $\Psi(\mathcal{V}) = \sum_{i=1}^{k+1} f(e_i)$ is less informative than $S = \{f(e_i)\}_{i=1}^{k+1}$ as it is the summation of elements in S (in fact, $\Psi(\mathcal{V})$ cannot capture the non-linear combinations of positional encodings of vertices, while S can). Accordingly, the provided information by applying $\Psi(\cdot)$ on \mathcal{G} cannot be more informative than applying $\Psi(\cdot)$ on the \mathcal{G} 's k -projected hypergraph. \square

Remark 3. Note that the pooling function $\Psi(\cdot)$ is defined on a (hyper)graph and gets only (hyper)edges as input.

Remark 4. Although $\Psi(\cdot) = \text{MEAN}(\cdot)$ cannot be written as Equation 13, we can simply see that the above proof works for this pooling function as well.

E.3 Proof of Theorem 3

Theorem 3. SETMIXER is permutation invariant and is a universal approximator of invariant multi-set functions. That is, SETMIXER can approximate any invariant multi-set function.

Proof. Let $\pi(S)$ be a given permutation of set S , we aim to show that $\Psi(S) = \Psi(\pi(S))$. We first recall the SETMIXER and its two phases: Let $S = \{\mathbf{v}_1, \dots, \mathbf{v}_d\}$, where $\mathbf{v}_i \in \mathbb{R}^{d_1}$, be the input set and $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_d]^T \in \mathbb{R}^{d \times d_1}$ be its matrix representation:

$$\Psi(\mathbf{V}) = \text{MEAN}\left(\mathbf{H}_{\text{token}} + \sigma\left(\text{LayerNorm}(\mathbf{H}_{\text{token}}) \mathbf{W}_s^{(1)}\right) \mathbf{W}_s^{(2)}\right), \quad (\text{Channel Mixer})$$

where

$$\mathbf{H}_{\text{token}} = \mathbf{V} + \sigma\left(\text{Softmax}\left(\text{LayerNorm}(\mathbf{V})^T\right)\right)^T. \quad (\text{Token Mixer})$$

Let $\pi(\mathbf{V}) = [\mathbf{v}_{\pi(1)}, \dots, \mathbf{v}_{\pi(d)}]^T$ be a permutation of the input matrix \mathbf{V} . In the token mixer phase, None of `LayerNorm`, `Softmax`, and activation function $\sigma(\cdot)$ can affect the order of elements (note that `Softmax` is applied row-wise). Accordingly, we can see the output of the token mixer is permuted by $\pi(\cdot)$:

$$\begin{aligned} \mathbf{H}_{\text{token}}(\pi(\mathbf{V})) &= \pi(\mathbf{V}) + \sigma\left(\text{Softmax}\left(\text{LayerNorm}(\pi(\mathbf{V}))^T\right)\right)^T \\ &= \pi(\mathbf{V}) + \pi\left(\sigma\left(\text{Softmax}\left(\text{LayerNorm}(\mathbf{V})^T\right)\right)^T\right) \\ &= \pi\left(\mathbf{V} + \sigma\left(\text{Softmax}\left(\text{LayerNorm}(\mathbf{V})^T\right)\right)^T\right) \\ &= \pi(\mathbf{H}_{\text{token}}(\mathbf{V})). \end{aligned} \quad (14)$$

Next, in the channel mixer, by using Equation 14 we have:

$$\begin{aligned} \Psi(\pi(\mathbf{V})) &= \text{MEAN}\left(\pi(\mathbf{H}_{\text{token}}) + \sigma\left(\text{LayerNorm}(\pi(\mathbf{H}_{\text{token}})) \mathbf{W}_s^{(1)}\right) \mathbf{W}_s^{(2)}\right) \\ &= \text{MEAN}\left(\pi(\mathbf{H}_{\text{token}}) + \pi\left(\sigma\left(\text{LayerNorm}(\mathbf{H}_{\text{token}}) \mathbf{W}_s^{(1)}\right)\right) \mathbf{W}_s^{(2)}\right) \\ &= \text{MEAN}\left(\pi(\mathbf{H}_{\text{token}}) + \pi\left(\sigma\left(\text{LayerNorm}(\mathbf{H}_{\text{token}}) \mathbf{W}_s^{(1)}\right) \mathbf{W}_s^{(2)}\right)\right) \\ &= \text{MEAN}\left(\pi\left(\mathbf{H}_{\text{token}} + \mathbf{W}_s^{(2)} \sigma\left(\text{LayerNorm}(\mathbf{H}_{\text{token}}) \mathbf{W}_s^{(1)}\right)\right)\right) \\ &= \Psi(\mathbf{V}). \end{aligned} \quad (15)$$

In the last step, we use the fact that `MEAN`(\cdot) is permutation invariant. Based on Equation 15 we can see that `SETMIXER` is permutation invariant.

Since the token mixer is just normalization it is inevitable and cannot affect the expressive power of `SETMIXER`. Also, channel mixer is a 2-layer MLP, which is the universal approximator of any function. Therefore, `SETMIXER` is a universal approximator. \square

E.4 Proof of Theorem 4

Theorem 4. *The set-based anonymization method is more expressive than any existing anonymization strategies on the CE of the hypergraph. More precisely, there exists a pair of hypergraphs $\mathcal{G}_1 = (\mathcal{V}_1, \mathcal{E}_1)$ and $\mathcal{G}_2 = (\mathcal{V}_2, \mathcal{E}_2)$ with different structures (i.e., $\mathcal{G}_1 \not\cong \mathcal{G}_2$) that are distinguishable by our anonymization process and are not distinguishable by the CE-based methods.*

Proof. To the best of our knowledge, there exist two anonymization processes for random walks by Wang et al. [33] and Micali and Zhu [45]. Both of these methods are designed for graphs and to adapt them to hypergraphs we need to apply them to the (weighted) CE. Here, we focus on the process designed by Wang et al. [33], which is more informative than the other. The proof for the Micali and Zhu [45] process is the same. Note that the goal of this theorem is to investigate whether a method can distinguish a hypergraph from its CE. Accordingly, this theorem does not provide any information about the expressivity of these methods in terms of the isomorphism test.

The proposed 2-step anonymization process can be seen as a positional encoding for both vertices and hyperedges. Accordingly, it is expected to assign different positional encodings to vertices and hyperedges of two non-isomorphism hypergraphs. To this end, we construct the same hypergraphs as in the proof of Theorem 1. Let M be the number of sampled `SETWALKS` with length L . We construct a temporal hypergraph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ as a hypergraph with $\mathcal{V} = \{u_1, u_2, \dots, u_{M(L+1)+1}\}$ and $\mathcal{E} = \{(e, t_i)\}_{i=1}^L$, where $e = \{u_1, u_2, \dots, u_{M(L+1)+1}\}$ and $t_1 < t_2 < \dots < t_L$. We further construct $\mathcal{G}' = (\mathcal{V}, \mathcal{E}')$ with the same set of vertices but with $\mathcal{E}' = \bigcup_{k=1}^L \{(u_i, u_j), t_k\}_{u_i, u_j \in \mathcal{V}}$. As we have seen in Theorem 1, random walks on the CE of the hypergraph cannot distinguish these two hypergraphs. Since CAW [33] also uses simple random walks, it cannot distinguish these two hypergraphs. Accordingly, after its anonymization process, it again cannot distinguish these two hypergraphs.

The main part of the proof is to show that in our method, the assigned positional encodings are different in these hypergraphs. The first step is to assign each node a positional encoding. Masking

the timestamps, there is only one SETWALK in the \mathcal{G} . Accordingly, the positional encodings of nodes in \mathcal{G} are the same and non-zero. Given a SETWALK with length L we might see at most $L \times (d_{\max} - 1) + 1$ nodes, where d_{\max} is the maximum size of hyperedges in the hypergraph. Accordingly, with M samples on \mathcal{G}' , which $d_{\max} = 2$, we can see at most $M \times (L + 1)$ vertices. Therefore, in any case, we assign a zero vector to at least one vertex. This proves that the positional encodings by SETWALKS are different in these two hypergraphs, and if the assigned hidden identities to counterpart nodes are different, clearly, feeding them to the SETMIXER results in different hyperedge encodings.

Note that each SETWALK can be decomposed into a set of causal anonymous walks [33]. Accordingly, it includes the information provided by these walks, so its expressiveness is not less than the CAW method on hypergraphs, which completes the proof of the theorem. \square

Although the above statement completes the proof, next we discuss that even given the same positional encodings for vertices in these two hypergraphs, SETMIXER can capture higher-order interactions by capturing the size of the hyperedge. Recall token mixer phase in SETMIXER:

$$\mathbf{H}_{\text{token}} = \mathbf{V} + \sigma \left(\text{Softmax} \left(\text{LayerNorm}(\mathbf{V})^T \right) \right)^T,$$

where $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_{M(L+1)+1}]^T \in \mathbb{R}^{(M(L+1)+1) \times d_1}$ and $\mathbf{v}_i \neq \mathbf{0}_{1 \times d_1}$ represents the positional encoding of u_i in \mathcal{G} . We assumed that the positional encoding of u_i in \mathcal{G}' is the same. The input of the token mixer phase on \mathcal{G} is \mathcal{V} as all of them are connected by a hyperedge. Then we have:

$$(\mathbf{H}_{\text{token}})_{i,j} = \mathbf{v}_{i,j} + \sigma \left(\frac{\exp(\mathbf{v}_{i,j})}{\sum_{k=1}^{M(L+1)+1} \exp(\mathbf{v}_{k,j})} \right). \quad (16)$$

On the other hand, when applied to hypergraph \mathcal{G}' and (u_{k_1}, u_{k_2}) . We have:

$$(\mathbf{H}'_{\text{token}})_{i,j} = \mathbf{v}_{i,j} + \sigma \left(\frac{\exp(\mathbf{v}_{i,j})}{\exp(\mathbf{v}_{k_1,j}) + \exp(\mathbf{v}_{k_2,j})} \right), \quad i \in \{k_1, k_2\}. \quad (17)$$

Since we use zero padding, for any $i \geq 3$, $(\mathbf{H}_{\text{token}})_{i,j} \neq 0$ and $(\mathbf{H}'_{\text{token}})_{i,j} = 0$. These zero rows, which capture the size of the hyperedge, result in different encodings for each connection.

Remark 5. *To the best of our knowledge, the only anonymization process that is used on hypergraphs is by Liu et al. [78], which uses simple walks on the CE and is the same as Wang et al. [33]. Accordingly, it also suffers from the above limitation. Also, note that this theorem shows the limitation of these anonymization procedures when simply adopted to hypergraphs.*

E.5 Proof of Proposition 2

Proposition 2. *Edge-independent random walks on hypergraphs [37], edge-dependent random walks on hypergraphs [40], and simple random walks on the CE of hypergraphs are all special cases of r -SETWALK, when applied to the 2-projected graph, UP graph, and 2-projected graph, respectively. Furthermore, all the above methods are less expressive than r -SETWALKS.*

Proof. For the first part, we discuss each walk separately:

- ① Simple random walks on the CE of the hypergraphs: We perform 2-SETWALKS on the (weighted) 2-projected hypergraph with $\Gamma(\cdot) = 1$. Accordingly, for every two adjacent edges in the 2-Projected graph like e and e' , we have $\varphi(e, e') = 1$. Therefore, it is equivalent to a simple random walk on the CE (2-projected graph).
- ② Edge-independent random walks on hypergraphs: As is shown by Chitra and Raphael [40], each edge-independent random walk on hypergraphs is equivalent to a simple random walk on the (weighted) CE of the hypergraph. Therefore, as discussed in ①, these walks are a special case of r -SETWALKS, when $r = 2$ and applied to (weighted) 2-Projected hypergraph.
- ③ Edge-dependent random walks on hypergraphs: Let $\Gamma'(e, u)$ be an edge-dependent weight function used in the hypergraph random walk sampling. For each node u in the UP hypergraph, we store the set of $\Gamma'(e, u)$ that e is a maximal hyperedge that u belongs to. Note that there might be several maximal hyperedges that u belongs to. Now, we perform 2-SETWALK sampling on the UP hypergraph with these weights and in each step, we sample each hyperedge with weight $\Gamma(u, e) = \Gamma'(e, u)$. It is

straightforward to show that given this procedure, the sampling probability of a hyperedge is the same in both cases. Therefore, edge-dependent random walks on hypergraphs are equivalent to 2-SETWALKS when applied to the UP hypergraph.

As discussed above, all these walks are special cases of r -SETWALKS and cannot be more expressive than r -SETWALKS. Also, as discussed in [Theorem 1](#), all these walks are less expressive than SETWALKS, which are also special cases of r -SETWALKS, when $r = \infty$. Accordingly, all these methods are less expressive than r -SETWALKS. \square

F Experimental Setup Details

F.1 Datasets

We use 10 publicly available³ benchmark datasets, whose descriptions are as follows:

- **NDC Class** [6]: The NDC Class dataset is a temporal higher-order network, in which each hyperedge corresponds to an individual drug, and the nodes contained within the hyperedges represent class labels assigned to these drugs. The timestamps, measured in days, indicate the initial market entry of each drug. Here, hyperedge prediction aims to predict future drugs.
- **NDC Substances** [6]: The NDC Substances is a temporal higher-order network, where each hyperedge represents an NDC code associated with a specific drug, while the nodes represent the constituent substances of the drug. The timestamps, measured in days, indicate the initial market entry of each drug. The hyperedge prediction task is the same as NDC Classes dataset.
- **High School** [6, 92]: The High School is a temporal higher-order network dataset constructed from interactions recorded by wearable sensors in a high school setting. The dataset captures a high school contact network, where each student/teacher is represented as a node and each hyperedge shows Face-to-face contact among individuals. Interactions were recorded at a resolution of 20 seconds, capturing all interactions that occurred within the previous 20 seconds. Node labels in this data are the class of students, and we focus on the node class "PSI" in our classification tasks.
- **Primary School** [6, 93]: The primary school dataset resembles the high school dataset, differing only in terms of the school level from which the data is collected. Node labels in this data are the class of students, and we focus on the node class "Teachers" in our classification tasks.
- **Congress Bill** [6, 94, 95]: Each node in this dataset represents a US Congressperson. Each hyperedge is a legislative bill in both the House of Representatives and the Senate, connecting the sponsors and co-sponsors of each respective bill. The timestamps, measured in days, indicate the date when each bill was introduced.
- **Email Enron** [6]: In this dataset nodes are email addresses at Enron and hyperedges are formed by emails, connecting the sender and recipients of each email. The timestamps have a resolution of milliseconds.
- **Email Eu** [6, 96]: In this dataset, the nodes represent email addresses associated with a European research institution. Each hyperedge consists of the sender and all recipients of the email. The timestamps in this dataset are measured with a resolution of 1 second.
- **Question Tags M (Math sx)** [6]: This dataset consists of nodes representing tags and hyperedges representing sets of tags applied to questions on math.stackexchange.com. The timestamps in the dataset are recorded at millisecond resolution and have been normalized to start at 0.
- **Question Tags U (Ask Ubuntu)** [6]: In this dataset, the nodes represent tags, and the hyperedges represent sets of tags applied to questions on askubuntu.com. The timestamps in the dataset are recorded with millisecond resolution and have been normalized to start at 0.

³<https://www.cs.cornell.edu/~arb/data/>

Table 3: Dataset statistics. HeP: **H**yperedge **P**rediction, NC: **N**ode **C**lassification

Dataset	NDC Class	High School	Primary School	Congress Bill	Email Enron	Email Eu	Question Tags M	Users-Threads	NDC Substances	Question Tags U
V	1,161	327	242	1,718	143	998	1,629	125,602	5,311	3,029
E	49,724	172,035	106,879	260,851	10,883	234,760	822,059	192,947	112,405	271,233
#Timestamps	5,891	7,375	3,100	5,936	10,788	232,816	822,054	189,917	7,734	271,233
Task	HeP	HeP& NC	HeP & NC	HeP	HeP	HeP	HeP	HeP	HeP	HeP

- **Users-Threads** [6]: In this dataset, the nodes represent users on askubuntu.com, and a hyperedge is formed by users participating in a thread that lasts for a maximum duration of 24 hours. The timestamps in the dataset denote the time of each post, measured in milliseconds but normalized such that the earliest post begins at 0.

The statistics of these datasets can be found in [Table 3](#).

F.2 Baselines

We compare our method to eight previous state-of-the-art methods and baselines on the hyperedge prediction task:

- **CHESHIRE** [11]: Chebyshev spectral hyperlink predictor (CHESHIRE), is a hyperedge prediction methods that initializes node embeddings by directly passing the incidence matrix through a one-layer neural network. CHESHIRE treats a hyperedge as a fully connected graph (clique) and uses a Chebyshev spectral GCN to refine the embeddings of the nodes within the hyperedge. The Chebyshev spectral GCN leverages Chebyshev polynomial expansion and spectral graph theory to learn localized spectral filters. These filters enable the extraction of local and composite features from graphs that capture complex geometric structures. The model with code provided is [here](#).
- **HYPER-SAGCN** [26]: Self-attention-based graph convolutional network for hypergraphs (HyperSAGCN) utilizes a Spectral Aggregated Graph Convolutional Network (SAGCN) to refine the embeddings of nodes within each hyperedge. HyperSAGCN generates initial node embeddings by hypergraph random walks and combines node embeddings by MEAN (.) pooling to compute the embedding of hyperedge. The model with code provided is [here](#).
- **NHP** [87]: Neural Hyperlink Predictor (NHP), is an enhanced version of HyperSAGCN. NHP initializes node embeddings using Node2Vec on the CE graph and then uses a novel maximum minimum-based pooling function that enables adaptive weight learning in a task-specific manner, incorporating additional prior knowledge about the nodes. The model with code provided is [here](#).
- **HPLSF** [89]: Hyperlink Prediction using Latent Social Features (HPLSF) is a probabilistic method. It leverages the homophily property of the networks and introduces a latent feature learning approach, incorporating the use of entropy in computing hyperedge embedding. The model with code provided is [here](#).
- **HPRA** [88]: Hyperlink Prediction Using Resource Allocation (HPRA) is a hyperedge prediction method based on the resource allocation process. HPRA calculates a hypergraph resource allocation (HRA) index between two nodes, taking into account direct connections and shared neighbors. The HRA index of a candidate hyperedge is determined by averaging all pairwise HRA indices between the nodes within the hyperedge. The model with code provided is [here](#).
- **CE-CAW**: This model is a baseline that we apply CAW [33] on the CE of the hypergraph. CAW is a temporal edge prediction method that uses causal anonymous random walks to capture the dynamic laws of the network in an inductive manner. The model with code provided is [here](#).
- **CE-EVOLVEGCN**: This is a snapshot-based temporal graph learning method that we apply EVOLVEGCN [90], which uses RNNs to estimate the GCN parameters for the future snapshots, on the CE of the hypergraph. The model with code provided is [here](#).
- **CE-GCN**: We apply Graph Convolutional Networks [91] to the CE of the hypergraph to obtain node embeddings. Next, we use MLP to predict edges. The implementation is provided in the Pytorch Geometric library.

Table 4: Hyperparameters used in the grid search.

Datasets	Sampling Number M	Sampling Time Bias α	SETWALK Length m	Hidden dimensions
NDC Class	4, 8, 16, 32, 64, 128	$\{0.5, 2.0, 20, 200\} \times 10^{-7}$	2, 3, 4, 5	32, 64, 128
High School	4, 8, 16, 32, 64, 128	$\{0.5, 2.0, 20, 200\} \times 10^{-7}$	2, 3, 4, 5	32, 64, 128
Primary School	4, 8, 16, 32, 64, 128	$\{0.5, 2.0, 20, 200\} \times 10^{-7}$	2, 3, 4, 5	32, 64, 128
Congress Bill	8, 16, 32, 64	$\{0.5, 2.0, 20, 200\} \times 10^{-7}$	2, 3, 4, 5	32, 64, 128
Email Enron	8, 16, 32, 64	$\{0.5, 2.0, 20, 200\} \times 10^{-7}$	2, 3, 4, 5	32, 64, 128
Email Eu	8, 16, 32, 64	$\{0.5, 2.0, 20, 200\} \times 10^{-7}$	2, 3, 4	32, 64, 128
Question Tags M	8, 16, 32, 64	$\{0.5, 2.0, 20, 200\} \times 10^{-7}$	2, 3, 4	32, 64, 128
Users-Threads	8, 16, 32, 64	$\{0.5, 2.0, 20, 200\} \times 10^{-7}$	2, 3, 4	32, 64, 128
NDC Substances	8, 16, 32, 64	$\{0.5, 2.0, 20, 200\} \times 10^{-7}$	2, 3, 4	32, 64, 128
Question Tags U	8, 16, 32, 64	$\{0.5, 2.0, 20, 200\} \times 10^{-7}$	2, 3, 4	32, 64, 128

For node classification, we use additional five state-of-the-art deep hypergraph learning methods and a CE-based baseline:

- **HYPERGCN** [19]: This is a generalization of GCNs to hypergraphs, where it uses hypergraph Laplacian to define convolution.
- **ALLDEEPSETS** and **ALLSETTRANSFORMER** [25]: These two methods are two variants of the general message passing framework, Allset, on hypergraphs, which are based on the aggregation of messages from nodes to hyperedges and from hyperedges to nodes.
- **UNIGCNII** [28]: Is an advanced variant of **UNIGNN**, a general framework for message passing on hypergraphs.
- **ED-HNN** [124]: Inspired by hypergraph diffusion algorithms, this method uses star expansions of hypergraphs with standard message passing neural networks.
- **CE-GCN**: We apply Graph Convolutional Networks [91] to the CE of the hypergraph to obtain node embeddings. Next, we use MLP to predict the labels of nodes. The implementation is provided in the Pytorch Geometric library. [91]

For all the baselines, we set all sensitive hyperparameters (e.g., learning rate, dropout rate, batch size, etc.) to the values given in the paper that describes the technique. Following [60], for deep learning methods, we tune their hidden dimensions via grid search to be consistent with what we did for CAT-WALK. We exclude HPLSF [89] and HPRA [88] from inductive hyperedge prediction as it does not apply to them.

F.3 Implementation and Training Details

In addition to hyperparameters and modules (activation functions) mentioned in the main paper, here, we report the training hyperparameters of CAT-WALK: On all datasets, we use a batch size of 64 and set learning rate = 10^{-4} . We also use an early stopping strategy to stop training if the validation performance does not increase for more than 5 epochs. We use the maximum training epoch number of 30 and dropout layers with rate = 0.1. Other hyperparameters used in the implementation can be found in the README file in the supplement.

Also, for tuning the model’s hyperparameters, we systematically tune them using grid search. The search domains of each hyperparameter are reported in Table 4. Note that, the last column in Table 4 reports the search domain for hidden dimensions of modules in CAT-WALK, including SETMIXER, MLP-MIXER, and MLPs. Also, we tune the last layer pooling strategy with two options: SETMIXER or MEAN(.) whichever leads to a better performance.

We implemented our method in Python 3.7 with *PyTorch* and run the experiments on a Linux machine with *nvidia RTX A4000* GPU with 16GB of RAM.

Table 5: Performance on node classification: Mean ACC (%) \pm standard deviation. Boldfaced letters shaded blue indicate the best result, while gray shaded boxes indicate results within one standard deviation of the best result.

	Methods	High School	Primary School	Average Performance
Inductive	CE-GCN	76.24 \pm 2.99	79.03 \pm 3.16	77.63 \pm 3.07
	HYPERGCN	83.91 \pm 3.05	86.17 \pm 3.40	85.04 \pm 3.23
	HYPERAGCN	84.89 \pm 3.80	82.13 \pm 3.69	83.51 \pm 3.75
	ALLDEEPSSETS	85.67 \pm 4.17	81.43 \pm 6.77	83.55 \pm 5.47
	UNIGCNII	88.36 \pm 3.78	88.27 \pm 3.52	88.31 \pm 3.63
	ALLSETTRANSFORMER	91.19 \pm 2.85	90.00 \pm 4.35	90.59 \pm 3.60
	ED-HNN	89.23 \pm 2.98	90.83 \pm 3.02	90.03 \pm 3.00
	CAT-WALK	88.99 \pm 4.76	93.28 \pm 2.41	91.13 \pm 3.58
Transductive	CE-GCN	78.93 \pm 3.11	77.46 \pm 2.97	78.20 \pm 3.04
	HYPERGCN	84.90 \pm 3.59	85.23 \pm 3.06	85.07 \pm 3.33
	HYPERAGCN	84.52 \pm 3.18	83.27 \pm 2.94	83.90 \pm 3.06
	ALLDEEPSSETS	85.97 \pm 4.05	80.20 \pm 10.18	83.09 \pm 7.12
	UNIGCNII	89.16 \pm 4.37	90.29 \pm 4.01	89.73 \pm 4.19
	ALLSETTRANSFORMER	90.75 \pm 3.13	89.80 \pm 2.55	90.27 \pm 2.84
	ED-HNN	91.41 \pm 2.36	91.74 \pm 2.62	91.56 \pm 2.49
	CAT-WALK	90.66 \pm 4.96	93.20 \pm 2.45	91.93 \pm 3.71

G Additional Experimental Results

G.1 Results on More Datasets

Due to the space limit, we report the AUC results on only eight datasets in Section 4. Table 6 reports both AUC and average precision (AP) results on all 10 datasets in both inductive and transductive hyperedge prediction tasks.

G.2 Node Classification

In the main text, we focus on the hyperedge prediction task. Here we describe how CAT-WALK can be used for node classification tasks.

For each node u_0 in the training set, we sample $\max\{\deg(u_0), 10\}$ hyperedges such as $e_0 = \{u_0, u_1, \dots, u_k\}$. Next, for each sampled hyperedge we sample M SETWALKS with length m starting from each $u_i \in e_0$ to construct $\mathcal{S}(u_i)$. Next, we anonymize each hyperedge that appears in at least one SETWALK in $\bigcup_{i=0}^k \mathcal{S}(u_i)$ by Equation 3 and then use the MLP-MIXER module to encode each $sw \in \bigcup_{i=0}^k \mathcal{S}(u_i)$. To encode each node $u_i \in e_0$, we use MEAN(.) pooling over SETWALKS in $\mathcal{S}(u_i)$. Finally, for node classification task, we use a 2-layer perceptron over the node encodings to make the final prediction.

Table 5 reports the results of dynamic node classification tasks on High School and Primary School datasets. CAT-WALK achieves the best or on-par performance on dynamic node classification tasks. While all baselines are specifically designed for node classification tasks, CAT-WALK achieves superior results due to ① its ability to incorporate temporal properties (both from SETWALKS and our time encoding module), which helps to learn underlying dynamic laws of the network, and ② its two-step set-based anonymization process that hides node identities from the model. Accordingly, CAT-WALK can learn underlying patterns needed for the node classification task, instead of using node identities, which might cause memorizing vertices.

G.3 Performance in Average Precision

In addition to the AUC, we also compare our model with baselines with respect to Average Precision (AP). Table 6 reports both AUC and AP results on all 10 datasets in inductive and transductive hyperedge prediction tasks. As discussed in Section 4, CAT-WALK due to its ability to capture both temporal and higher-order properties of the hypergraphs, achieves superior performance and outperforms all baselines in both transductive and inductive settings with a significant margin.

Table 6: Performance on hyperedge prediction: AUC and Average Precision (%) \pm standard deviation. Boldfaced letters shaded blue indicate the best result, while gray shaded boxes indicate results within one standard deviation of the best result. N/A: the method has computational issues.

Metric	NDC Class		High School		Primary School		Congress Bill		Email Enron	
	AUC	AP								
Strongly Inductive										
CE-GCN	52.31 \pm 2.99	54.33 \pm 2.48	60.54 \pm 2.06	59.92 \pm 2.25	52.34 \pm 2.75	56.41 \pm 2.06	49.18 \pm 3.61	53.85 \pm 3.92	63.04 \pm 1.80	57.70 \pm 2.27
CE-EVOLVEGCN	49.78 \pm 3.13	55.24 \pm 3.56	46.12 \pm 3.83	52.87 \pm 3.48	58.01 \pm 2.56	55.68 \pm 2.41	54.00 \pm 1.84	50.27 \pm 1.76	57.31 \pm 4.19	54.52 \pm 3.79
CE-CAW	76.45 \pm 0.29	78.58 \pm 1.32	83.73 \pm 1.42	82.96 \pm 1.04	80.31 \pm 1.46	82.84 \pm 1.71	75.38 \pm 1.25	77.19 \pm 1.38	70.81 \pm 1.13	72.07 \pm 1.52
NHP	70.53 \pm 4.95	68.18 \pm 4.31	65.29 \pm 3.80	62.86 \pm 3.74	70.86 \pm 3.42	71.31 \pm 3.51	69.82 \pm 2.19	64.09 \pm 2.87	49.71 \pm 6.09	50.01 \pm 4.87
HYPER-SAGCN	79.05 \pm 2.48	77.24 \pm 2.05	88.12 \pm 3.01	82.72 \pm 2.93	80.13 \pm 1.38	76.32 \pm 2.96	79.51 \pm 1.27	80.58 \pm 2.61	73.09 \pm 2.60	72.29 \pm 3.69
CHESHIRE	72.24 \pm 2.63	70.31 \pm 2.26	82.54 \pm 0.88	80.34 \pm 1.19	77.26 \pm 1.01	77.72 \pm 0.76	79.43 \pm 1.58	78.63 \pm 1.25	70.03 \pm 2.55	72.97 \pm 1.81
CAT-WALK	98.89 \pm 1.82	98.97 \pm 1.69	96.03 \pm 1.50	96.41 \pm 0.70	95.32 \pm 0.89	96.03 \pm 0.84	93.54 \pm 0.56	93.93 \pm 0.36	73.45 \pm 2.92	74.66 \pm 3.87
Weakly Inductive										
CE-GCN	51.80 \pm 3.29	50.94 \pm 3.77	50.33 \pm 3.40	48.54 \pm 3.92	52.19 \pm 2.54	53.21 \pm 3.59	52.38 \pm 2.75	50.81 \pm 2.68	50.81 \pm 2.87	55.38 \pm 2.79
CE-EVOLVEGCN	55.39 \pm 5.16	57.24 \pm 4.98	57.85 \pm 3.51	63.26 \pm 4.01	51.50 \pm 4.07	52.59 \pm 4.53	55.63 \pm 3.41	5.19 \pm 3.56	45.66 \pm 2.10	50.93 \pm 2.57
CE-CAW	77.61 \pm 1.05	80.03 \pm 1.65	83.77 \pm 1.41	83.41 \pm 1.19	82.98 \pm 1.06	80.84 \pm 1.57	79.51 \pm 0.94	80.39 \pm 1.07	80.54 \pm 1.02	77.41 \pm 1.28
NHP	75.17 \pm 2.02	77.23 \pm 3.11	67.25 \pm 5.19	66.73 \pm 4.94	71.92 \pm 1.83	72.30 \pm 1.89	69.58 \pm 4.07	72.48 \pm 4.83	60.38 \pm 4.45	55.62 \pm 4.67
HYPER-SAGCN	79.45 \pm 2.18	80.32 \pm 2.23	88.53 \pm 1.26	87.26 \pm 1.49	85.08 \pm 1.45	86.84 \pm 1.60	80.12 \pm 2.00	73.48 \pm 2.77	78.86 \pm 0.63	79.14 \pm 1.51
CHESHIRE	79.03 \pm 1.24	78.98 \pm 1.17	88.40 \pm 1.06	86.53 \pm 1.82	83.55 \pm 1.27	79.42 \pm 2.03	79.67 \pm 0.83	80.03 \pm 1.38	74.53 \pm 0.91	75.88 \pm 1.14
CAT-WALK	99.16 \pm 1.08	99.33 \pm 0.89	94.68 \pm 2.37	96.54 \pm 0.82	96.53 \pm 1.39	96.83 \pm 1.16	98.38 \pm 0.21	98.48 \pm 0.18	64.11 \pm 7.96	67.68 \pm 6.93
Transductive										
HPRA	70.83 \pm 0.01	67.40 \pm 0.00	94.91 \pm 0.00	89.17 \pm 0.00	89.86 \pm 0.06	88.11 \pm 0.02	79.48 \pm 0.03	77.16 \pm 0.03	78.62 \pm 0.00	76.74 \pm 0.00
HPLSF	76.19 \pm 0.82	77.62 \pm 1.42	92.14 \pm 0.29	92.79 \pm 0.15	88.57 \pm 1.09	87.69 \pm 1.61	79.31 \pm 0.52	75.88 \pm 0.05	75.88 \pm 0.05	75.32 \pm 0.08
CE-GCN	66.83 \pm 3.74	65.83 \pm 3.61	62.99 \pm 3.02	59.76 \pm 3.78	59.14 \pm 3.87	55.59 \pm 3.46	64.42 \pm 3.11	63.19 \pm 3.34	58.06 \pm 3.80	55.27 \pm 3.12
CE-EVOLVEGCN	67.08 \pm 3.51	66.51 \pm 3.80	65.19 \pm 2.26	59.27 \pm 2.19	63.15 \pm 1.32	65.18 \pm 1.89	69.30 \pm 2.27	64.38 \pm 2.66	69.98 \pm 5.38	67.76 \pm 5.16
CE-CAW	76.30 \pm 0.84	77.73 \pm 1.42	81.63 \pm 0.97	79.37 \pm 0.53	86.53 \pm 0.84	87.03 \pm 1.15	76.99 \pm 1.02	77.05 \pm 1.14	79.57 \pm 1.14	78.37 \pm 1.15
NHP	82.39 \pm 2.81	80.72 \pm 2.04	76.85 \pm 3.08	75.37 \pm 3.12	80.04 \pm 3.42	80.24 \pm 3.49	80.27 \pm 2.53	77.82 \pm 1.91	63.17 \pm 3.79	66.87 \pm 3.19
HYPER-SAGCN	80.76 \pm 2.64	80.50 \pm 2.73	94.98 \pm 1.30	89.73 \pm 1.21	90.77 \pm 2.05	88.64 \pm 2.09	82.84 \pm 1.61	81.12 \pm 1.79	83.59 \pm 0.98	80.54 \pm 1.66
CHESHIRE	84.91 \pm 1.05	82.24 \pm 1.49	95.11 \pm 0.94	94.29 \pm 1.23	91.62 \pm 1.18	92.72 \pm 1.07	86.81 \pm 1.24	83.66 \pm 1.90	82.27 \pm 0.86	81.39 \pm 0.81
CAT-WALK	98.72 \pm 1.38	98.71 \pm 1.36	95.30 \pm 0.43	95.90 \pm 0.44	97.91 \pm 3.30	97.92 \pm 2.95	88.15 \pm 1.46	88.66 \pm 1.57	80.47 \pm 5.30	82.87 \pm 3.50
Strongly Inductive										
CE-GCN	52.76 \pm 2.41	50.37 \pm 2.59	56.10 \pm 1.88	54.15 \pm 1.94	57.91 \pm 1.56	59.45 \pm 1.21	55.70 \pm 2.91	54.29 \pm 2.78	51.97 \pm 2.91	55.03 \pm 2.72
CE-EVOLVEGCN	44.16 \pm 1.27	49.15 \pm 1.23	64.08 \pm 2.75	60.64 \pm 2.78	52.00 \pm 2.32	52.69 \pm 2.15	58.17 \pm 2.24	57.35 \pm 2.13	54.57 \pm 2.25	57.16 \pm 2.55
CE-CAW	72.99 \pm 0.20	73.45 \pm 0.68	70.14 \pm 1.89	70.26 \pm 1.77	73.12 \pm 1.06	72.64 \pm 1.18	75.87 \pm 0.77	73.19 \pm 0.86	74.21 \pm 2.04	76.52 \pm 2.06
NHP	65.35 \pm 2.07	64.24 \pm 1.61	68.23 \pm 3.34	69.82 \pm 3.41	71.83 \pm 2.64	71.09 \pm 2.83	70.43 \pm 3.64	73.22 \pm 3.03	72.52 \pm 2.90	71.56 \pm 2.26
HYPER-SAGCN	78.01 \pm 1.24	80.04 \pm 1.87	73.66 \pm 1.95	73.98 \pm 1.35	73.94 \pm 2.57	72.97 \pm 2.45	75.85 \pm 2.21	73.24 \pm 2.75	78.88 \pm 2.69	77.53 \pm 2.28
CHESHIRE	69.98 \pm 2.71	70.10 \pm 3.05	N/A	N/A	76.99 \pm 2.82	74.03 \pm 2.78	76.60 \pm 2.19	74.91 \pm 2.71	75.04 \pm 3.39	75.46 \pm 2.90
CAT-WALK	91.68 \pm 2.78	91.75 \pm 2.82	88.03 \pm 3.38	88.46 \pm 3.09	89.84 \pm 6.02	91.58 \pm 4.37	93.29 \pm 1.55	94.26 \pm 1.21	97.59 \pm 2.21	97.71 \pm 2.07
Weakly Inductive										
CE-GCN	49.60 \pm 3.96	55.01 \pm 3.25	55.13 \pm 2.76	51.48 \pm 2.66	57.06 \pm 3.16	58.37 \pm 2.86	60.92 \pm 2.81	55.93 \pm 2.03	56.85 \pm 2.73	57.19 \pm 2.52
CE-EVOLVEGCN	52.44 \pm 2.38	50.61 \pm 2.32	61.79 \pm 1.63	59.61 \pm 1.12	55.81 \pm 2.54	50.63 \pm 2.46	58.48 \pm 2.49	55.90 \pm 2.51	54.10 \pm 1.21	56.13 \pm 2.32
CE-CAW	73.54 \pm 1.19	74.10 \pm 1.41	77.29 \pm 0.86	77.67 \pm 1.94	80.79 \pm 0.82	81.88 \pm 0.63	77.28 \pm 1.30	79.24 \pm 1.19	76.51 \pm 1.26	77.17 \pm 1.39
NHP	67.19 \pm 4.33	66.53 \pm 4.21	70.46 \pm 3.52	65.66 \pm 3.94	76.44 \pm 1.90	75.23 \pm 3.96	73.37 \pm 3.51	70.62 \pm 3.71	78.15 \pm 4.41	79.64 \pm 4.32
HYPER-SAGCN	77.26 \pm 2.09	74.05 \pm 2.12	78.15 \pm 1.41	76.19 \pm 1.53	75.38 \pm 1.43	70.35 \pm 1.63	80.82 \pm 2.18	76.67 \pm 2.06	74.22 \pm 1.91	70.57 \pm 1.02
CHESHIRE	77.31 \pm 0.95	76.01 \pm 0.98	N/A	N/A	81.27 \pm 0.85	82.96 \pm 1.41	80.68 \pm 1.31	80.78 \pm 1.13	77.60 \pm 1.57	79.48 \pm 1.79
CAT-WALK	91.98 \pm 2.41	92.22 \pm 2.40	90.28 \pm 2.81	90.56 \pm 2.62	97.15 \pm 1.81	97.55 \pm 1.49	95.65 \pm 1.82	96.18 \pm 1.52	98.11 \pm 1.31	98.25 \pm 1.13
Transductive										
HPRA	72.51 \pm 0.00	71.08 \pm 0.00	83.18 \pm 0.00	80.12 \pm 0.00	70.49 \pm 0.02	72.83 \pm 0.00	77.94 \pm 0.01	75.78 \pm 0.01	81.05 \pm 0.00	81.71 \pm 0.00
HPLSF	75.27 \pm 0.31	77.95 \pm 0.14	83.45 \pm 0.93	82.29 \pm 1.06	74.38 \pm 1.11	73.81 \pm 1.45	82.12 \pm 0.71	84.51 \pm 0.62	80.89 \pm 1.51	75.62 \pm 1.38
CE-GCN	64.19 \pm 2.79	65.93 \pm 2.52	55.18 \pm 5.12	55.84 \pm 4.53	62.78 \pm 2.69	59.71 \pm 2.25	63.08 \pm 2.19	65.37 \pm 2.48	66.79 \pm 2.88	60.51 \pm 2.26
CE-EVOLVEGCN	64.36 \pm 4.17	66.98 \pm 3.72	72.56 \pm 1.72	69.38 \pm 1.51	68.55 \pm 2.26	67.86 \pm 2.61	70.09 \pm 3.42	66.37 \pm 3.17	71.31 \pm 2.92	70.36 \pm 2.72
CE-CAW	78.19 \pm 1.10	77.95 \pm 0.98	81.73 \pm 2.48	83.27 \pm 2.34	80.86 \pm 0.45	80.57 \pm 1.08	84.72 \pm 1.65	84.93 \pm 1.26	80.37 \pm 1.77	83.14 \pm 0.97
NHP	78.90 \pm 4.39	76.95 \pm 5.08	79.14 \pm 3.36	78.79 \pm 3.15	82.33 \pm 1.02	81.44 \pm 1.53	81.38 \pm 1.42	82.17 \pm 1.38	78.99 \pm 4.16	80.06 \pm 4.33
HYPER-SAGCN	79.61 \pm 2.35	75.99 \pm 2.23	84.07 \pm 2.50	84.22 \pm 2.43	79.62 \pm 2.04	79.38 \pm 2.55	85.07 \pm 2.46	85.32 \pm 2.20	85.18 \pm 2.64	80.99 \pm 3.04
CHESHIRE	86.38 \pm 1.23	87.39 \pm 1.07	N/A	N/A	82.75 \pm 1.99	81.96 \pm 1.75	86.30 \pm 1.57	83.18 \pm 1.92	87.83 \pm 2.15	88.62 \pm 1.76
CAT-WALK	96.74 \pm 1.28	97.08 \pm 1.20	91.63 \pm 1.41	92.28 \pm 1.26	93.51 \pm 1.27	94.98 \pm 0.98	90.64 \pm 0.44	91.96 \pm 0.41	96.59 \pm 4.39	97.06 \pm 3.72

G.4 More Results on RNN v.s. MLP-MIXER in Walk Encoding

Most existing methods on (temporal) random walk encoding see a walk as a sequence of vertices and uses sequence encoders like RNNs or TRANSFORMERS to encode each walk. The main drawback of these methods is that they fail to directly process temporal walks with irregular gaps between timestamps. That is, sequential encoders can be seen as discrete approximations of dynamic systems; however, this discretization often fails if we have irregularly observed data [125]. This is the main motivation of recent studies to develop methods on continuous-time temporal networks [34, 126]. Most of these methods are too complicated and sometimes fail to generalize [127]. In CAT-WALK, we suggest a simple architecture to encode temporal walks by a time-encoding module along with a MIXER module (see Section 3.4 for the details). In this part, we evaluate the power of our MIXER module and compare its performance when we replace it with RNNs [82]. Figure 7 reports the results on all datasets. We observe that using MLP-MIXER with the time-encoding module in CAT-WALK can always outperform CAT-WALK when we replace MLP-MIXER with a RNN, and mostly this improvement is more on datasets with high variance in their timestamps. We relate this superiority to the importance of using continuous-time encoding instead of sequential encoders.

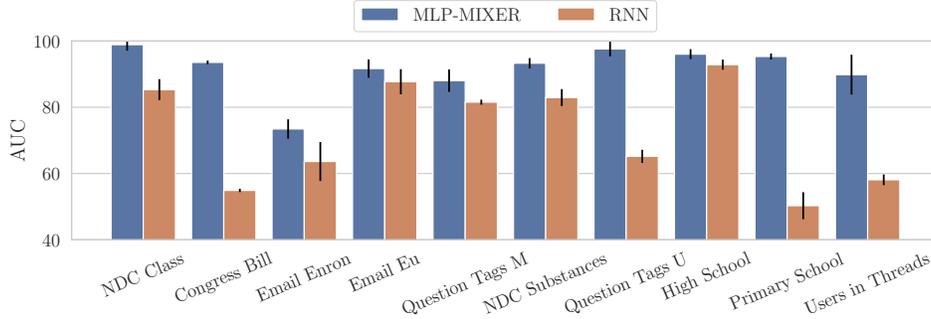


Figure 7: The importance of using MLP-MIXER in CAT-WALK. Using an RNN instead of MLP-MIXER can damage the performance in *all* datasets. RNNs are sequential encoders and are not able to encode continuous time in the data.

H Broader Impacts

Temporal hypergraph learning methods, such as CAT-WALK, benefit a wide array of real-world applications, including but not limited to social network analysis, recommender systems, brain network analysis, drug discovery, stock price prediction, and anomaly detection (e.g. bot detection in social media or abnormal human brain activity). However, there might be some potentially negative impacts, which we list as: ① Learning underlying biased patterns in the training data, which may result in stereotyped predictions. Since CAT-WALK learns underlying dynamic laws in the training data, given biased training data, the predictions of CAT-WALK can be biased. ② Also, powerful dynamic hypergraph models can be used for manipulation in the abovementioned applications (e.g., stock price manipulation). Accordingly, to prevent the potential risks in sensitive tasks, e.g., decision-making from graph-structured data in health care, interpretability and explainability of machine learning models on hypergraphs is a critical area for future work.

Furthermore, this work does not perform research on human subjects as part of the study and all used datasets are anonymized and publicly available.