

# VilLain: Self-Supervised Learning on Hypergraphs without Features via Virtual Label Propagation

Anonymous Author(s)

## ABSTRACT

Group interactions arise in various scenarios in real-world systems: collaborations of researchers, co-purchases of products, and discussions in online Q&A sites, to name a few. Such higher-order relations are naturally modeled as hypergraphs, which consist of hyperedges (i.e., any-sized subsets of nodes). For hypergraphs, the challenge to learn node representation when features or labels are not available is imminent, given that (a) most real-world hypergraphs are not equipped with external features while (b) most existing approaches for hypergraph learning resort to additional information. Thus, in this work, we propose VilLain, a novel self-supervised hypergraph representation learning method based on the propagation of virtual labels (v-labels). Specifically, we learn for each node a sparse probability distribution over v-labels as its feature vector, and we propagate the vectors to construct the final node embeddings. Inspired by higher-order label homogeneity, which we discover in real-world hypergraphs, we design novel self-supervised loss functions for the v-labels to reproduce the higher-order structure-label pattern. We demonstrate that VilLain is: (a) **Requirement-free**: learning node embeddings without relying on node labels and features, (b) **Versatile**: giving embeddings that are not specialized to specific tasks but generalizable to diverse downstream tasks, and (c) **Accurate**: more accurate than its competitors for node classification, hyperedge prediction, node clustering, and node retrieval tasks.

## ACM Reference Format:

Anonymous Author(s). 2018. VilLain: Self-Supervised Learning on Hypergraphs without Features via Virtual Label Propagation. In *Proceedings of Make sure to enter the correct conference title from your rights confirmation email (Conference acronym 'XX)*. ACM, New York, NY, USA, 18 pages. <https://doi.org/XXXXXXX.XXXXXXX>

## 1 INTRODUCTION

In many real-world complex systems, interactions often occur in groups: research collaborations, email communications, group discussions, and protein interactions, to name a few. Representing such group interactions (i.e., higher-order relationships) as edges in an ordinary pairwise graph impairs the semantics of the interactions, often leading to considerable information loss [13, 37, 76].

Hypergraphs address the limitations of ordinary graphs by modeling group interactions as hyperedges, the non-empty subsets of

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

*Conference acronym 'XX, June 03–05, 2018, Woodstock, NY*

© 2018 Association for Computing Machinery.

ACM ISBN 978-1-4503-XXXX-X/18/06...\$15.00

<https://doi.org/XXXXXXX.XXXXXXX>

nodes. Specifically, the flexibility in hyperedge sizes enables each hyperedge to naturally represent an interaction among any number of nodes. Hypergraphs have been used to model data from various fields, including bioinformatics [31], social network analysis [73], circuit design [32], and computer vision [29, 33, 66]. Notably, hypergraph modeling has demonstrated its effectiveness over ordinary graphs in diverse applications, such as recommendation [67, 68], medical prediction [6], and crime prediction [40].

A popular approach for analyzing such complex relations is to learn node embeddings (i.e., vector representations of nodes) through *self-supervision*. In the context of hypergraphs, self-supervised learning has been applied for node classification [27, 35, 64], hyperedge prediction [62, 82], recommendation [69, 79], and user location prediction in social media [73]. Self-supervised learning enjoys several key advantages. It does not require external node labels, which are scarce in many real-world scenarios due to substantial costs in their acquisition [26]. Moreover, the learned embeddings often demonstrate considerable *versatility*, maintaining their utility across a broad range of tasks.

Many self-supervised node embedding methods require external features. Hypergraph Neural Networks (HNNs) [11, 17, 20, 28, 35, 64] and Graph Neural Networks (GNNs) [24, 34, 50, 60, 61, 70, 85], for instance, heavily rely on the external node features. As such, most of them are only tested on attributed benchmark datasets [18, 25, 49, 54, 56, 74], and their performances strongly depend on the feature quality [15, 19, 41, 46].

Despite their usefulness, external features are often entirely or partially missing in real-world hypergraphs [10, 15, 19, 53, 75, 82]. In fact, only 3.03% of the graphs at a popular graph database are given with node features [54],<sup>1</sup> and none of the hypergraphs at the largest hypergraph database is attributed.<sup>2</sup> Such a problem, in combination with the issue of label scarcity, poses an imminent challenge for hypergraph representation learning.

While some self-supervised approaches do not require external features, their embeddings are hardly versatile. Some link prediction HNNs and GNNs leverage the structural or identity features [7, 62, 78, 80, 86] without the external ones, and random walk (RW) [23, 27, 51]- or matrix factorization (MF) [47, 52, 58]-based methods (i.e. Hyper2Vec) only need graph structure for their node embeddings. However, they arguably only preserve structural properties, since their input and objective functions are *solely structural*. Such models are, thus, less applicable to tasks where the importance of structural property is less prominent, such as node classification.

Thus, in this paper, we aim to learn versatile node embeddings for hypergraphs without relying on external labels or features. To this end, we propose VilLain (**V**irtual **L**abel **P**ropagation). VilLain constructs for each node a sparse probability distribution over virtual labels (v-labels) as its feature. The probabilistic v-label assignment vectors are propagated to construct the final node embeddings.

<sup>1</sup>Out of 6,659 graph datasets, 202 are given with node attributes.

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

At each propagation step, the  $v$ -labels are optimized with a novel self-supervised loss function, inspired by higher-order label homogeneity in real-world hypergraphs (see Section 4). Thus, ViLLain learns potential (higher-order) structure-label relationships, beyond purely structural properties.

Through extensive experiments using eight real-world hypergraphs and three downstream tasks (specifically, node classification, node retrieval, node clustering, and hyperedge prediction), we demonstrate the superiority of ViLLain over 15 baseline approaches. We summarize its strengths as follows:

- **Minimum Requirements:** ViLLain learns node embeddings without any supervision (e.g., node labels) or extra information (e.g., node features and the number of labels).
- **Versatile Embedding:** ViLLain learns general-purpose node embeddings that are not specialized to specific tasks but generalized to diverse downstream tasks.
- **Accurate Embedding:** ViLLain achieves up to 71.6%, 72.3%, and 6.7% better accuracy than unsupervised and (semi-)supervised baseline approaches for node classification, node retrieval, and hyperedge prediction tasks, respectively.

**Reproducibility.** Our code and dataset are available at <https://anonymous.4open.science/r/ViLLain-C18B> (anonymous).

## 2 RELATED WORK

In this section, we briefly review related works on node representation learning, focusing on learning without labels or features.

**Node embedding with propagation.** Propagation has been widely applied and shown effective for both hypergraph and graph representation learning. GNNs typically have each node propagate its features to the direct neighbors [9, 22, 34], whereas for HNNs, the propagation is conducted on hypergraph structure. Specifically, HGNN [20] has each node propagate to its hyperedges, where the node feature are aggregated and propagated back to the nodes that belong to the hyperedges. HNHN [17] uses non-linear aggregation functions to update both node and hyperedge embeddings, alternatingly. AllSet [11] uses permutation-invariant functions to propagate on hyperedges. Other simplified GNNs [12, 16, 21, 65] first learn soft label vectors from feature vectors, which are propagated to learn the final node embeddings. Note that all the described methods require external labels or features.

**Node embedding without external labels.** Self-supervision has been widely adopted for representation learning without external labels. Self-supervised HNNs and GNNs often utilize contrastive losses. Given both original and perturbed features or structures, the models maximize the mutual information between them [35, 61, 85]. For hypergraphs, HyperGCL [64] uses node- and hyperedge-level perturbation, and TriCL [35] conducts tri-directional contrasts that maximize the agreement between two augmented views of nodes, groups, and memberships. Intuitively, such self-supervised loss functions are designed to learn node embeddings that denoise the input features and structure. It, then, implies that these self-supervised models can only learn structural properties if their input node features are random or structural.

Given random walk sequences, RW-based embedding methods [23, 27, 51] typically use Skip-Gram [44] to optimize the embeddings to maximize the likelihood of the visited nodes. MF-based approaches [47, 52, 58], on the other hand, factorize proximity matrices into low-rank matrices. As such, most RW- and MF-based embedding methods specifically preserve structural proximity.

**Node embedding without external features.** If external features are not available, HNNs and GNNs require derived features for their prediction. For structural prediction, some models have leveraged only structural information as the input features [7, 15, 62, 78, 80, 86]. Specifically, structural [4, 7, 15, 78], positional [15, 39, 63], and identity [1, 55, 62, 77, 80, 81, 86] encoding methods have been developed. Such encoding methods generally aim to enhance model expressivity beyond 1-WL test [70]. On the other hand, the majority of RW- and MF-based approaches do not require any features or labels [23, 27, 47, 51, 52, 58]. It is, however, worth noting that all the described methods *over-emphasize structural properties*, since their features and objective loss functions are solely structural. Thus, predictions from their embeddings hardly generalize to less structure-dependent tasks, such as node classification.<sup>3</sup>

**Relating ViLLain to the prior works.** In comparison to (hyper) graph learning models without external features or labels, we present the novelty of ViLLain in the subsequent sections as follows:

- **Novel Self-Supervised Loss:** Only ViLLain has loss function that learns *beyond structural information* for embedding versatility.
- **Novel Input Feature Learning:** ViLLain’s motivation and mechanism of input feature learning are distinguished from the prior methods.

## 3 PROBLEM STATEMENT

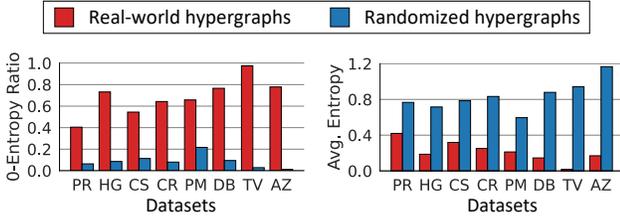
In this section, we formulate hypergraph representation learning without features or labels. A hypergraph  $G = (V, E)$  consists of a set of nodes  $V = \{v_1, \dots, v_{|V|}\}$  and a set of hyperedges  $E = \{e_1, \dots, e_{|E|}\}$ . Each hyperedge  $e_j \in E$  is a non-empty subset of nodes, i.e.,  $\emptyset \subsetneq e_j \subseteq V$ . In the incidence matrix  $\mathbf{H} \in \{0, 1\}^{|V| \times |E|}$  of  $G$ ,  $\mathbf{H}_{ij} = 1$ , if  $v_i \in e_j$ , and  $\mathbf{H}_{ij} = 0$  otherwise.

Given a hypergraph  $G = (V, E)$ , the objective of self-supervised hypergraph representation learning is to learn a node embedding  $\mathbf{Z}_i \in \mathbb{R}^d$  of each node  $v_i \in V$ , or equivalently, a node embedding matrix  $\mathbf{Z} \in \mathbb{R}^{|V| \times d}$  that captures meaningful proximity between nodes in  $G$ . Specifically, we aim to learn node embeddings that are generally useful for various tasks (e.g., node classification and hyperedge prediction), without relying on any kind of supervision (e.g., ground-truth semantic labels or even the number of unique labels) or external information (e.g., node attributes).

## 4 MOTIVATING OBSERVATIONS

In this section, we present our observation in real-world hypergraphs, which motivate the design of ViLLain in Section 5. Inspired by pervasive *homophily* [2, 43] in real-world graphs, we postulate that hypergraphs also exhibit a similar tendency. For example, researchers from the same area tend to co-author a paper, and e-mails are likely to be exchanged within the same department. To substantiate this hypothesis, we examine label homogeneity in eight different real-world hypergraphs.

<sup>3</sup>See the low performances of such methods (e.g. Hyper2Vec, HyperGCL) in Table 2.



**Figure 1: Hyperedges in real-world hypergraphs (statistics in Appendix B) exhibit label homogeneity (Obs. 1).**

Using the ground-truth node labels, for each hyperedge, we measure the entropy of its soft label assignment vector, which is obtained by averaging the label assignment one-hot vectors of the nodes in the hyperedge. If the entropy is 0, all nodes in the hyperedge are labeled identically (high homogeneity). The higher the entropy is, the more diverse labels the nodes in the hyperedge have (low homogeneity). As shown in Figure 1, the entropy in real-world hypergraphs tends to be lower than that in hypergraphs that are randomized as described in [36]. Moreover, the ratio of the hyperedges with entropy 0 is much higher in real-world hypergraphs than in the randomized hypergraphs, and the average entropy is lower in real-world hypergraphs than in the randomized hypergraphs.

**OBSERVATION 1.** *Hyperedges in real-world hypergraphs exhibit label homogeneity, i.e., they tend to contain the same labeled nodes.*

In addition, we examine higher-order homogeneity in real-world hypergraphs. To this end, we measure the entropy of the higher-order label assignment vectors (or  $\ell$ -step labels in short) of hyperedges. For each  $\ell \geq 0$ , the  $\ell$ -step label of a hyperedge is obtained by averaging the  $\ell$ -step labels of the nodes in it. The  $\ell$ -step label of each node is given if  $\ell = 0$ , or obtained by averaging  $(\ell - 1)$ -step labels of the incident hyperedges (the detailed procedure can be found in Section 5.1). Figure 2 demonstrates that (a) the entropy of 50-step labels of hyperedges in a real-world hypergraph (spec., Trivago) is lower than those in the randomized counterpart, and (b) regardless of the step count  $\ell$ , hyperedges in the real-world hypergraph exhibit higher homogeneity than those in the randomized hypergraph. These findings provide concrete evidence supporting the presence of higher-order homogeneity in real-world hypergraphs. Refer to Appendix C for results from other real-world hypergraphs.

**OBSERVATION 2.** *Real-world hypergraphs exhibit higher-order label homogeneity, i.e., the node labels in each hyperedge tend to be homogeneous even after multiple steps of propagation.*

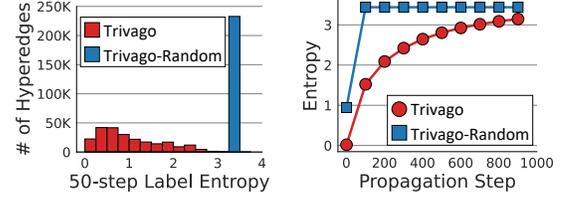
## 5 PROPOSED METHOD

In this section, we propose VilLain (Figure 3), a self-supervised node representation learning method for hypergraphs. Notably, VilLain does not require external labels or features.

### 5.1 VilLain: Virtual Label Propagation

We first present how VilLain obtains node embeddings through virtual label (v-labels) propagation, *without external features*.

**Virtual Labels.** Since node labels or features are not given, VilLain assumes the presence of  $d$  v-labels and leverages the soft v-label assignment vector of each node as its learnable feature. Specifically,



**Figure 2: Real-world hypergraphs exhibit higher-order label homogeneity (Obs. 2).**

VilLain employs a learnable matrix  $\tilde{\mathbf{X}} \in \mathbb{R}^{|V| \times d}$  where each  $i^{\text{th}}$  row  $\tilde{\mathbf{X}}_i$  is used to obtain the soft assignment vector  $\mathbf{X}_i^{(0)} \in [0, 1]^d$  of the node  $v_i$  to  $d$  v-labels as follows:

$$\mathbf{X}_{ij}^{(0)} = \frac{e^{(\tilde{\mathbf{X}}_{ij} + g_j)}}{\sum_{j'=1}^d e^{(\tilde{\mathbf{X}}_{ij'} + g_{j'})}}, \quad \text{for } j = 1, \dots, d, \quad (1)$$

where  $g_j = -\log(\log(\frac{1}{u_i}))$  is random noise and  $u_i \sim \text{Uniform}(0, 1)$ . The above equation transforms the vector into a probability vector and encourages it to be biased towards a single v-label. As described later, the v-label assignment vectors are optimized to reproduce higher-order label homogeneity (Observations 1 and 2).

**Hypergraph V-label Propagation.** After obtaining the v-label matrix  $\mathbf{X}^{(0)}$ , VilLain conducts v-label propagation on the input hypergraph to obtain  $\mathbf{X}^{(\ell)}$ . At each step, v-labels are propagated alternately between nodes and hyperedges. Specifically, the v-label assignment matrices of hyperedges and nodes at step  $\ell$  are:

$$\mathbf{Y}^{(\ell)} = \mathbf{D}_E^{-1} \mathbf{H}^T \mathbf{X}^{(\ell-1)} \quad \text{and} \quad \mathbf{X}^{(\ell)} = \mathbf{D}_V^{-1} \mathbf{H} \mathbf{Y}^{(\ell)}, \quad (2)$$

where  $\mathbf{D}_V$  and  $\mathbf{D}_E$  are the diagonal matrices with node degrees and hyperedge sizes, respectively. To capture higher-order dependencies among nodes, VilLain computes node embeddings  $\mathbf{Z} \in [0, 1]^{|V| \times d}$  by averaging the v-label assignment vectors obtained at propagation steps  $1, \dots, k'$ :

$$\mathbf{Z} = \frac{1}{k'} \sum_{\ell=1}^{k'} \mathbf{X}^{(\ell)}. \quad (3)$$

Namely, the embedding  $\mathbf{Z}_i$  of node  $v_i$  is a probability vector averaging its v-label assignment vector at each step.

**Multi-V-label Propagation.** In real-world hypergraphs, nodes may have multiple labels, each representing different aspects. For instance, in a social network, socioeconomic status and political inclination can both serve as labels, albeit their independent homogeneity w.r.t. hypergraph topology. The same goes for the number of labels. Learning a single set of v-labels, then, can be *insufficient* to capture their complex structure-label patterns.

Thus, VilLain learns multi-v-labels for the final node embedding  $\mathbf{Z}^*$ . Specifically, we partition the  $d$ -dimensional embedding space into  $D$  subspaces of potentially different dimensions, allowing for independent v-label propagation within each subspace. Then, VilLain concatenates the outputs from each subspace as follows:

$$\mathbf{Z}_i^* = \left[ \mathbf{Z}_i^{(1)} \parallel \mathbf{Z}_i^{(2)} \parallel \dots \parallel \mathbf{Z}_i^{(D)} \right],$$

where  $\parallel$  is the concatenation operation, and  $\mathbf{Z}_i^{(\cdot)}$  is the embedding obtained from each subspace.

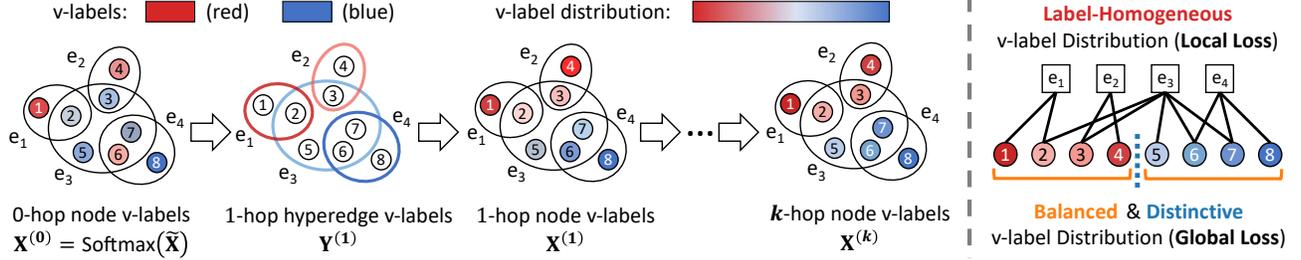


Figure 3: (Left) Two v-labels (red and blue) are propagated between nodes and hyperedges on a hypergraph (Sec. 5.1). Note that hyperedges are colored to indicate  $Y^{(1)}$ . (Right) By minimizing the proposed local and global losses, the v-label distributions are learned to exhibit higher-order label homogeneity while being balanced and distinctive at each propagation step (Sec. 5.2).

## 5.2 Self-Supervision Objectives

The learning objectives of ViLLain are designed to reproduce higher-order label homogeneity by effectively capturing structural properties and also potential higher-order structure-label relationships. Recall that the entries of the matrix  $\bar{X}$  are the only learnable parameters in ViLLain that the objective function updates.

**Capturing Local Information.** Motivated by Observations 1 and 2 in Section 4, we design an objective to capture the higher-order homogeneity of nodes and hyperedges. Specifically, ViLLain minimizes the entropy of the v-label assignment vectors of each node and hyperedge obtained at propagation steps  $1, \dots, k$ :

$$\mathcal{L}_{\text{local}} = \sum_{\ell=1}^k \left( \frac{1}{|V|} \sum_{i=1}^{|V|} \mathcal{E}(\mathbf{x}_i^{(\ell)}) + \frac{1}{|E|} \sum_{j=1}^{|E|} \mathcal{E}(\mathbf{y}_j^{(\ell)}) \right), \quad (4)$$

where  $\mathcal{E}(p) = -\sum_i p_i \log p_i$  is the entropy measure of  $p$ . That is, we induce structurally close nodes (or hyperedges) to be assigned to the same v-label. Beyond capturing the homogeneity at the hyperedge level, i.e.,  $\ell = 1$  (Observation 1), the loss function is designed to reproduce the higher-order homogeneity of nodes and hyperedges by minimizing the entropy of v-label assignment vectors at each propagation step  $\ell \in \{1, \dots, k\}$  (Observation 2). For training speed, the number of steps  $k$  for training can be smaller than  $k'$  for inference.

**Capturing Global Information.** ViLLain also considers the global distribution of labels. To this end, we give v-label-level supervision to ViLLain so that v-labels are properly distributed over the entire hypergraph. First, since Eq. (4) is trivially minimized when all nodes and hyperedges are assigned to a single v-label, we use the following term to prevent this problem:

$$\mathcal{J}_{\text{cls}} = - \sum_{\ell=1}^k \left( \mathcal{E}(\mathbf{x}^{(\ell)}) + \mathcal{E}(\mathbf{y}^{(\ell)}) \right) \quad (5)$$

$$\text{where } \mathbf{x}_i^{(\ell)} = \frac{\|\mathbf{X}_{:,i}^{(\ell)}\|_1}{\sum_{j=1}^d \|\mathbf{X}_{:,j}^{(\ell)}\|_1} \text{ and } \mathbf{y}_j^{(\ell)} = \frac{\|\mathbf{Y}_{:,j}^{(\ell)}\|_1}{\sum_{i=1}^d \|\mathbf{Y}_{:,i}^{(\ell)}\|_1}.$$

Here,  $\mathbf{x}^{(\ell)} = [\mathbf{x}_1^{(\ell)}, \dots, \mathbf{x}_d^{(\ell)}]$  and  $\mathbf{y}^{(\ell)} = [\mathbf{y}_1^{(\ell)}, \dots, \mathbf{y}_d^{(\ell)}]$  denote the weighted ratios of nodes and hyperedges for each v-label at step  $\ell$ . Note that  $\mathbf{X}_{:,i}^{(\ell)}$  and  $\mathbf{Y}_{:,j}^{(\ell)}$ , which are the  $i$ th columns of  $\mathbf{X}^{(\ell)}$  and  $\mathbf{Y}^{(\ell)}$ , correspond to the vectors of v-label  $i$  for nodes and hyperedges,

respectively. That is, we maximize the *entropy of the global distribution* of the v-labels at each step, restraining any single v-label from dominating the entire hypergraph.

In addition, we aim to make v-labels distinctive by making the sets of nodes and hyperedges assigned to each v-label nearly disjoint from those with another v-label. To this end, we minimize the following cross-entropy-based objective:

$$\mathcal{J}_{\text{dst}} = - \sum_{\ell=1}^k \sum_{i=1}^d \left( \log \bar{\mathbf{x}}_i^{(\ell)} + \log \bar{\mathbf{y}}_i^{(\ell)} \right) \quad (6)$$

$$\text{where } \bar{\mathbf{x}}_i^{(\ell)} = \frac{e^{\mathcal{S}(\mathbf{x}_{:,i}^{(\ell)}, \mathbf{x}_{:,i}^{(\ell)})}}{\sum_{j=1}^d e^{\mathcal{S}(\mathbf{x}_{:,i}^{(\ell)}, \mathbf{x}_{:,j}^{(\ell)})}} \text{ and } \bar{\mathbf{y}}_i^{(\ell)} = \frac{e^{\mathcal{S}(\mathbf{y}_{:,i}^{(\ell)}, \mathbf{y}_{:,i}^{(\ell)})}}{\sum_{j=1}^d e^{\mathcal{S}(\mathbf{y}_{:,i}^{(\ell)}, \mathbf{y}_{:,j}^{(\ell)})}}.$$

Here,  $\bar{\mathbf{x}}_i^{(\ell)}$  and  $\bar{\mathbf{y}}_j^{(\ell)}$  indicate the distinctiveness of v-label  $i$  at each propagation step  $\ell$  in nodes and hyperedges, respectively, and  $\mathcal{S}(\cdot, \cdot)$  measures the cosine similarity of two input vectors. Minimizing Eq. (6) reinforces the *distinctiveness of each v-label* at each step.

Finally, we minimize the global-level loss, defined as the sum of Eq. (5) and Eq. (6), to let v-labels be properly distributed across the entire hypergraph:

$$\mathcal{L}_{\text{global}} = \mathcal{J}_{\text{cls}} + \mathcal{J}_{\text{dst}} \quad (7)$$

**Objective Function.** To exhibit both the local and global structure-label patterns, ViLLain minimizes both objectives, Eq. (4) and (7):

$$\mathcal{L} = \mathcal{L}_{\text{local}} + \mathcal{L}_{\text{global}}.$$

While we can introduce a hyperparameter for balancing  $\mathcal{L}_{\text{local}}$  and  $\mathcal{L}_{\text{global}}$ , we simply add the two losses since hyperparameter tuning based on external supervision is strictly restricted in our setting.

Note that, by reproducing higher-order label homogeneity, ViLLain captures not only structural properties but also potential higher-order structure-label relationships. Consequently, compared to self-supervised methods that exclusively focus on structural aspects (see Section 2 for further discussion), ViLLain can learn effective embeddings for less structure-dependent tasks, such as node classification, as confirmed empirically in Section 6.

**Complexity Analysis.** We analyze the time and space complexity of ViLLain for computing the final embedding  $Z^*$ , as well as the computational cost associated with optimizing their losses. Specifically,

when the dimension of each subspace is  $d/D$ , it takes:

$$O\left(kd \sum_{e \in E} |e| + \frac{kd^2}{D} (|V| + |E|)\right)_{\text{time}} \text{ and } O\left(\frac{kd^2}{D} (|V| + |E|)\right)_{\text{space}}$$

for propagating v-labels (Eq. (2)) and computing losses  $\mathcal{L}_{\text{local}}$  (Eq. (4)),  $\mathcal{J}_{\text{cls}}$  (Eq. (5)), and  $\mathcal{J}_{\text{dst}}$  (Eq. (6)) for  $1, \dots, k$  steps. To generate node embeddings (Eq. (3)), the losses are not necessarily computed, and thus it takes  $O(k'd \sum_{e \in E} |e|)$  time and requires  $O(k'd(|V| + |E|))$  space. The details can be found in Appendix A. Importantly, introducing multi-v-labels (i.e.,  $D > 1$ ) leads to the reduction in time and space complexity, thereby indicating an additional advantage of learning v-labels in multiple subspaces. This is empirically supported in Section 6.4.

### 5.3 Extension to Unobserved Nodes

Heretofore, we described how VilLain learns node embeddings  $\mathbf{Z}$  from a static hypergraph. However, in many scenarios, hypergraphs evolve over time (e.g., new members in the group), introducing new nodes and hyperedges to the hypergraph. This motivates us to extend VilLain to generate embeddings also for newly introduced, unobserved nodes and hyperedges. In this subsection, we extend VilLain to embed such unobserved nodes and hyperedges.

**Settings.** Consider a connected hypergraph  $G_S = (V_S, E_S)$ , which is a subset of a connected hypergraph  $G = (V, E)$ , where  $V_S \subseteq V$  and  $E_S \subseteq E$ . Using the incidence matrix  $\mathbf{H}_S \in \{0, 1\}^{|V_S| \times |E_S|}$  of  $G_S$ , VilLain has generated v-labels and embeddings  $\mathbf{X}_S^{(0)}, \mathbf{Z}_S \in \mathbb{R}^{|V_S| \times d}$ , respectively, for the observed nodes  $V_S$ . Nodes  $V \setminus V_S$  and hyperedges  $E \setminus E_S$  are introduced after VilLain training.

**Embedding Unobserved Nodes.** To embed nodes including the unobserved ones  $V \setminus V_S$ , VilLain propagates learned v-labels  $\mathbf{X}_S^{(0)}$  of the observed nodes  $V_S$  on hypergraph  $G$  containing the unobserved nodes and hyperedges. Specifically, v-label assignment matrices for all nodes  $\mathbf{X}^{(\ell)}$  and hyperedges  $\mathbf{Y}^{(\ell)}$  at step  $\ell \geq 1$  are obtained like in Eq. (2) as follows:

$$\mathbf{Y}^{(\ell)} = \mathbf{D}_E^{-1} \mathbf{H}^T \mathbf{X}^{(\ell-1)} \quad \text{and} \quad \mathbf{X}^{(\ell)} = \mathbf{D}_V^{-1} \mathbf{H} \mathbf{Y}^{(\ell)},$$

where  $\mathbf{X}^{(0)} \in \mathbb{R}^{|V| \times d}$  is  $\mathbf{X}_S^{(0)}$  with zero-paddings at row indices of the nodes  $V \setminus V_S$ . Since we assume a connected hypergraph  $G$ , there always exists  $\ell'$  such that all nodes  $V$  are assigned non-zero v-labels. Then, using Eq. (3),  $\mathbf{X}^{(\ell')}, \dots, \mathbf{X}^{(k')}$  are used to generate embeddings  $\mathbf{Z}$  for all nodes  $V$ , where  $k' \geq \ell'$ . We empirically show that VilLain generates informative embeddings for unobserved nodes in Section 6.4.

## 6 EXPERIMENTAL RESULTS

In this section, we present experimental results for four downstream tasks utilizing node embeddings. We first assess the accuracy of VilLain by comparing it with the state-of-the-art (hyper)graph representation learning methods (Section 6.2). Then, we demonstrate the effectiveness of each design choice of VilLain (Section 6.3). Lastly, we conduct additional analyses on VilLain (Section 6.4).

### 6.1 Experimental Settings

In this subsection, we report the experimental settings.

**Table 1: Summary statistics of eight real-world hypergraphs: the number of nodes  $|V|$ , the number of hyperedges  $|E|$ , the size of the hypergraph  $\sum_{e \in E} |e|$ , the number of edges  $|\mathcal{E}|$  in the clique expansion, and the number of ground-truth labels.**

Dataset	$ V $	$ E $	$\sum_{e \in E}  e $	$ \mathcal{E} $	# Labels
Primary (PR) [57]	242	12,704	30,729	8,317	11
High (HG) [42]	327	7,818	18,192	5,818	9
Citeseer (CS) [71]	1,019	819	2,808	3,867	6
Cora (CR) [71]	1,330	1,503	4,599	4,144	7
Pubmed (PM) [71]	3,824	7,951	34,605	123,819	3
DBLP (DB) [71]	36,188	18,924	90,868	425,669	6
Trivago (TV) [14]	172,738	233,202	726,861	1,095,204	160
Amazon (AZ) [45]	260,209	31,964	422,076	14,142,811	10

**Datasets.** We use eight publicly available real-world hypergraphs summarized in Table 1. All datasets are derived from group interactions that arise in real-world scenarios (e.g., coauthorship and co-purchase). For details regarding the preprocessing method and descriptions for each dataset, refer to Appendix B.1.

**Baselines.** We consider 15 unsupervised and (semi-)supervised graph and hypergraph embedding methods as competitors. Deepwalk [51], Node2vec [23], DGI [61], GRACE [85], GMI [50], Hyper2vec [27], LBSN [73], and TriCL [35] are unsupervised methods, and GCN [34], GAT [60], HGNN [20], HNHN [17], AllSet [11], UniGNN [28], and HyperGCL [64] are (semi-)supervised methods. For graph embedding methods (i.e., GCN, GAT, Deepwalk, Node2vec, DGI, GRACE, and GMI), we use the clique expansion of the hypergraph.<sup>4</sup> For all methods that require node features (i.e., GCN, GAT, DGI, GRACE, GMI, HGNN, HNHN, AllSet, UniGNN, HyperGCL, and TriCL), we use the embeddings obtained by Hyper2vec,<sup>5</sup> which lead to the best performance among three alternatives (see Appendix C for detailed results).

**Implementations.** We simply use  $k = 4$  for VilLain and all its variants and use  $k' = 10$  for small datasets (s.t.,  $|V| < 10,000$ ) and  $k' = 100$  for large datasets (s.t.,  $|V| \geq 10,000$ ). As discussed in Section 5.1, to capture diverse structural-label information, we aggregate embeddings obtained with various numbers of v-labels. Specifically, we concatenate embeddings obtained using different numbers of v-labels. For each number  $\lceil \frac{d}{D} \rceil \in \{2, 3, \dots, 8\}$  of v-labels, we learn  $D$  subspaces and then perform PCA to ensure that the final embedding is of the target dimension  $d$ . Refer to Appendix B.2 for the detailed settings of other baselines.

### 6.2 Accuracy of VilLain

To verify the quality of the VilLain's node embeddings, we consider four downstream tasks on hypergraphs: node classification, node retrieval, node clustering, and hyperedge prediction. The embedding dimension of all methods, including VilLain, is fixed to 128. Results including standard deviation is provided in Appendix C.2.

**Node Classification.** We perform node classification by randomly and disjointly splitting the dataset into training, validation, and test sets. For training and validation sets, the labels of 20 nodes per class are given for all datasets except for Primary and High, where the labels of 2 nodes are given per class. The remaining

<sup>4</sup>The clique expansion is the pairwise graph obtained by replacing each hyperedge with the clique formed by the nodes in the hyperedge.

<sup>5</sup>For Amazon, we use Node2vec since Hyper2vec ran out of time ( $> 24$  hours).

**Table 2: VilLain performs best on node classification in terms of accuracy. Each baseline method is designed for either graphs or hypergraphs and for either semi-supervised or unsupervised settings.**

Method	DBLP	Trivago	Amazon	Primary	High	Citeseer	Cora	Pubmed	Rank
GCN	67.37 ± 1.45	38.06 ± 1.49	28.73 ± 4.73	75.63 ± 5.08	96.25 ± 2.55	60.64 ± 3.47	72.96 ± 1.82	77.56 ± 2.58	7.00 ± 2.91
GAT	61.74 ± 1.97	51.52 ± 0.68	30.94 ± 2.13	66.79 ± 4.73	90.58 ± 2.76	49.57 ± 2.64	58.09 ± 2.14	73.67 ± 1.78	11.75 ± 3.83
Deepwalk	29.03 ± 1.43	16.85 ± 0.45	25.43 ± 1.72	84.89 ± 3.67	99.31 ± 0.48	45.10 ± 3.18	56.58 ± 1.88	68.58 ± 2.60	11.62 ± 4.71
Node2vec	29.21 ± 1.89	16.88 ± 0.44	25.27 ± 2.36	83.53 ± 3.09	99.38 ± 0.45	45.37 ± 3.17	59.15 ± 1.84	69.05 ± 3.00	11.00 ± 4.35
DGI	62.37 ± 3.32	73.46 ± 1.22	31.80 ± 1.45	86.66 ± 4.51	92.49 ± 0.60	61.36 ± 2.91	71.23 ± 2.04	77.51 ± 1.38	7.25 ± 3.59
GRACE	71.86 ± 2.51	OOM	OOM	63.78 ± 5.12	99.03 ± 0.30	61.16 ± 2.78	73.43 ± 1.81	77.70 ± 1.81	5.50 ± 4.75
GMI	64.19 ± 1.63	OOM	OOM	80.10 ± 4.94	96.61 ± 2.63	58.67 ± 2.68	71.31 ± 1.69	75.51 ± 2.77	9.16 ± 1.57
HGNN	66.60 ± 2.18	OOM	OOM	88.28 ± 5.02	92.19 ± 3.84	60.91 ± 2.32	72.90 ± 2.00	76.58 ± 2.86	7.50 ± 3.09
HNNH	63.99 ± 2.21	59.52 ± 1.64	28.99 ± 2.63	91.31 ± 2.47	96.83 ± 1.25	59.02 ± 1.63	68.81 ± 1.26	75.33 ± 1.77	7.50 ± 2.39
AllSet	63.67 ± 1.89	36.58 ± 0.93	21.75 ± 1.67	85.94 ± 3.02	95.70 ± 1.66	56.08 ± 1.95	67.73 ± 1.81	74.11 ± 2.04	10.75 ± 1.08
UniGNN	67.16 ± 2.15	69.98 ± 1.60	33.77 ± 3.22	88.88 ± 3.58	95.12 ± 3.97	59.10 ± 2.76	71.44 ± 1.03	74.37 ± 2.10	7.12 ± 3.09
HyperGCL	58.72 ± 1.54	74.99 ± 1.23	22.86 ± 2.01	74.07 ± 6.06	85.79 ± 8.92	57.54 ± 1.61	74.99 ± 1.33	78.44 ± 3.33	8.87 ± 5.18
Hyper2vec	67.18 ± 1.78	75.82 ± 1.45	OOT	92.52 ± 2.45	96.34 ± 1.34	61.50 ± 2.60	71.79 ± 1.63	77.04 ± 1.51	4.85 ± 2.35
LBSN	22.63 ± 2.20	47.99 ± 0.82	11.56 ± 0.90	86.71 ± 3.71	95.87 ± 2.28	45.43 ± 2.15	59.70 ± 1.31	54.89 ± 2.38	11.87 ± 3.21
TriCL	68.18 ± 1.36	OOM	OOM	92.67 ± 2.50	98.10 ± 1.02	59.17 ± 3.35	72.35 ± 1.53	78.57 ± 1.88	4.16 ± 1.95
<b>VilLain</b>	<b>77.16 ± 1.26</b>	<b>79.43 ± 1.63</b>	<b>57.95 ± 2.47</b>	<b>93.66 ± 3.93</b>	<b>99.19 ± 0.41</b>	<b>61.53 ± 3.17</b>	<b>75.03 ± 1.38</b>	<b>78.82 ± 1.47</b>	<b>1.25 ± 0.66</b>

**Table 3: VilLain performs overall best on hyperedge prediction (in terms of accuracy), node clustering (in terms of normalized mutual information), and node retrieval (in terms of mean average precision).**

Method	Hyperedge Prediction (Acc.)									Node Clustering (NMI)									Node Retrieval (MAP)								
	DB	TV	AZ	PR	HG	CS	CR	PM	Rank	DB	TV	AZ	PR	HG	CS	CR	PM	Rank	DB	TV	AZ	PR	HG	CS	CR	PM	Rank
Deepwalk	63.9	61.3	69.4	83.8	85.9	69.6	67.2	65.9	6.25	0.7	16.7	7.8	85.2	100.0	14.6	23.9	34.4	5.00	21.3	7.5	27.7	81.6	98.7	27.6	29.2	49.0	6.37
Node2vec	64.2	61.4	69.3	83.2	85.4	70.4	66.9	65.8	6.62	0.9	17.0	7.7	83.5	100.0	14.5	23.8	32.8	5.87	21.6	7.1	27.8	81.1	98.6	27.3	29.4	49.4	6.50
DGI	86.1	83.8	90.8	79.1	84.4	79.2	76.3	80.9	3.75	16.6	44.5	13.0	84.4	73.9	29.1	32.1	31.3	5.62	36.1	37.3	31.1	89.7	97.8	43.8	50.6	61.7	3.25
GRACE	85.4	OOM	OOM	80.3	87.4	77.9	74.5	79.1	4.00	43.0	OOM	OOM	67.6	98.2	33.0	46.0	31.6	5.00	50.2	OOM	OOM	61.4	99.5	41.1	54.2	60.9	4.00
GMI	75.6	OOM	OOM	82.4	85.9	74.4	69.4	72.3	5.50	27.8	OOM	OOM	84.1	93.1	25.3	42.6	18.7	6.50	34.6	OOM	OOM	80.0	97.8	35.9	41.6	55.1	6.33
Hyper2vec	71.2	72.4	OOT	76.4	79.6	78.1	71.7	71.5	6.14	43.4	66.3	OOT	92.5	99.3	34.3	45.5	33.6	2.27	35.5	43.1	OOT	85.7	90.7	41.2	46.7	55.6	4.85
LBSN	48.7	89.1	63.7	79.4	87.1	74.3	69.6	66.1	5.87	1.1	39.4	2.7	85.5	97.8	12.1	29.0	4.6	6.50	21.0	19.1	29.1	81.3	93.2	30.6	40.1	43.5	6.62
TriCL	77.4	OOM	OOM	84.0	87.8	82.0	76.7	80.5	2.33	38.0	OOM	OOM	87.8	98.7	34.4	44.8	33.7	3.00	45.1	OOM	OOM	89.9	97.6	42.4	55.0	61.9	3.16
<b>VilLain</b>	<b>81.6</b>	<b>95.1</b>	<b>94.9</b>	<b>83.2</b>	<b>87.8</b>	<b>82.1</b>	<b>79.0</b>	<b>82.8</b>	<b>1.50</b>	<b>46.6</b>	<b>69.4</b>	<b>35.2</b>	<b>85.7</b>	<b>98.7</b>	<b>34.5</b>	<b>50.4</b>	<b>32.7</b>	<b>2.25</b>	<b>60.2</b>	<b>67.2</b>	<b>53.6</b>	<b>91.3</b>	<b>99.0</b>	<b>46.4</b>	<b>58.0</b>	<b>64.4</b>	<b>1.12</b>

nodes are used as the test set. For un- or self-supervised methods including VilLain, we evaluate the accuracy of logistic regression using the embeddings obtained from each method. Table 2 shows the accuracy of all methods in all datasets. VilLain ranks first on average, showing the best performance. We conjecture that v-label propagation inherits rich structural properties and also potential higher-order structure-label relationships, generating high-quality representations of nodes.

**Hyperedge Prediction.** The problem of hyperedge prediction is formulated as a binary classification task, predicting whether the given hyperedge is real or fake [30, 48, 76]. Given a set  $E$  of real hyperedges, we generate a set  $E'$  of fake hyperedges with the same hyperedge size distribution by randomly sampling subsets of nodes. To obtain the embedding of each hyperedge, we apply maxmin pooling<sup>6</sup> to the embeddings of the nodes in it. For more training details on hyperedge prediction, refer to Appendix B.4. As shown in Table 3, VilLain performs the best on average. We conjecture that VilLain, which captures potential structure-label relations, is effective for this task because it indirectly relates to labels due to the high label homogeneity of real hyperedges.

**Node Clustering.** For the clustering task, we group nodes into the number of unique ground-truth labels, applying k-means to the

<sup>6</sup>We compute maxmin pooling by: elementwise max pooling - elementwise min pooling. An alternative pooling method is compared in Appendix C.8.

learned embeddings. Then, we compute the Normalized Mutual Information (NMI) to assess the quality of clustering. As shown in Table 3, VilLain outperforms all baseline methods in terms of average ranks. This indicates that the embeddings learned by VilLain exhibit meaningful semantic similarities in their distribution.

**Node Retrieval.** The problem of node retrieval aims to search for similar nodes of a given query node, using the learned embeddings. Specifically, we retrieve nodes based on the cosine similarity between their embeddings and the embedding of the query node. Then, we compute the Mean Average Precision (MAP), to measure the retrieval quality. Intuitively, the retrieval is considered to be successful if the nodes of the same class as the query node are highly ranked. For more details regarding the task, refer to Appendix B.3. As shown in Table 3, VilLain outperforms baseline methods, with a large margin. These results imply that v-labels, which are *virtual* and learned without any ground-truth node labels, are useful for finding similar nodes of the same class.

### 6.3 Ablation Study

In this subsection, we conduct ablation studies to verify the effectiveness of each component of VilLain by comparing its performance to that of its variants.

**Table 4: VilLain outperforms its three variants, VilLain-S, VilLain-M, and VilLain-L, in four downstream tasks, implying that VilLain benefits from (1) propagating v-labels in multiple subspaces, (2) aggregating embeddings from various numbers of v-labels, and (3) reproducing both local and global structure-label patterns for self-supervision.**

Method	Node Classification (Accuracy)								Hyperedge Prediction (Accuracy)								Node Clustering (NMI)								Node Retrieval (MAP)							
	DB	TV	AZ	PR	HG	CS	CR	PM	DB	TV	AZ	PR	HG	CS	CR	PM	DB	TV	AZ	PR	HG	CS	CR	PM	DB	TV	AZ	PR	HG	CS	CR	PM
VilLain-S	69.5	OOM	OOM	83.1	96.8	59.9	71.0	77.1	79.7	OOM	OOM	79.2	86.4	80.9	75.4	78.9	35.7	OOM	OOM	<b>88.4</b>	97.9	21.3	30.3	25.8	45.3	OOM	OOM	65.3	98.4	41.4	47.2	60.9
VilLain-M	74.2	75.1	54.8	<u>91.6</u>	<u>98.6</u>	61.1	73.7	<u>78.5</u>	80.7	<u>95.0</u>	<u>94.7</u>	<u>83.0</u>	<u>87.5</u>	<u>82.4</u>	<b>79.0</b>	<u>82.6</u>	43.5	65.4	<u>35.3</u>	<u>87.3</u>	<b>98.8</b>	31.9	<u>46.2</u>	<b>34.7</b>	49.2	46.7	50.4	<u>86.1</u>	<u>98.7</u>	44.0	53.7	62.1
VilLain-L	<u>76.9</u>	<u>79.3</u>	<u>56.7</u>	64.5	97.6	<b>61.9</b>	<u>74.1</u>	78.1	<u>81.4</u>	94.9	94.2	76.6	87.4	<b>82.9</b>	78.8	82.1	42.7	<u>66.6</u>	<b>36.2</b>	64.0	96.6	<b>37.1</b>	43.9	<u>33.0</u>	<u>59.3</u>	<u>66.1</u>	<u>51.6</u>	66.1	97.5	<b>46.5</b>	<u>54.9</u>	<u>63.3</u>
VilLain	<b>77.2</b>	<b>79.4</b>	<b>58.0</b>	<b>93.7</b>	<b>99.2</b>	<u>61.5</u>	<b>75.0</b>	<b>78.8</b>	<b>81.6</b>	<b>95.1</b>	<b>94.9</b>	<b>83.2</b>	<b>87.8</b>	82.1	<b>79.0</b>	<b>82.8</b>	<b>46.6</b>	<b>69.4</b>	35.2	85.7	<u>98.7</u>	<u>34.5</u>	<b>50.4</b>	32.7	<b>60.2</b>	<b>67.2</b>	<b>53.6</b>	<b>91.3</b>	<b>99.0</b>	<u>46.4</u>	<b>58.0</b>	<b>64.4</b>

**Table 5: VilLain benefits from the long-range propagation of v-labels. Increasing both the number of v-label propagation ( $k$  for loss computation and  $k'$  for embedding generation) tends to improve the node classification accuracy.**

	DB	TV	AZ	PR	HG	CS	CR	PM	Rank
$k = 1$	74.25	78.14	52.16	<b>94.67</b>	<b>99.51</b>	60.48	74.96	78.21	3.00
$k = 2$	75.76	78.44	55.09	93.43	<u>99.29</u>	60.17	<b>75.15</b>	<u>78.97</u>	2.62
$k = 4$	<u>77.16</u>	<u>79.43</u>	<u>57.95</u>	<u>93.66</u>	99.19	<u>61.53</u>	<u>75.03</u>	78.82	<u>2.25</u>
$k = 8$	<b>78.22</b>	<b>80.24</b>	<b>59.12</b>	92.47	98.78	<b>62.05</b>	74.24	<b>79.22</b>	<b>2.12</b>
$k' = 1$	64.71	60.60	48.46	<b>96.74</b>	<b>99.58</b>	60.62	74.68	77.94	5.62
$k' = 2$	65.29	61.59	49.42	<u>96.36</u>	<u>99.57</u>	60.44	74.70	78.18	5.50
$k' = 4$	66.64	63.25	50.52	96.33	99.39	60.54	74.77	78.29	5.00
$k' = 8$	67.88	65.08	53.22	93.91	99.26	61.29	<b>75.06</b>	78.75	4.50
$k' = 16$	70.83	68.28	54.65	94.49	98.86	61.62	<u>74.86</u>	79.12	<u>3.75</u>
$k' = 32$	73.20	72.31	55.80	<u>92.57</u>	98.59	61.96	74.68	<u>79.22</u>	4.00
$k' = 64$	<u>76.47</u>	<u>76.77</u>	<u>56.42</u>	94.21	98.50	<u>62.42</u>	74.25	78.98	4.00
$k' = 128$	<b>77.62</b>	<b>80.63</b>	<b>57.46</b>	88.68	98.09	<b>63.67</b>	74.41	<b>79.37</b>	<b>3.50</b>

**Effectiveness of Multi-V-label Learning.** To demonstrate the effectiveness of using multiple subspaces, we consider two variants of VilLain: (a) **VilLain-S** learns  $d$  v-labels in a single embedding space and (b) **VilLain-M** learns  $\lceil d/D \rceil$  v-labels in  $D$  subspaces. In Table 4, we compare VilLain with its two variants on the four considered tasks. Regarding VilLain-M, we report the average accuracy when  $\lceil d/D \rceil = \{2, 3, \dots, 8\}$ . We first observe that VilLain-M consistently outperforms VilLain-S, indicating the effectiveness of the multi-v-label propagation. Additionally, introducing multiple subspaces enhances the space complexity, as VilLain-M avoids out-of-memory issues in large hypergraphs like Amazon and Trivago, in contrast to VilLain-S. This aligns with our space complexity analysis presented in Section 5. Furthermore, the superior performance of VilLain over VilLain-M implies that aggregating embeddings from various numbers of v-labels (see Section 6.1 for details) captures more informative potential structure-label relations.

**Effectiveness of Loss Functions.** To examine the effectiveness of the designed loss functions, we consider another variant of VilLain, **VilLain-L**, which only uses the local loss  $\mathcal{L}_{\text{local}}$  to learn v-label distributions. As shown in Table 4, VilLain, which jointly optimizes  $\mathcal{L}_{\text{local}}$  and  $\mathcal{L}_{\text{global}}$  and thus captures both local and global information of the input hypergraph, outperforms VilLain-L, demonstrating the effectiveness of the proposed loss functions. In Appendix C.3, we analyze when  $\mathcal{L}_{\text{global}}$  is particularly beneficial.

**Effects of Long-Range V-label Propagation.** To examine the effects of the long-range propagation of v-labels, we test how the

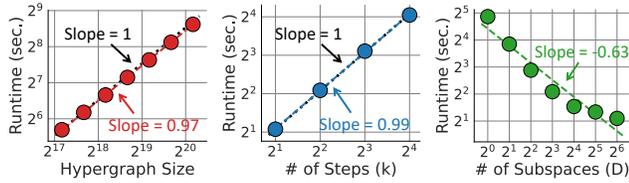
**Table 6: VilLain yields informative embeddings even for unobserved nodes. Fully observed hypergraphs consist of the entire set  $V$  of nodes, whereas partially observed hypergraphs only contain the subset  $V_S \subseteq V$  of nodes after removing 50% of the hyperedges. Despite a performance decrease compared to its fully observable settings, VilLain outperforms its strongest baseline, TriCL in node classification, even for the set  $V \setminus V_S$  of nodes are not observed in VilLain but observed in TriCL during training.**

Learning/Node Type			DB	TV	AZ	CS	CR	PM
Fully Observed	VilLain	$V_S$	78.01	80.03	56.77	62.75	75.43	79.21
		$V \setminus V_S$	66.19	76.07	58.74	57.52	73.46	69.62
Partially Observed	VilLain	$V_S$	76.45	78.66	53.72	62.49	73.79	77.78
		$V \setminus V_S$	65.86	74.71	55.23	56.42	72.27	69.71
Fully Observed	TriCL	$V_S$	69.20	OOM	OOM	60.83	72.94	78.99
		$V \setminus V_S$	55.11	OOM	OOM	53.74	70.02	68.60

number of steps  $k$  (during loss computation) and  $k'$  (during embedding generation) affect the performance of VilLain in node classification. As shown in Table 5, except for Primary and High, which are the smallest datasets, adopting long-range propagation of v-labels is beneficial. In particular, we can see that large datasets (e.g., DBLP, Trivago, and Amazon) benefit from large  $k$ s and  $k'$ s. This tendency holds in other tasks (i.e., hyperedge prediction, node clustering, and node retrieval) as shown in Appendix C.6. This implies that the higher-order label homogeneity, which VilLain aims to reproduce, positively affects the performance in downstream tasks.

## 6.4 Further Analysis of VilLain

In this subsection, we summarize additional experimental results, a part of which is provided in Appendices C and D. Here, we consider the node classification task for evaluation, unless otherwise stated. **Scalability of VilLain.** We test the scalability of VilLain by measuring its training time. In order to test scalability on larger hypergraphs, we upscale Cora using HyperCL [36] by  $2^{\{5, 0.5, 5, \dots, 8, 0\}}$  times. As seen in Figure 4, VilLain scales linearly with the size of the hypergraph and also the number of propagation steps. In addition, the training time decreases with an increased number of subspaces, which is consistent to our time complexity analysis in Section 5. **Performances on Unobserved Nodes.** In Section 5.3, we discussed how VilLain can generate embeddings for nodes that are not observed during training. Instead of using the original hypergraph  $G = (V, E)$ , we evaluate how VilLain, after learning embeddings for the subset  $V_S$  of nodes from a partial hypergraph  $G_S = (V_S, E_S)$ , effectively generates node embeddings for both sets  $V_S$  and  $V \setminus V_S$



**Figure 4: The training time (for 100 epochs) of ViLLain is linear in the hypergraph size (i.e.,  $\sum_{e \in E} |e|$ ) and the number of steps of v-label propagation (i.e.,  $k$ ). The training time decreases with respect to the number  $D$  of subspaces, implying the efficiency of multi-space v-label propagation which is consistent with the complexity analysis in Section 5.**

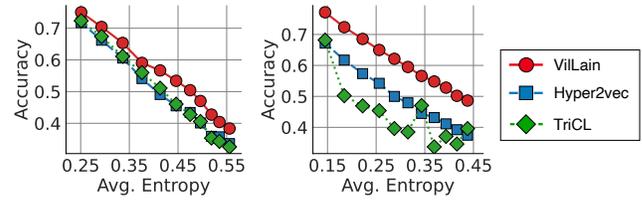
of nodes. Indeed, due to the utilization of reduced structural information, it is natural to expect a degraded quality of node embeddings for both  $V_S$  and  $V \setminus V_S$  sets of nodes in this scenario. This degradation is empirically shown in Table 6 in comparison to the fully-observable setting. However, ViLLain outperforms its strongest baseline, TriCL, across six datasets,<sup>7</sup> even when utilizing partial hypergraphs with 50% of hyperedges removed. TriCL, on the other hand, employs complete hypergraphs to learn embeddings for both sets of nodes. This demonstrates the effectiveness of ViLLain in generating informative embeddings for unobserved nodes, as well as its robustness to the removed hyperedges.

**Performance on Less Homophilic Hypergraphs.** While ViLLain is rooted in the insights gained from the observations of higher-order label homogeneity across various real-world hypergraphs (refer to Section 4), it demonstrates a comparable level of performance also in less homophilic hypergraphs. In Figure 5, we generated semi-real hypergraphs by (1) selecting two hyperedges uniformly at random, and (2) interchanging a single node from each. We repeat this process  $\{100, 200, \dots, 1000\}$  and  $\{1000, 2000, \dots, 10000\}$  times in Cora and DBLP, respectively, resulting in hypergraphs with a diverse range of increased hyperedge entropy (i.e., heterophilicity) and thus less homophilic. From the results, we can observe that the node classification accuracies of ViLLain in Cora and DBLP degrade with the degree of heterophilicity in the hypergraph. Nonetheless, its performance remains superior to that of the two strongest baselines, Hyper2vec and TriCL, demonstrating its effectiveness in less homophilic hypergraphs as well.

**Sensitivity of Multi-V-label Parameters.** We analyze how the parameters related to multi-v-label propagation affect the performance of ViLLain, specifically the number  $D$  of v-label subspaces and the number  $\lceil d/D \rceil$  of v-labels in each subspace. As we can see in Figure 6, both the number of subspaces ( $D$ ) and the number of v-labels in each subspace ( $\lceil d/D \rceil$ ) contribute to the improvement in embedding quality. Empirically, we find that the number of v-labels per subspace has a stronger impact on the performance of ViLLain.

**Additional Experimental Results.** Due to the space limit, other experimental results are provided in Appendix C including (1) usefulness as input features, (2) improvements from external node features, (3) alternative aggregation methods for embedding generation, and (4) comparisons with graph-modeling-based baselines. Furthermore, in Appendix D, we develop ViLLain<sub>B</sub>, a space-efficient

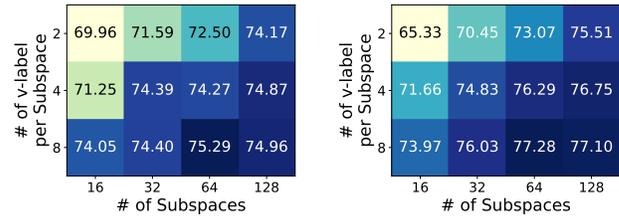
<sup>7</sup>We did not evaluate on Primary and High. Due to their high density, even removing 90% of their hyperedges did not result in any unobserved nodes.



(a) Cora

(b) DBLP

**Figure 5: ViLLain consistently outperforms Hyper2vec and TriCL in node classification across varying levels of average hyperedge entropy (i.e., heterophilicity).**



(a) Cora

(b) DBLP

**Figure 6: Both the number of subspaces ( $D$ ) and the number of v-labels in each subspace ( $\lceil d/D \rceil$ ) are positively correlated to the node classification accuracy.**

variant of ViLLain that generates binary node embeddings for hypergraphs. Empirical results demonstrate its superior performance compared to baseline methods while requiring only  $1/32$  of the bits for encoding the node embedding vectors.

## 7 CONCLUSIONS, LIMITATIONS, AND FUTURE DIRECTIONS

In this work, we propose ViLLain for self-supervised node representation learning on hypergraphs. ViLLain learns node embeddings that reproduces higher-order label homogeneity in real-world hypergraphs, without requiring external node labels or features. We summarize our contributions as follows:

- **Empirical Findings:** We discover the higher-order homogeneity in real-world hypergraphs, which serves as a guiding principle in the design of ViLLain (Section 4).
- **Algorithm Design:** We develop ViLLain, a node embedding method for hypergraphs that does not require external information such as labels or features. It produces versatile embeddings that are effective for various tasks (Section 5).
- **Extensive Experiments:** We demonstrate the overall superiority of ViLLain over 15 unsupervised and (semi-)supervised competitors on eight datasets in four tasks (Section 6).

While higher-order label homogeneity is observed in a majority of real-world hypergraphs, this may not hold in certain hypergraphs with heterophilic characteristics. Extending ViLLain for heterophilic hypergraphs, thus, can be a promising future work.

## REFERENCES

- [1] Ralph Abboud, Ismail Ilkan Ceylan, Martin Grohe, and Thomas Lukasiewicz. 2021. The surprising power of graph neural networks with random node initialization. In *IJCAI*.
- [2] Réka Albert and Albert-László Barabási. 2002. Statistical mechanics of complex networks. *Reviews of Modern Physics* 74, 1 (2002), 47.
- [3] Florian Boudin, Ygor Gallina, and Akiko Aa Aizawa. 2020. Keyphrase Generation for Scientific Document Retrieval. In *ACL*.
- [4] Giorgos Bouritsas, Fabrizio Frasca, Stefanos Zafeiriou, and Michael M Bronstein. 2022. Improving graph neural network expressivity via subgraph isomorphism counting. *TPAMI* 45, 1 (2022), 657–668.
- [5] Shaked Brody, Uri Alon, and Eran Yahav. 2021. How Attentive are Graph Attention Networks?. In *ICLR*.
- [6] Derun Cai, Chenxi Sun, Moxian Song, Baofeng Zhang, Shenda Hong, and Hongyan Li. 2022. Hypergraph contrastive learning for electronic health records. In *SDM*.
- [7] Benjamin Paul Chamberlain, Sergey Shirobokov, Emanuele Rossi, Fabrizio Frasca, Thomas Markovich, Nils Hammerla, Michael M Bronstein, and Max Hansmire. 2023. Graph neural networks for link prediction with subgraph sketching. In *ICLR*.
- [8] Abhra Chaudhuri, Ayan Kumar Bhunia, Yi-Zhe Song, and Anjan Dutta. 2023. Data-Free Sketch-Based Image Retrieval. In *CVPR*.
- [9] Ming Chen, Zhewei Wei, Zengfeng Huang, Bolin Ding, and Yaliang Li. 2020. Simple and deep graph convolutional networks. In *ICML*. PMLR.
- [10] Xu Chen, Siheng Chen, Jiangchao Yao, Huangjie Zheng, Ya Zhang, and Ivor W Tsang. 2020. Learning on attribute-missing graphs. *TPAMI* 44, 2 (2020), 740–757.
- [11] Eli Chien, Chao Pan, Jianhao Peng, and Olga Milenkovic. 2021. You are AllSet: A multitask function framework for hypergraph neural networks. In *ICLR*.
- [12] Eli Chien, Jianhao Peng, Pan Li, and Olga Milenkovic. 2021. Adaptive universal generalized pagerank graph neural network. In *ICLR*.
- [13] Uthsav Chitra and Benjamin Raphael. 2019. Random walks on hypergraphs with edge-dependent vertex weights. In *ICML*.
- [14] Philip S Chodrow, Nate Veldt, and Austin R Benson. 2021. Generative hypergraph clustering: From blockmodels to modularity. *Science Advances* 7, 28 (2021), eabh1303.
- [15] Hejie Cui, Zijie Lu, Pan Li, and Carl Yang. 2022. On positional and structural node features for graph neural networks on non-attributed graphs. In *CIKM*.
- [16] Hande Dong, Jiawei Chen, Fuli Feng, Xiangnan He, Shuxian Bi, Zhaolin Ding, and Peng Cui. 2021. On the equivalence of decoupled graph convolution network and label propagation. In *WWW*.
- [17] Yihe Dong, Will Sawin, and Yoshua Bengio. 2020. HNH: Hypergraph networks with hyperedge neurons. *arXiv preprint arXiv:2006.12278* (2020).
- [18] Dheeru Dua, Casey Graff, et al. 2017. UCI machine learning repository. (2017).
- [19] Chi Thang Duong, Thanh Dat Hoang, Ha The Hien Dang, Quoc Viet Hung Nguyen, and Karl Aberer. 2019. On node features for graph neural networks. *arXiv preprint arXiv:1911.08795* (2019).
- [20] Yifan Feng, Haoxuan You, Zizhao Zhang, Rongrong Ji, and Yue Gao. 2019. Hypergraph neural networks. In *AAAI*.
- [21] Johannes Gasteiger, Aleksandar Bojchevski, and Stephan Günnemann. 2019. Predict then propagate: Graph neural networks meet personalized pagerank. In *ICLR*.
- [22] Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. 2017. Neural message passing for quantum chemistry. In *ICML*. PMLR.
- [23] Aditya Grover and Jure Leskovec. 2016. Node2vec: Scalable feature learning for networks. In *KDD*.
- [24] William L Hamilton, Rex Ying, and Jure Leskovec. 2017. Inductive representation learning on large graphs. In *NeurIPS*.
- [25] Weihua Hu, Matthias Fey, Marinka Zitnik, Yuxiao Dong, Hongyu Ren, Bowen Liu, Michele Catasta, and Jure Leskovec. 2020. Open graph benchmark: Datasets for machine learning on graphs. In *NeurIPS*.
- [26] Weihua Hu, Bowen Liu, Joseph Gomes, Marinka Zitnik, Percy Liang, Vijay Pande, and Jure Leskovec. 2020. Strategies for pre-training graph neural networks. In *ICLR*.
- [27] Jie Huang, Chuan Chen, Fanghua Ye, Jiajing Wu, Zibin Zheng, and Guohui Ling. 2019. Hyper2vec: Biased random walk for hyper-network embedding. In *DASFAA Workshops*.
- [28] Jing Huang and Jie Yang. 2021. Unignn: a unified framework for graph and hypergraph neural networks. In *IJCAI*.
- [29] Yuchi Huang, Qingshan Liu, and Dimitris Metaxas. 2009. Video object segmentation by hypergraph cut. In *CVPR*.
- [30] Hyunjin Hwang, Seungwoo Lee, Chanyoung Park, and Kijung Shin. 2022. Ahp: Learning to negative sample for hyperedge prediction. In *SIGIR*.
- [31] TaeHyun Hwang, Ze Tian, Rui Kuangy, and Jean-Pierre Kocher. 2008. Learning on weighted hypergraphs to integrate protein interactions and gene expressions for cancer outcome prediction. In *ICDM*.
- [32] George Karypis, Rajat Aggarwal, Vipin Kumar, and Shashi Shekhar. 1999. Multi-level hypergraph partitioning: Applications in VLSI domain. *VLSI* 7, 1 (1999), 69–79.
- [33] Eun-Sol Kim, Woo Young Kang, Kyoung-Woon On, Yu-Jung Heo, and Byoung-Tak Zhang. 2020. Hypergraph attention networks for multimodal learning. In *CVPR*.
- [34] Thomas N Kipf and Max Welling. 2017. Semi-supervised classification with graph convolutional networks. In *ICLR*.
- [35] Dongjin Lee and Kijung Shin. 2023. I’m me, we’re us, and I’m us: Tri-directional contrastive learning on hypergraphs. In *AAAI*.
- [36] Geon Lee, Minyoung Choe, and Kijung Shin. 2021. How do hyperedges overlap in real-world hypergraphs?—patterns, measures, and generators. In *WWW*.
- [37] Geon Lee, Jihoon Ko, and Kijung Shin. 2020. Hypergraph motifs: concepts, algorithms, and discoveries. *Vldb* 13, 11 (2020), 2256–2269.
- [38] Seongwon Lee, Suhyeon Lee, Hongje Seong, and Euntai Kim. 2023. Revisiting Self-Similarity: Structural Embedding for Image Retrieval. In *CVPR*.
- [39] Pan Li, Yanbang Wang, Hongwei Wang, and Jure Leskovec. 2020. Distance encoding: Design provably more powerful neural networks for graph representation learning. In *NeurIPS*.
- [40] Zhonghang Li, Chao Huang, Lianghao Xia, Yong Xu, and Jian Pei. 2022. Spatial-temporal hypergraph self-supervised learning for crime prediction. In *ICDE*.
- [41] Xiaorui Liu, Jiayuan Ding, Wei Jin, Han Xu, Yao Ma, Zitao Liu, and Jiliang Tang. 2021. Graph neural networks with adaptive residual. In *NeurIPS*.
- [42] Rossana Mastrandrea, Julie Fournet, and Alain Barrat. 2015. Contact patterns in a high school: a comparison between data collected using wearable sensors, contact diaries and friendship surveys. *PLoS one* 10, 9 (2015), e0136497.
- [43] Miller McPherson, Lynn Smith-Lovin, and James M Cook. 2001. Birds of a feather: Homophily in social networks. *Annual Review of Sociology* 27, 1 (2001), 415–444.
- [44] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781* (2013).
- [45] Jianmo Ni, Jiacheng Li, and Julian McAuley. 2019. Justifying recommendations using distantly-labeled reviews and fine-grained aspects. In *EMNLP*.
- [46] Hoang Nt and Takanori Maehara. 2019. Revisiting graph neural networks: All we have is low-pass filters. *arXiv preprint arXiv:1905.09550* (2019).
- [47] Mingdong Ou, Peng Cui, Jian Pei, Ziwei Zhang, and Wenwu Zhu. 2016. Asymmetric transitivity preserving graph embedding. In *KDD*.
- [48] Prasanna Patil, Govind Sharma, and M Narasimha Murty. 2020. Negative sampling for hyperlink prediction in networks. In *PAKDD*.
- [49] Hongbin Pei, Bingzhe Wei, Kevin Chen-Chuan Chang, Yu Lei, and Bo Yang. 2020. Geom-gcn: Geometric graph convolutional networks. In *ICLR*.
- [50] Zhen Peng, Wenbing Huang, Minnan Luo, Qinghua Zheng, Yu Rong, Tingyang Xu, and Junzhou Huang. 2020. Graph representation learning via graphical mutual information maximization. In *WWW*.
- [51] Bryan Perozzi, Rami Al-Rfou, and Steven Skiena. 2014. Deepwalk: Online learning of social representations. In *KDD*.
- [52] Jiezhong Qiu, Yuxiao Dong, Hao Ma, Jian Li, Chi Wang, Kuansan Wang, and Jie Tang. 2019. NetSMF: Large-scale network embedding as sparse matrix factorization. In *WWW*.
- [53] Emanuele Rossi, Henry Kenlay, Maria I Gorinova, Benjamin Paul Chamberlain, Xiaowen Dong, and Michael M Bronstein. 2022. On the unreasonable effectiveness of feature propagation in learning on graphs with missing node features. In *LoG*. PMLR.
- [54] Ryan Rossi and Nesreen Ahmed. 2015. The network data repository with interactive graph analytics and visualization. In *AAAI*.
- [55] Ryoma Sato, Makoto Yamada, and Hisashi Kashima. 2021. Random features strengthen graph neural networks. In *SDM*. SIAM.
- [56] Aleksandr Schuch, Maximilian Mumme, Aleksandar Bojchevski, and Stephan Günnemann. 2018. Pitfalls of graph neural network evaluation. *arXiv preprint arXiv:1811.05868* (2018).
- [57] Juliette Stehlé, Nicolas Voirin, Alain Barrat, Ciro Cattuto, Lorenzo Isella, Jean-François Pinton, Marco Quaggiotto, Wouter Van den Broeck, Corinne Régis, Bruno Lina, et al. 2011. High-resolution measurements of face-to-face contact patterns in a primary school. *PLoS one* 6, 8 (2011), e23176.
- [58] Jiankai Sun, Bortik Bandyopadhyay, Armin Bashizade, Jiongqian Liang, P Sadayappan, and Srinivasan Parthasarathy. 2019. Atp: Directed graph embedding with asymmetric transitivity preservation. In *AAAI*.
- [59] Shuo Sun, Suzanna Sia, and Kevin Duh. 2020. Clireval: Evaluating machine translation as a cross-lingual information retrieval task. In *ACL*.
- [60] Petar Velčković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. 2018. Graph Attention Networks. In *ICLR*.
- [61] Petar Velčković, William Fedus, William L Hamilton, Pietro Liò, Yoshua Bengio, and R Devon Hjelm. 2019. Deep graph infomax. In *ICLR*.
- [62] Changlin Wan, Muhun Zhang, Wei Hao, Sha Cao, Pan Li, and Chi Zhang. 2021. Principled hyperedge prediction with structural spectral features and neural networks. *arXiv preprint arXiv:2106.04292* (2021).
- [63] Haorui Wang, Haoteng Yin, Muhun Zhang, and Pan Li. 2022. Equivariant and stable positional encoding for more powerful graph neural networks. In *ICLR*.
- [64] Tianxin Wei, Yuning You, Tianlong Chen, Yang Shen, Jingrui He, and Zhangyang Wang. 2022. Augmentations in hypergraph contrastive learning: Fabricated and

- generative. In *NeurIPS*.
- [65] Felix Wu, Amauri Souza, Tianyi Zhang, Christopher Fifty, Tao Yu, and Kilian Weinberger. 2019. Simplifying graph convolutional networks. In *ICML*. PMLR.
- [66] Xiangping Wu, Qingcai Chen, Wei Li, Yulun Xiao, and Baotian Hu. 2020. AdaHGNN: Adaptive hypergraph neural networks for multi-label image classification. In *MM*.
- [67] Lianghao Xia, Chao Huang, Yong Xu, Jiashu Zhao, Dawei Yin, and Jimmy Huang. 2022. Hypergraph contrastive collaborative filtering. In *SIGIR*.
- [68] Lianghao Xia, Chao Huang, and Chuxu Zhang. 2022. Self-supervised hypergraph transformer for recommender systems. In *KDD*.
- [69] Xin Xia, Hongzhi Yin, Junliang Yu, Qinyong Wang, Lizhen Cui, and Xiangliang Zhang. 2021. Self-supervised hypergraph convolutional networks for session-based recommendation. In *AAAI*.
- [70] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. 2018. How powerful are graph neural networks?. In *ICLR*.
- [71] Naganand Yadati, Madhav Nimishakavi, Prateek Yadav, Vikram Nitin, Anand Louis, and Partha Talukdar. 2019. HyperGCN: a new method of training graph convolutional networks on hypergraphs. In *NeurIPS*.
- [72] Naganand Yadati, Vikram Nitin, Madhav Nimishakavi, Prateek Yadav, Anand Louis, and Partha Talukdar. 2020. Nhp: Neural hypergraph link prediction. In *CIKM*.
- [73] Dingqi Yang, Bingqing Qu, Jie Yang, and Philippe Cudre-Mauroux. 2019. Revisiting user mobility and social relationships in lbsns: a hypergraph embedding approach. In *WWW*.
- [74] Zhilin Yang, William Cohen, and Ruslan Salakhudinov. 2016. Revisiting semi-supervised learning with graph embeddings. In *ICML*. PMLR.
- [75] Jaemin Yoo, Hyunseok Jeon, Jinhong Jung, and U Kang. 2022. Accurate node feature estimation with structured variational graph autoencoder. In *KDD*.
- [76] Se-eun Yoon, Hyungseok Song, Kijung Shin, and Yung Yi. 2020. How much and when do we need higher-order information in hypergraphs? a case study on hyperedge prediction. In *WWW*.
- [77] Jiakuan You, Jonathan M Gomes-Selman, Rex Ying, and Jure Leskovec. 2021. Identity-aware graph neural networks. In *AAAI*.
- [78] Seongjun Yun, Seoyoon Kim, Junhyun Lee, Jaewoo Kang, and Hyunwoo J Kim. 2021. Neo-gnns: Neighborhood overlap-aware graph neural networks for link prediction. In *NeurIPS*.
- [79] Junwei Zhang, Min Gao, Junliang Yu, Lei Guo, Jundong Li, and Hongzhi Yin. 2021. Double-scale self-supervised hypergraph learning for group recommendation. In *CIKM*.
- [80] Muhan Zhang and Yixin Chen. 2018. Link prediction based on graph neural networks. In *NeurIPS*.
- [81] Muhan Zhang, Pan Li, Yinglong Xia, Kai Wang, and Long Jin. 2021. Labeling trick: A theory of using graph neural networks for multi-node representation learning. In *NeurIPS*.
- [82] Ruoqi Zhang, Yuesong Zou, and Jian Ma. 2020. Hyper-SAGNN: a self-attention based graph neural network for hypergraphs. In *ICLR*.
- [83] Dengyong Zhou, Jiayuan Huang, and Bernhard Schölkopf. 2006. Learning with hypergraphs: Clustering, classification, and embedding. In *NeurIPS*.
- [84] Yu Zhu, Ziyu Guan, Shulong Tan, Haifeng Liu, Deng Cai, and Xiaofei He. 2016. Heterogeneous hypergraph embedding for document recommendation. *Neurocomputing* 216 (2016), 150–162.
- [85] Yanqiao Zhu, Yichen Xu, Feng Yu, Qiang Liu, Shu Wu, and Liang Wang. 2020. Deep graph contrastive representation learning. *arXiv preprint arXiv:2006.04131* (2020).
- [86] Zhaocheng Zhu, Zuobai Zhang, Louis-Pascal Xhonneux, and Jian Tang. 2021. Neural bellman-ford networks: A general graph neural network framework for link prediction. In *NeurIPS*.

## A DETAILS ON TIME/SPACE COMPLEXITY

In this section, we provide details on the time and space complexity analysis provided in Section 5.

### A.1 Details on Time Complexity

**Time complexity for v-label propagation.** Since we adopt mean-pooling aggregation of  $\frac{d}{D}$  v-labels for both nodes and hyperedges in each subspace, it takes  $O\left(\frac{d}{D} \sum_{e \in E} |e|\right)$ . Thus, for  $D$  subspaces, it takes:

$$O\left(d \sum_{e \in E} |e|\right) \text{ time} \quad (8)$$

for each step of propagation.

**Time complexity for loss computation.** In ViLLain, there are three losses,  $\mathcal{L}_{\text{local}}$ ,  $\mathcal{J}_{\text{cls}}$ , and  $\mathcal{J}_{\text{dst}}$  that are computed to optimize  $\tilde{X}$ .

- For  $\mathcal{L}_{\text{local}}$ , the entropy of the assignment over  $\frac{d}{D}$  v-labels at each node and each hyperedge at each step needs to be computed, and this takes  $O\left(\frac{d}{D}(|V| + |E|)\right)$  time for each subspace. Thus, for  $D$  subspaces, it takes:

$$O(d(|V| + |E|)) \text{ time} \quad (9)$$

for each propagation step.

- For  $\mathcal{J}_{\text{cls}}$ , the entropy of the global assignment over  $\frac{d}{D}$  v-labels needs to be computed at each step, and this takes  $O\left(\frac{d}{D}(|V| + |E|)\right)$  time for each subspace. Thus, for  $D$  subspaces, it takes:

$$O(d(|V| + |E|)) \text{ time} \quad (10)$$

for each propagation step.

- For  $\mathcal{J}_{\text{dst}}$ ,  $\bar{x}_1^{(\ell)}, \dots, \bar{x}_{d/D}^{(\ell)}$  and  $\bar{y}_1^{(\ell)}, \dots, \bar{y}_{d/D}^{(\ell)}$  are required, and it takes  $O\left(\left(\frac{d}{D}\right)^2 |V|\right)$  time and  $O\left(\left(\frac{d}{D}\right)^2 |E|\right)$  time, respectively, to compute them for each subspace. Thus, for  $D$  subspaces, it takes:

$$O\left(\frac{d^2}{D}(|V| + |E|)\right) \text{ time} \quad (11)$$

for each propagation step.

Thus, from Eq. (8)-(11), the time complexity including (a) v-label propagation and (b) loss computation is:

$$O\left(kd \sum_{e \in E} |e| + \frac{kd^2}{D}(|V| + |E|)\right).$$

**Time complexity for embedding generation.** To generate node embeddings using Eq. (3), which requires the mean-pooling propagation of v-labels for  $k'$  steps, it takes:

$$O(k'd(|V| + |E|)) \text{ time.}$$

### A.2 Details on Space Complexity

**Space complexity for v-label propagation.** During its v-label propagation, ViLLain stores assignment matrices of nodes  $X^{(\ell)}$  and hyperedges  $Y^{(\ell)}$  of  $\frac{d}{D}$  v-labels in  $D$  subspaces which requires:

$$O(kd(|V| + |E|)) \text{ space}$$

for  $\ell = 1, \dots, k$  steps.

**Space complexity for loss computation.** The losses  $\mathcal{L}_{\text{local}}^{(\ell)}$  and  $\mathcal{J}_{\text{cls}}^{(\ell)}$  at the  $\ell^{\text{th}}$  step can be computed directly from  $X^{(\ell)}$  and  $Y^{(\ell)}$ , without requiring additional storage space. On the other hand, to compute  $\mathcal{J}_{\text{dst}}^{(\ell)}$  at the  $\ell^{\text{th}}$  step,  $\bar{x}_1^{(\ell)}$  and  $\bar{y}_1^{(\ell)}$  are used, which are computed based on the pairwise cosine similarity between  $d/D$  v-labels, requiring  $O\left(\frac{d^2}{D}(|V| + |E|)\right)$  space for  $D$  subspaces. Thus, the total space required for  $\ell = 1, \dots, k$  steps is;

$$O\left(\frac{kd^2}{D}(|V| + |E|)\right).$$

Note that unlike GNN-based methods [17, 20], ViLLain does not have any additional learnable parameters in each layer.

1161 **Space complexity for embedding generation.** To generate node  
 1162 embeddings,  $X^{(\ell)}$  and hyperedge embeddings  $Y^{(\ell)}$  for  $\ell = 1, \dots, k'$   
 1163 steps are used, and thus  $O(k'd(|V| + |E|))$  space is required.

## 1165 B DETAILS ON EXPERIMENTAL SETTINGS

1166 Here, we provide detailed information on experimental settings.

### 1168 B.1 Details of Datasets

1169 The statistics of the datasets we used are shown in Table 1.

1170 **Preprocessing** For all datasets, we use the largest connected com-  
 1171 ponent of the original hypergraph. We process the huge Amazon  
 1172 by remaining nodes that are from the 10 most frequently appeared  
 1173 labels. Then, we randomly sample 1% of the nodes from each label.

1174 **Ground-truth labels** Here, we provide how the ground-truth la-  
 1175 bels of each dataset are assigned. In Primary and High, each node  
 1176 is a person (e.g., student or teacher), and each hyperedge indicate  
 1177 a group interaction among them. If a person is a teacher, then he or  
 1178 she is labeled as a teacher. Otherwise, students are labeled based  
 1179 on the classroom they belong to. In Citeseer, Cora, and Pubmed,  
 1180 which are co-citation hypergraphs, each node is a paper and each  
 1181 hyperedge is a paper that cited the paper. In these hypergraphs,  
 1182 nodes are assigned by their categories. In DBLP, which is a collabo-  
 1183 ration hypergraph, each node is a paper and each hyperedge is the  
 1184 set of papers written by the same author. Nodes are labeled by their  
 1185 categories. In Trivago, each node is a hotel and each hyperedge is a  
 1186 set of hotels that were clicked in a Web browsing session. Each node  
 1187 is labeled by the location, specifically, the country where the hotel  
 1188 is located. In Amazon, each node is a product, and each hyperedge  
 1189 is a set of products that were co-purchased. Labels of the nodes are  
 1190 assigned by the product categories.

### 1192 B.2 Baselines & Hyperparameters

1194 In this subsection, we discuss the hyperparameters that are used  
 1195 for each method. The implementations we used to run baseline  
 1196 methods are listed in Table 7. Since we consider the unsupervised  
 1197 setting, specifically, without using any labels, the models used for  
 1198 evaluation should be selected without validating on hold-out la-  
 1199 beled data. Thus, for unsupervised baseline methods, we either  
 1200 used their default hyperparameter settings or try to find the set-  
 1201 tings that generally work well across all datasets. However, for  
 1202 (semi-)supervised methods, we use the validation set to tune their  
 1203 hyperparameters.

1204 In VilLain, we fix the number of propagation steps for training  
 1205 to  $k = 4$ , and for inference, we use  $k' = 10$  for small datasets (i.e.,  
 1206 Primary, High, Cora, Citeseer, Pubmed) and  $k' = 100$  for large  
 1207 datasets (i.e., DBLP, Amazon, and Trivago). The learning rate is  
 1208 fixed to 0.01, and the explained variance ratio of the PCA used in  
 1209 VilLain is fixed to 0.99, throughout the experiments.

1210 For Deepwalk [51] and Node2vec [23], we use the default hy-  
 1211 perparameters. Specifically, we set the number of walks to 10, the  
 1212 length of each walk to 80, the window size to 5, and the learning  
 1213 rate to 0.05. For  $p$  and  $q$  in Node2vec, we use 1 for both.

1214 For DGI [61], we use the PReLU for the activation function and  
 1215 set the learning rate to 0.001, as given as default.

1216 For GRACE [85], we use the ReLU for the activation function  
 1217 and the number of GCN layers is set to 2. The learning rate and

1219 **Table 7: Open source links to the baseline source codes.**

Method	Github Link
GCN	<a href="https://pytorch-geometric.readthedocs.io">https://pytorch-geometric.readthedocs.io</a>
GAT	<a href="https://pytorch-geometric.readthedocs.io">https://pytorch-geometric.readthedocs.io</a>
Deepwalk	<a href="https://github.com/benedekrozemberczki/karateclub">https://github.com/benedekrozemberczki/karateclub</a>
Node2vec	<a href="https://github.com/benedekrozemberczki/karateclub">https://github.com/benedekrozemberczki/karateclub</a>
DGI	<a href="https://github.com/PetarV-/DGI">https://github.com/PetarV-/DGI</a>
GRACE	<a href="https://github.com/CRIPAC-DIG/GRACE">https://github.com/CRIPAC-DIG/GRACE</a>
GMI	<a href="https://github.com/zpeng27/GMI">https://github.com/zpeng27/GMI</a>
HGNN	<a href="https://github.com/iMoonLab/HGNN">https://github.com/iMoonLab/HGNN</a>
HNNH	<a href="https://github.com/twistedcubic/HNNH">https://github.com/twistedcubic/HNNH</a>
AllSet	<a href="https://github.com/jianhao2016/AllSet">https://github.com/jianhao2016/AllSet</a>
UniGNN	<a href="https://github.com/OneForward/UniGNN">https://github.com/OneForward/UniGNN</a>
HyperGCL	<a href="https://github.com/weitianxin/HyperGCL">https://github.com/weitianxin/HyperGCL</a>
Hyper2vec	<a href="https://github.com/jeffhj/NHNE">https://github.com/jeffhj/NHNE</a>
TriCL	<a href="https://github.com/wooner49/TriCL">https://github.com/wooner49/TriCL</a>

1235 the weight decay rate are set to 0.001 and 0.00001, respectively.  
 1236 Regarding augmentations (e.g., edge drop and feature drop), all  
 1237 rates are set to 0.2. The dimension of the projection head is set to  
 1238 be the same as the hidden dimension.

1239 For GMI [50], we use the PReLU for the activation function. The  
 1240 learning rate is set to 0.001 without weight decaying. There are  
 1241 three additional hyperparameters  $\alpha$ ,  $\beta$ , and  $\gamma$  that determine the  
 1242 weights of the local and global mutual information, and they are  
 1243 set to  $\alpha = 0.8$ ,  $\beta = 1.0$ , and  $\gamma = 1.0$ , as the default values provided  
 1244 by the authors.

1245 For HyperGCL [64], we use their default hyperparameters. The  
 1246 number of epochs is set to 500, the augmentation ratio is set to 0.3,  
 1247 the temperature is set to 0.3, and the dropout rate is set to 0.2.

1248 For Hyper2vec [27], we use their default hyperparameters. The  
 1249 number of walks is set to 10 and the length of each walk is set to  
 1250 20. The size of the window is 5 and two additional parameters  $p$   
 1251 and  $q$  are both set to 1.

1252 For LBSN [73], the number of negative samples and the learning  
 1253 rate are set to 10 and 0.01, respectively.

1254 For TriCL [35], we set the number of GCN layers to 1 since  
 1255 it was given as default hyperparameters for most datasets. The  
 1256 learning rate and the weight decaying rate are set to 0.0005 and  
 1257 0.00001, respectively. Regarding the data augmentation, the drop  
 1258 rates for node features and the incidence matrix are both set to 0.4.  
 1259 Three temperature hyperparameters,  $\tau_n$ ,  $\tau_g$ , and  $\tau_m$  are all set to  
 1260 0.5, and two weight hyperparameters  $w_g$  and  $w_m$  are set to 4 and 1,  
 1261 respectively.

### 1264 B.3 Node Retrieval Protocol

1265 To perform the node retrieval task, we sample  $\min(|V|, 1000)$  query  
 1266 nodes from the hypergraph uniformly at random. For each query  
 1267 node, we rank the nodes, excluding the query node, based on the  
 1268 cosine similarity between their learned embeddings and the that  
 1269 of the query node. Then, we measure the Mean Average Precision  
 1270 (MAP), which is commonly employed in information retrieval tasks  
 1271 (e.g., computer vision [8, 38] or natural language processing [3, 59]).  
 1272 Here, we define nodes with labels same as that of the query node as  
 1273 the ground-truth. Thus, the MAP yields a higher score when nodes  
 1274 belonging to the same class as the query node are ranked highly.

**Table 8: The number of label propagation steps required for the average entropy of the hyperedges in the real-world hypergraphs to reach  $\epsilon$  of that of the hyperedges in the randomized hypergraphs.**

	PR	HG	CS	CR	PM	DB	TV	AZ
$\epsilon = 0.9$	6	15	140	31	16	832	812	22
$\epsilon = 0.99$	15	34	456	102	43	2,813	2,234	47
$\epsilon = 0.999$	22	47	$\infty$	161	70	4,046	3,409	59

## B.4 Hyperedge Prediction Protocol

To perform the hyperedge prediction task, we first split the original hypergraph  $G = (V, E)$  into two sub-hypergraphs  $G_{\text{train}} = (V_{\text{train}}, E_{\text{train}})$  and  $G_{\text{test}} = (V_{\text{test}}, E_{\text{test}})$  where  $E = E_{\text{train}} \cup E_{\text{test}}$  and  $E_{\text{train}} \cap E_{\text{test}} = \emptyset$ . We also ensure that all nodes are contained in  $G_{\text{train}}$  (i.e.,  $V_{\text{train}} = V$ ) so that embeddings of all nodes in  $G$  are learned. Given a train ratio  $\gamma$ , we set the number of hyperedges in  $G_{\text{train}}$  and  $G_{\text{test}}$  to be divided based on it, i.e.,  $|E_{\text{train}}| : |E_{\text{test}}| = \gamma : 1 - \gamma$ . Specifically, we set  $\gamma = 0.80$  for all datasets except for Amazon, which is relatively very sparse, and thus we set  $\gamma = 0.95$ .

Once we obtain node embeddings of all nodes  $V$ , we generate sets of fake hyperedges  $E_{\text{train}}^{\text{fake}}$  and  $E_{\text{test}}^{\text{fake}}$  as counterparts of true hyperedges  $E_{\text{train}}$  and  $E_{\text{test}}$ . Specifically, for each true hyperedge  $e \in E_{\text{train}}$  (or  $E_{\text{test}}$ ), we randomly sample  $|e|$  nodes from  $V$  and create  $e' \in E_{\text{train}}^{\text{fake}}$  (or  $E_{\text{test}}^{\text{fake}}$ ). Then, a logistic regression classifier is trained on the  $E_{\text{train}} \cup E_{\text{train}}^{\text{fake}}$  and the performance of the hyperedge prediction is evaluated on  $E_{\text{test}} \cup E_{\text{test}}^{\text{fake}}$ .

## C ADDITIONAL EXPERIMENTAL RESULTS

In this section, we provide additional experimental results that are not covered in the main context.

### C.1 Higher-Order Homogeneity

We examine the number of steps of label propagation required for the average entropy of the hyperedges in the real-world hypergraphs to reach  $\epsilon$  of that of the hyperedges in the randomized hypergraphs. For example, it requires 3,409 steps of label propagation to reach 0.999 of the average entropy of random hypergraphs, as shown in Table 8. These results support Observation 2, i.e., real-world hypergraphs exhibit not only the hyperedge-level label homogeneity but also the higher-order homogeneity.

### C.2 Full Results

We provide the full results on the three considered downstream tasks: node classification (Table 16), hyperedge prediction (Table 17), node clustering (Table 18), and node retrieval (Table 19). In these tables, we include the results of the space-efficient version, ViLLain<sub>B</sub> (see Appendix D). For ViLLain<sub>B</sub>, we consider two variants, ViLLain<sub>B</sub><sup>128</sup> and ViLLain<sub>B</sub><sup>256</sup>, which generate binary embeddings that cost 128 and 256 bits, respectively, for each embedding vector. We set the number of v-labels in each subspace to 4 for both variants. Note that ViLLain<sub>B</sub><sup>128</sup> and ViLLain<sub>B</sub><sup>256</sup> require only 1/32 and 1/16 of the space used by the other methods, respectively. We include the standard deviation in the tables. In node classification, hyperedge prediction, and node retrieval tasks, on average,

**Table 9: Density (i.e.,  $|E|/|V|$ ) and overlapness (i.e.,  $\sum_{e \in E} |e|/|V|$ ) of each dataset. Primary exhibits exceptionally high density and overlapness compared to other datasets.**

	PR	HG	CS	CR	PM	DB	TV	AZ
Density	<b>52.495</b>	23.908	0.803	1.130	2.079	0.522	1.350	0.122
Overlapness	<b>126.979</b>	55.633	2.755	3.457	9.049	2.510	4.207	1.622

**Table 10: ViLLain benefits from input node features in node classification. When utilizing node features, it ranks highest on average among its feature-requiring baselines across four datasets where node features are provided.**

Method	DBLP	Citeseer	Cora	Pubmed	Rank
GCN	84.45 $\pm$ 1.25	64.60 $\pm$ 3.00	76.06 $\pm$ 2.29	74.92 $\pm$ 2.90	5.25 $\pm$ 2.62
GAT	77.07 $\pm$ 1.63	50.39 $\pm$ 3.40	59.79 $\pm$ 2.08	73.96 $\pm$ 2.13	11.00 $\pm$ 1.15
DGI	85.64 $\pm$ 1.13	<b>68.53 <math>\pm</math> 2.91</b>	<b>77.50 <math>\pm</math> 2.04</b>	75.62 $\pm$ 2.82	3.50 $\pm$ 2.38
GRACE	85.63 $\pm$ 1.05	61.33 $\pm$ 2.78	71.16 $\pm$ 1.81	77.47 $\pm$ 1.57	5.75 $\pm$ 3.30
GMI	80.85 $\pm$ 1.49	57.09 $\pm$ 2.68	74.73 $\pm$ 1.69	76.38 $\pm$ 2.21	8.00 $\pm$ 2.44
HGNN	84.36 $\pm$ 1.70	64.28 $\pm$ 2.53	75.63 $\pm$ 1.39	76.63 $\pm$ 2.44	5.00 $\pm$ 0.81
HNHN	74.44 $\pm$ 1.98	58.53 $\pm$ 3.31	67.87 $\pm$ 3.51	69.38 $\pm$ 3.47	10.75 $\pm$ 1.25
AllSet	83.67 $\pm$ 1.53	57.88 $\pm$ 3.14	70.07 $\pm$ 3.23	75.24 $\pm$ 2.93	9.00 $\pm$ 1.15
UniGNN	84.22 $\pm$ 1.57	63.79 $\pm$ 3.72	74.44 $\pm$ 2.50	76.99 $\pm$ 2.82	5.75 $\pm$ 1.89
HyperGCL	76.12 $\pm$ 6.04	63.30 $\pm$ 2.11	73.01 $\pm$ 3.68	<b>82.62 <math>\pm</math> 3.25</b>	6.75 $\pm$ 4.19
TriCL	<b>86.59 <math>\pm</math> 0.88</b>	64.53 $\pm$ 3.17	<b>79.03 <math>\pm</math> 0.63</b>	76.60 $\pm$ 1.71	<b>2.75 <math>\pm</math> 2.06</b>
ViLLain	<b>85.68 <math>\pm</math> 0.85</b>	<b>68.77 <math>\pm</math> 1.82</b>	76.54 $\pm$ 1.44	<b>78.25 <math>\pm</math> 2.41</b>	<b>2.00 <math>\pm</math> 0.81</b>

ViLLain and ViLLain<sub>B</sub> show the best performance. In the node clustering task, ViLLain show the second-best performance. Notably, ViLLain<sub>B</sub><sup>128</sup> and ViLLain<sub>B</sub><sup>256</sup>, which require substantially less number of bits for embeddings than the other, highly rank on average. Moreover, it is worthwhile to notice that the proposed methods outperform (semi-)supervised methods (e.g., HGNN and AllSet), which are trained specifically for the node classification task. We conjecture that v-label propagation inherits rich structural properties and also potential higher-order structure-label relationships, generating high-quality representations of nodes.

### C.3 When $\mathcal{L}_{\text{global}}$ is Important

As shown in Table 4 in Section 6.3, ViLLain outperforms ViLLain-L in most datasets. Notably, this performance advantage is particularly significant in Primary, and in this subsection, we analyze the reasons behind this improvement and explore when the inclusion of  $\mathcal{L}_{\text{global}}$  is particularly beneficial. We hypothesize that ViLLain-L faces difficulty in learning distinctive v-label distributions, with a single v-label accounting for nearly 100% of nodes in Primary, regardless of the predefined number of v-labels. This challenge may arise due to the dataset’s unique characteristic of densely connected nodes. This is supported by the measured density (i.e.,  $|E|/|V|$ ) and overlapness (i.e.,  $\sum_{e \in E} |e|/|V|$ ) of the hypergraphs in Table 9.

### C.4 Improvements from Node Features

External node features, if available, are useful and typically enhance method performance. ViLLain can be extended to incorporate node features by introducing  $|V|$  additional hyperedges, where each hyperedge is a group of the  $k$ -nearest neighbors of each node based on cosine similarity between node features. Then, it learns v-label

**Table 11: The average accuracy over all feature-requiring methods (e.g., GCN, HGNN, and TriCL) using different input features. Hyper2vec is the most useful input feature, compared to learnable embeddings and Node2vec.**

Method	DBLP	Trivago	Amazon	Primary	High	Citeseer	Cora	Pubmed	Rank
Learnable	21.87 ± 3.72	6.75 ± 6.64	13.88 ± 5.73	61.71 ± 24.20	74.19 ± 31.24	41.73 ± 9.26	44.80 ± 13.02	55.11 ± 11.28	2.87 ± 0.33
Node2vec	36.21 ± 5.07	25.11 ± 11.47	26.63 ± 5.75	81.87 ± 10.21	95.97 ± 4.92	53.37 ± 4.18	54.81 ± 6.48	71.32 ± 5.34	1.62 ± 0.48
Hyper2vec	63.63 ± 6.06	55.51 ± 23.03	OOT	81.79 ± 11.54	92.60 ± 6.48	57.94 ± 3.27	70.05 ± 3.93	74.60 ± 4.85	1.28 ± 0.45

**Table 12: HGNN and TriCL yield unsatisfactory performance with learnable features and random features, while input features learned by Hyper2vec demonstrate significantly better accuracy. VilLain outperforms them by a large margin.**

Method		DBLP	Primary	High	Citeseer	Cora	Pubmed	Rank
HGNN	Learnable Features	21.47 ± 2.28	76.71 ± 3.36	79.58 ± 3.48	42.34 ± 2.11	43.02 ± 2.33	54.66 ± 2.76	4.67 ± 0.47
	Random Features	21.24 ± 1.91	78.43 ± 2.36	83.12 ± 3.65	43.04 ± 3.39	43.26 ± 2.44	54.41 ± 3.46	4.33 ± 0.47
	Hyper2vec	66.60 ± 2.18	88.28 ± 5.02	92.19 ± 3.84	60.91 ± 2.32	72.90 ± 2.00	76.58 ± 2.86	2.67 ± 0.45
TriCL	Learnable Features	19.31 ± 1.11	31.86 ± 2.64	30.34 ± 3.75	24.94 ± 1.62	25.10 ± 2.22	38.74 ± 2.25	6.50 ± 0.50
	Random Features	18.96 ± 1.20	31.84 ± 3.49	38.33 ± 4.42	25.89 ± 2.28	24.29 ± 1.74	39.69 ± 1.93	6.50 ± 0.50
	Hyper2vec	68.18 ± 1.36	92.67 ± 2.50	98.10 ± 1.02	59.17 ± 3.35	72.35 ± 1.53	78.57 ± 1.88	2.33 ± 0.45
VilLain		77.16 ± 1.26	93.66 ± 3.93	99.19 ± 0.41	61.53 ± 3.17	75.03 ± 1.38	78.82 ± 1.47	1.00 ± 0.00

distributions on an augmented hypergraph with  $|V|$  nodes and  $|V| + |E|$  hyperedges. As shown in Table 10, VilLain benefits from using node features, outperforming its feature-requiring baselines in terms of average ranks when using  $k = 3$ .

We would like to emphasize that our simple approach to utilizing external node features is distinguished from how other baseline methods utilize them (i.e., projecting and propagating them through edges), potentially making it a suboptimal choice. However, it is crucial to note that VilLain is primarily designed for scenarios where node features are unavailable and thus is tailored to perform best in such cases. Furthermore, it is important to note that in our other experiments, we used topological node features obtained through Hyper2vec, instead of external features for the baselines that require input node features.

### C.5 Usefulness as Input Features

We evaluate the usefulness of the methods as an input of the feature-requiring methods (i.e., GCN, GAT, DGI, GRACE, GMI, HGNN, HNHN, AllSet, UniGNN, HyperGCL, and TriCL). Specifically, we train these models using three different input features including a learnable one, which is trained together with the models. As shown in Table 11, using Hyper2vec yields the best accuracy in node classification, and thus we use their embeddings for input features of feature-requiring methods.

In addition, in Table 12, we present a comparison of node classification accuracies of HGNN and TriCL, which are semi-supervised and self-supervised GNN methods for hypergraphs, respectively. We utilize different input features across the considered datasets, except for those that result in out-of-memory issues. We can see that GNNs with learnable features and random features yield unsatisfactory performance, while input features learned by Hyper2vec demonstrate significantly better accuracy. Most importantly, VilLain outperforms them by a large margin.

### C.6 Effects of Long-Range V-label Propagation

To examine the effects of the long-range propagation of  $v$ -labels, we test how the number of steps  $k$  (during training) and  $k'$  (during inference) affect the performance of VilLain in the three considered tasks in Tables 20 and 21, respectively. Except for Primary and High, which are the smallest datasets, adopting long-range propagation of  $v$ -labels is beneficial for node classification, node retrieval, and hyperedge prediction. In particular, we can see that large datasets (e.g., DBLP, Trivago, and Amazon) benefit from large  $k$ s and  $k'$ s.

### C.7 Aggregation Method for Embedding Generation

As discussed in Section 5.1, we aggregate embeddings obtained with various numbers of  $v$ -labels. While the aggregation method is flexible, we concatenate embeddings obtained using different numbers of  $v$ -labels, specifically, for each number  $\lceil \frac{d}{B} \rceil \in \{2, 3, \dots, 8\}$  of  $v$ -labels, we learn  $D$  subspaces and then perform PCA to ensure that the final embedding is of the target dimension  $d$ . In Table 13, we compare the performance with VilLain when applying mean-pooling, instead of PCA, to the embeddings from different  $v$ -label numbers, for the embedding aggregation. Across three different downstream tasks, the concatenate-then-PCA outperforms mean-pooling on average.

### C.8 Hyperedge Embedding Method for Hyperedge Prediction

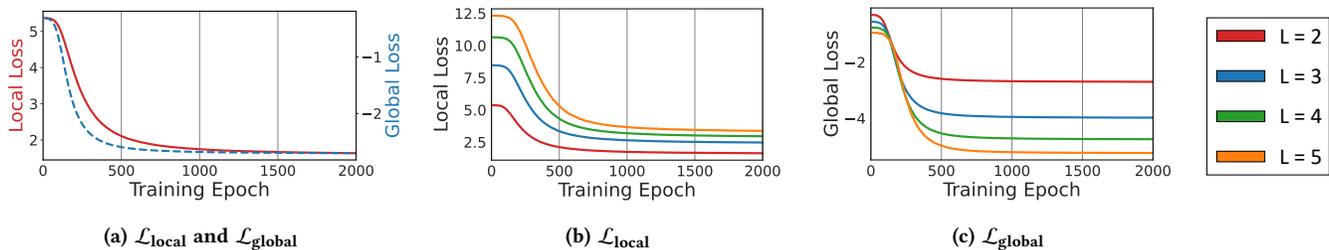
To obtain the embedding of each hyperedge, we apply maxmin pooling, i.e., elementwise max pooling - elementwise min pooling, to the embeddings of the nodes in it. In Table 14, we test the effectiveness of maxmin pooling compared to mean pooling for hyperedge prediction. For both VilLain and TriCL [35], which is the strongest baseline, maxmin pooling is more effective than mean pooling across all datasets.

**Table 13: To aggregate embeddings obtained from various numbers of  $v$ -labels in each subspace, concatenating the embeddings and applying PCA (PCA) outperforms averaging the embeddings (mean) in the three considered downstream tasks.**

	Method	DBLP	Trivago	Amazon	Primary	High	Citeseer	Cora	Pubmed	Rank
NCS	Mean	$74.56 \pm 1.14$	$77.23 \pm 1.35$	$56.36 \pm 2.23$	<b><math>93.88 \pm 3.94</math></b>	$98.95 \pm 0.70$	<b><math>62.70 \pm 2.78</math></b>	$74.38 \pm 1.31$	$79.03 \pm 1.64$	$1.62 \pm 0.48$
	PCA	<b><math>77.16 \pm 1.26</math></b>	<b><math>79.43 \pm 1.63</math></b>	<b><math>57.95 \pm 2.47</math></b>	$93.66 \pm 3.93$	<b><math>99.19 \pm 0.41</math></b>	$61.53 \pm 3.17$	<b><math>75.03 \pm 1.38</math></b>	<b><math>78.82 \pm 1.47</math></b>	<b><math>1.37 \pm 0.48</math></b>
HP	Mean	$80.37 \pm 0.97$	$95.11 \pm 0.55$	$94.81 \pm 0.37$	$82.40 \pm 0.89$	$87.21 \pm 0.67$	<b><math>82.66 \pm 0.95</math></b>	$79.44 \pm 0.57$	<b><math>83.10 \pm 0.70</math></b>	$1.62 \pm 0.48$
	PCA	<b><math>81.61 \pm 0.52</math></b>	<b><math>95.12 \pm 0.37</math></b>	<b><math>94.91 \pm 0.36</math></b>	<b><math>83.19 \pm 0.56</math></b>	<b><math>87.79 \pm 0.68</math></b>	$82.08 \pm 1.42$	$78.95 \pm 0.79$	$82.79 \pm 0.79$	<b><math>1.37 \pm 0.48</math></b>
NCT	Mean	$46.32 \pm 1.36$	$65.77 \pm 0.32$	$34.77 \pm 0.50$	<b><math>85.90 \pm 1.30</math></b>	<b><math>98.72 \pm 0.00</math></b>	$34.04 \pm 0.86$	$48.38 \pm 0.95$	$32.62 \pm 0.02$	$1.75 \pm 0.43$
	PCA	<b><math>46.58 \pm 0.62</math></b>	<b><math>69.35 \pm 0.32</math></b>	<b><math>35.24 \pm 0.48</math></b>	$85.67 \pm 1.88$	<b><math>98.72 \pm 0.00</math></b>	<b><math>34.53 \pm 0.45</math></b>	<b><math>50.38 \pm 2.25</math></b>	<b><math>32.73 \pm 0.00</math></b>	<b><math>1.12 \pm 0.22</math></b>
NR	Mean	$49.35 \pm 0.00$	$43.84 \pm 0.00$	$51.26 \pm 0.72$	$86.68 \pm 0.00$	$98.78 \pm 0.00$	$43.95 \pm 0.10$	$52.76 \pm 0.30$	$63.12 \pm 0.37$	$2.00 \pm 0.00$
	PCA	<b><math>60.15 \pm 0.55</math></b>	<b><math>67.23 \pm 0.72</math></b>	<b><math>53.64 \pm 0.47</math></b>	<b><math>91.26 \pm 0.00</math></b>	<b><math>98.99 \pm 0.00</math></b>	<b><math>46.37 \pm 0.00</math></b>	<b><math>57.96 \pm 0.07</math></b>	<b><math>64.43 \pm 0.07</math></b>	<b><math>1.00 \pm 0.00</math></b>

**Table 14: To obtain the embedding of each hyperedge, maxmin-pooling is more effective than mean-pooling in all datasets in both ViLLain and TriCL (the strongest considered baseline method).**

	Method	DBLP	Trivago	Amazon	Primary	High	Citeseer	Cora	Pubmed	Rank
ViLLain	Mean	$52.55 \pm 1.13$	$59.36 \pm 1.52$	$64.99 \pm 2.34$	$57.54 \pm 1.79$	$56.83 \pm 1.81$	$54.03 \pm 2.20$	$56.89 \pm 2.38$	$57.56 \pm 0.99$	$2.00 \pm 0.00$
	Maxmin	<b><math>81.61 \pm 0.52</math></b>	<b><math>95.12 \pm 0.37</math></b>	<b><math>94.91 \pm 0.36</math></b>	<b><math>83.19 \pm 0.56</math></b>	<b><math>87.79 \pm 0.68</math></b>	<b><math>82.08 \pm 1.42</math></b>	<b><math>78.95 \pm 0.79</math></b>	<b><math>82.79 \pm 0.79</math></b>	<b><math>1.00 \pm 0.00</math></b>
TriCL	Mean	$53.60 \pm 0.89$	OOM	OOM	$63.84 \pm 0.97$	$59.29 \pm 1.08$	$61.33 \pm 2.40$	$68.28 \pm 1.50$	$59.69 \pm 1.34$	$2.00 \pm 0.00$
	Maxmin	<b><math>77.40 \pm 0.68</math></b>	OOM	OOM	<b><math>83.00 \pm 0.63</math></b>	<b><math>87.78 \pm 0.52</math></b>	<b><math>81.96 \pm 0.91</math></b>	<b><math>76.69 \pm 0.70</math></b>	<b><math>80.45 \pm 0.75</math></b>	<b><math>1.00 \pm 0.00</math></b>



**Figure 7: Losses with respect to training epochs in ViLLain. (a) The two losses  $\mathcal{L}_{\text{local}}$  and  $\mathcal{L}_{\text{global}}$  are jointly optimized in ViLLain. (b) Optimization with smaller numbers of  $v$ -labels is easier to optimize  $\mathcal{L}_{\text{local}}$ . (c) On the other hand, optimization with larger numbers of  $v$ -labels is easier to optimize  $\mathcal{L}_{\text{global}}$ .**

**Table 15: ViLLain, applied to hypergraphs, outperforms the recent graph-based baseline methods (i.e., GATv2, GCNII, and GPRGNN) across benchmark graph datasets.**

Graph Type	Model	Cora	Citeseer	Pubmed	Avg. Rank
Graph	GATv2	$70.58 \pm 1.4$	$52.76 \pm 2.4$	$72.92 \pm 2.9$	$3.33 \pm 0.47$
	GCNII	$71.37 \pm 2.4$	$56.47 \pm 1.6$	$74.63 \pm 2.0$	$2.00 \pm 0.00$
	GPRGNN	$69.60 \pm 2.0$	$56.38 \pm 2.2$	$71.44 \pm 1.6$	$3.66 \pm 0.47$
Hypergraph	ViLLain	<b><math>75.03 \pm 1.38</math></b>	<b><math>61.53 \pm 3.17</math></b>	<b><math>78.82 \pm 1.47</math></b>	<b><math>1.00 \pm 0.00</math></b>

## C.9 Comparison with Graph-Modeling-Based Methods

In Section 6, we applied GNNs to pairwise graphs which are transformed from hypergraphs. For this transformation, we adopted clique expansion, which is a popular approach to transform hypergraphs into graphs [72, 83, 84]. However, such clique-expanded graphs are often different from graphs conventionally used for GNN benchmarks. Specifically, for the citation datasets (e.g., Cora, Pubmed, and Citeseer), each edge joins co-cited graphs in clique-expanded graphs, while each edge in GNN-benchmark graphs joins

a pair of citing and cited papers. Thus, we evaluate the performance of well-established GNN models, specifically GATv2 [5], GCNII [9], and GPRGNN [12], when applied to the original structures of the graph datasets. We consider the setting without features, which our paper focuses on and thus use embeddings obtained from Node2vec as their input features. As shown in Table 15, ViLLain outperforms GNN competitors in node classification, even when they use graphs modeled with the same semantics as hypergraphs, instead of clique expansion. This demonstrates the effectiveness of employing hypergraph modeling and ViLLain for learning embeddings from its structure.

## 1625 C.10 Loss of VilLain

1626 We examine how losses of VilLain decrease with training epochs.  
 1627 As discussed in Section 5, VilLain optimizes two losses  $\mathcal{L}_{\text{local}}$  and  
 1628  $\mathcal{L}_{\text{global}}$  that aim to capture the local and global structural informa-  
 1629 tion of the hypergraph, respectively. As shown in Figure 7a, the  
 1630 two losses  $\mathcal{L}_{\text{local}}$  and  $\mathcal{L}_{\text{global}}$  jointly decrease as VilLain is trained.  
 1631 In terms of the number  $L = \lceil d/D \rceil$  of v-labels in each subspace, the  
 1632 decrease of  $\mathcal{L}_{\text{local}}$  is facilitated by smaller  $L$ . Intuitively, a smaller  
 1633 number of v-labels is more likely to lead to homogeneous hyper-  
 1634 edges. On the other hand,  $\mathcal{L}_{\text{global}}$  decreases faster with a larger  
 1635 number of v-labels in each subspace since more diverse v-labels are  
 1636 more likely to be distinctive from each other.

## 1638 D VILLAIN<sub>B</sub>: SPACE-EFFICIENT BINARY 1639 EMBEDDING

1640 As hypergraphs grow in size, so does the space required to store  
 1641 the embeddings. Specifically, a continuous  $d$ -dimensional vector  
 1642 consisting of  $d$  real numbers requires  $32d$  bits if float-32 is used to  
 1643 represent each real number. To reduce the space requirement, we  
 1644  
 1645  
 1646  
 1647  
 1648  
 1649  
 1650  
 1651  
 1652  
 1653  
 1654  
 1655  
 1656  
 1657  
 1658  
 1659  
 1660  
 1661  
 1662  
 1663  
 1664  
 1665  
 1666  
 1667  
 1668  
 1669  
 1670  
 1671  
 1672  
 1673  
 1674  
 1675  
 1676  
 1677  
 1678  
 1679  
 1680  
 1681  
 1682

1683 propose VilLain<sub>B</sub>, a space-efficient version of VilLain that produces  
 1684 binary node embeddings for hypergraphs. Specifically, we binarize  
 1685 the continuous vector  $\mathbf{Z}_i^{(t)}$  of node  $v_i$  in each  $t^{\text{th}}$  subspace obtained  
 1686 by VilLain, which is a probabilistic distribution over  $d/D$  v-labels,  
 1687 to a one-hot vector  $\widehat{\mathbf{Z}}_i^{(t)} \in \{0, 1\}^{d/D}$  as:

$$1688 \widehat{\mathbf{Z}}_i^{(t)} = \text{one-hot} \left( \arg \max_j \left( \mathbf{Z}_{i,j}^{(t)} \right) \right).$$

1689 Then, the final binarized embedding  $\widehat{\mathbf{Z}}_i \in \{0, 1\}^d$  is obtained by  
 1690 concatenating the binarized embeddings from the  $D$  subspaces.

1691 To encode a  $d/D$ -dimensional one-hot vector in each subspace,  
 1692  $\lceil \log_2 \frac{d}{D} \rceil$  bits are required. Hence, encoding a final binarized vec-  
 1693 tor, which is the concatenation of  $D$  one-hot vectors, requires  
 1694  $\lceil D \log_2 \frac{d}{D} \rceil$  bits. Note that,  $D \log_2 \frac{d}{D} < 32d$  always holds for any  
 1695 positive integers  $d$  and  $D$  ( $\leq d$ ).

1696 In Tables 16, 17, 18, and 19, we include the performance of  
 1697 VilLain<sub>B</sub>. While requiring a substantially smaller number of bits to  
 1698 encode embeddings, VilLain<sub>B</sub> outperforms baseline methods in the  
 1699 considered four downstream tasks.  
 1700  
 1701  
 1702  
 1703  
 1704  
 1705  
 1706  
 1707  
 1708  
 1709  
 1710  
 1711  
 1712  
 1713  
 1714  
 1715  
 1716  
 1717  
 1718  
 1719  
 1720  
 1721  
 1722  
 1723  
 1724  
 1725  
 1726  
 1727  
 1728  
 1729  
 1730  
 1731  
 1732  
 1733  
 1734  
 1735  
 1736  
 1737  
 1738  
 1739  
 1740

**Table 16: Full results on node classification (in terms of accuracy). ViLLain and ViLLain<sub>B</sub> outperform the existing (hyper)graph representation learning methods.**

Method	DBLP	Trivago	Amazon	Primary	High	Citeseer	Cora	Pubmed	Rank
GCN	49.65 ± 2.91	18.53 ± 3.61	19.08 ± 2.43	89.54 ± 2.49	89.82 ± 4.99	53.35 ± 3.76	63.71 ± 2.73	70.89 ± 1.60	12.62 ± 2.64
GAT	OOM	OOM	OOM	58.48 ± 7.12	76.94 ± 9.60	51.06 ± 4.29	62.74 ± 3.07	61.66 ± 5.27	16.60 ± 1.35
Deepwalk	29.03 ± 1.43	16.85 ± 0.45	25.43 ± 1.72	84.89 ± 3.67	99.31 ± 0.48	45.10 ± 3.18	56.58 ± 1.88	68.58 ± 2.60	13.00 ± 5.04
Node2vec	29.21 ± 1.89	16.88 ± 0.44	25.27 ± 2.36	83.53 ± 3.09	<b>99.38 ± 0.45</b>	45.37 ± 3.17	59.15 ± 1.84	69.05 ± 3.00	12.37 ± 5.09
DGI	62.37 ± 3.32	73.46 ± 1.22	31.80 ± 1.45	86.66 ± 4.51	92.49 ± 0.60	61.36 ± 2.91	71.23 ± 2.04	77.51 ± 1.38	8.37 ± 3.87
GRACE	<u>71.86 ± 2.51</u>	OOM	OOM	63.78 ± 5.12	99.03 ± 0.30	61.16 ± 2.78	73.43 ± 1.81	77.70 ± 1.81	6.50 ± 4.85
GMI	64.19 ± 1.63	OOM	OOM	80.10 ± 4.94	96.61 ± 2.63	58.67 ± 2.68	71.31 ± 1.69	75.51 ± 2.77	10.66 ± 2.05
HGNN	66.60 ± 2.18	OOM	OOM	88.28 ± 5.02	92.19 ± 3.84	60.91 ± 2.32	72.90 ± 2.00	76.58 ± 2.86	8.66 ± 3.19
HNHN	63.99 ± 2.21	59.52 ± 1.64	28.99 ± 2.63	91.31 ± 2.47	96.83 ± 1.25	59.02 ± 1.63	68.81 ± 1.26	75.33 ± 1.77	8.87 ± 2.08
AllSet	63.67 ± 1.89	36.58 ± 0.93	21.75 ± 1.67	85.94 ± 3.02	95.70 ± 1.66	56.08 ± 1.95	67.73 ± 1.81	74.11 ± 2.04	11.75 ± 1.19
UniGNN	67.16 ± 2.15	69.98 ± 1.60	33.77 ± 3.22	88.88 ± 3.58	95.12 ± 3.97	59.10 ± 2.76	71.44 ± 1.03	74.37 ± 2.10	8.37 ± 2.91
HyperGCL	58.72 ± 1.54	74.99 ± 1.23	22.86 ± 2.01	74.07 ± 6.06	85.79 ± 8.92	57.54 ± 1.61	<u>74.99 ± 1.33</u>	78.44 ± 3.33	9.37 ± 5.67
Hyper2vec	67.18 ± 1.78	<u>75.82 ± 1.45</u>	OOT	92.52 ± 2.45	96.34 ± 1.34	<u>61.50 ± 2.60</u>	71.79 ± 1.63	77.04 ± 1.51	5.85 ± 2.84
LBSN	22.63 ± 2.20	47.99 ± 0.82	11.56 ± 0.90	86.71 ± 3.71	95.87 ± 2.28	45.43 ± 2.15	59.70 ± 1.31	54.89 ± 2.38	13.62 ± 3.27
TriCL	68.18 ± 1.36	OOM	OOM	92.67 ± 2.50	98.10 ± 1.02	59.17 ± 3.35	72.35 ± 1.53	<u>78.57 ± 1.88</u>	5.44 ± 2.13
<b>ViLLain<sub>B</sub><sup>128</sup></b>	67.99 ± 1.16	64.93 ± 1.76	52.37 ± 1.82	<b>95.63 ± 0.28</b>	<u>99.32 ± 0.17</u>	60.83 ± 2.82	74.40 ± 1.38	77.57 ± 1.61	4.25 ± 1.98
<b>ViLLain<sub>B</sub><sup>256</sup></b>	70.39 ± 1.76	69.26 ± 1.45	<u>52.40 ± 2.03</u>	<u>95.15 ± 2.04</u>	99.14 ± 0.29	60.27 ± 2.97	74.46 ± 1.88	78.00 ± 1.20	<u>4.00 ± 1.73</u>
<b>ViLLain</b>	<b>77.16 ± 1.26</b>	<b>79.43 ± 1.63</b>	<b>57.95 ± 2.47</b>	93.66 ± 3.93	99.19 ± 0.41	<b>61.53 ± 3.17</b>	<b>75.03 ± 1.38</b>	<b>78.82 ± 1.47</b>	<b>1.62 ± 1.11</b>

**Table 17: Full results on hyperedge prediction (in terms of accuracy). ViLLain and ViLLain<sub>B</sub> (see Appendix D) outperform the existing (hyper)graph representation learning methods.**

Method	DBLP	Trivago	Amazon	Primary	High	Citeseer	Cora	Pubmed	Rank
Deepwalk	63.90 ± 0.94	61.27 ± 1.14	69.36 ± 0.74	<u>83.79 ± 0.68</u>	85.87 ± 0.73	69.55 ± 1.63	67.18 ± 1.32	65.90 ± 0.62	8.00 ± 2.87
Node2vec	64.20 ± 0.79	61.43 ± 0.81	69.29 ± 0.70	<u>83.15 ± 0.86</u>	85.36 ± 0.64	70.35 ± 1.44	66.94 ± 1.57	65.75 ± 0.77	7.87 ± 2.36
DGI	<b>86.05 ± 0.60</b>	83.83 ± 0.70	90.82 ± 0.65	79.06 ± 0.99	84.38 ± 0.77	79.15 ± 0.94	76.33 ± 0.97	80.92 ± 0.74	5.37 ± 2.95
GRACE	<u>85.43 ± 0.76</u>	OOM	OOM	80.32 ± 0.77	87.42 ± 0.46	77.88 ± 1.31	74.52 ± 0.78	79.05 ± 0.75	5.00 ± 1.82
GMI	75.60 ± 0.71	OOM	OOM	82.43 ± 0.68	85.90 ± 0.60	74.41 ± 1.25	69.40 ± 1.38	72.34 ± 0.70	7.16 ± 1.21
Hyper2vec	71.19 ± 1.01	72.36 ± 1.08	OOT	76.41 ± 0.92	79.57 ± 0.85	78.05 ± 1.76	71.65 ± 1.54	71.48 ± 0.88	8.14 ± 1.95
LBSN	48.68 ± 1.08	89.08 ± 0.68	63.65 ± 1.60	79.43 ± 0.80	87.05 ± 0.60	74.29 ± 1.64	69.63 ± 0.98	66.10 ± 0.77	7.62 ± 2.34
TriCL	77.40 ± 0.76	OOM	OOM	<b>83.99 ± 0.70</b>	<u>87.78 ± 0.44</u>	81.96 ± 1.42	76.69 ± 0.79	80.45 ± 0.67	3.66 ± 1.69
<b>ViLLain<sub>B</sub><sup>128</sup></b>	79.39 ± 0.78	93.49 ± 0.66	<u>92.97 ± 0.60</u>	79.76 ± 0.54	86.05 ± 0.51	<u>82.48 ± 1.29</u>	78.92 ± 1.27	<u>81.42 ± 0.79</u>	3.87 ± 2.08
<b>ViLLain<sub>B</sub><sup>256</sup></b>	79.44 ± 0.68	<u>93.90 ± 0.75</u>	92.64 ± 0.52	79.81 ± 0.75	86.35 ± 0.65	<b>83.03 ± 1.18</b>	<b>79.95 ± 1.35</b>	80.69 ± 0.71	<u>3.37 ± 1.93</u>
<b>ViLLain</b>	81.61 ± 0.52	<b>95.12 ± 0.37</b>	<b>94.91 ± 0.36</b>	83.19 ± 0.56	<b>87.79 ± 0.68</b>	82.08 ± 1.42	<u>78.95 ± 0.79</u>	<b>82.79 ± 0.79</b>	<b>1.87 ± 0.92</b>

**Table 18: Full results on node clustering (in terms of normalized mutual information). ViLLain and ViLLain<sub>B</sub> (see Appendix D) outperform the existing (hyper)graph representation learning methods.**

Method	DBLP	Trivago	Amazon	Primary	High	Citeseer	Cora	Pubmed	Rank
Deepwalk	0.70 ± 0.05	16.70 ± 0.21	7.75 ± 0.03	85.15 ± 0.39	<b>100.00 ± 0.00</b>	14.62 ± 1.73	23.87 ± 1.83	<b>34.35 ± 0.33</b>	6.50 ± 3.60
Node2vec	0.91 ± 0.03	16.96 ± 0.54	7.73 ± 0.13	83.47 ± 0.70	<b>100.00 ± 0.00</b>	14.52 ± 0.82	23.80 ± 1.21	32.83 ± 0.07	7.50 ± 3.20
DGI	16.63 ± 0.01	44.50 ± 1.68	13.01 ± 0.95	84.43 ± 1.68	73.88 ± 0.95	29.09 ± 0.61	32.07 ± 0.79	31.27 ± 0.00	7.50 ± 2.06
GRACE	42.96 ± 0.11	OOM	OOM	67.59 ± 1.12	98.17 ± 0.18	33.04 ± 1.33	46.04 ± 2.44	31.55 ± 0.11	6.16 ± 3.02
GMI	27.80 ± 2.61	OOM	OOM	84.08 ± 0.54	93.10 ± 0.26	25.33 ± 1.72	42.60 ± 3.56	18.71 ± 0.01	8.50 ± 1.25
Hyper2vec	<u>43.40 ± 0.94</u>	<u>66.33 ± 0.27</u>	OOT	<b>92.48 ± 0.35</b>	99.34 ± 0.30	34.28 ± 0.30	45.53 ± 1.05	33.62 ± 0.13	<b>2.57 ± 0.90</b>
LBSN	1.05 ± 0.00	39.41 ± 0.12	2.68 ± 0.33	85.53 ± 0.55	97.80 ± 0.24	12.14 ± 0.43	28.96 ± 0.30	4.62 ± 0.59	8.50 ± 1.87
TriCL	38.00 ± 0.02	OOM	OOM	87.83 ± 1.22	98.74 ± 0.00	<u>34.41 ± 0.02</u>	44.75 ± 0.30	<u>33.74 ± 0.01</u>	3.66 ± 1.37
<b>ViLLain<sub>B</sub><sup>128</sup></b>	35.77 ± 1.92	56.51 ± 0.41	<u>31.94 ± 0.13</u>	89.40 ± 0.04	98.72 ± 0.00	31.60 ± 0.73	44.99 ± 1.84	32.40 ± 0.00	4.75 ± 1.56
<b>ViLLain<sub>B</sub><sup>256</sup></b>	35.90 ± 0.92	58.85 ± 0.40	31.23 ± 0.16	<u>89.76 ± 1.18</u>	98.72 ± 0.00	32.34 ± 1.79	<u>49.08 ± 1.23</u>	33.43 ± 0.02	3.62 ± 1.21
<b>ViLLain</b>	<b>46.58 ± 0.62</b>	<b>69.35 ± 0.32</b>	<b>35.24 ± 0.48</b>	85.67 ± 1.88	98.72 ± 0.00	<b>34.53 ± 0.45</b>	<b>50.38 ± 2.25</b>	32.73 ± 0.00	<u>2.62 ± 2.11</u>

**Table 19: Full results on node retrieval (in terms of mean average precision). VilLain and VilLain<sub>B</sub> (see Appendix D) outperform the existing (hyper)graph representation learning methods.**

Method	DBLP	Trivago	Amazon	Primary	High	Citeseer	Cora	Pubmed	Rank
Deepwalk	21.34 ± 0.24	7.47 ± 0.11	27.71 ± 0.17	81.62 ± 0.00	98.72 ± 0.00	27.61 ± 0.07	29.24 ± 0.14	48.95 ± 0.19	8.37 ± 1.99
Node2vec	21.61 ± 0.20	7.09 ± 0.07	27.78 ± 0.17	81.07 ± 0.00	98.58 ± 0.00	27.29 ± 0.05	29.42 ± 0.15	49.39 ± 0.21	8.50 ± 1.65
DGI	36.06 ± 0.37	37.32 ± 0.13	31.13 ± 0.38	89.73 ± 0.00	97.81 ± 0.00	43.78 ± 0.10	50.64 ± 0.30	61.65 ± 0.36	4.75 ± 1.29
GRACE	50.22 ± 0.60	OOM	OOM	61.41 ± 0.00	<b>99.49 ± 0.00</b>	41.09 ± 0.08	54.17 ± 0.36	60.94 ± 0.36	5.33 ± 3.19
GMI	34.63 ± 0.33	OOM	OOM	80.00 ± 0.00	97.78 ± 0.00	35.89 ± 0.07	41.62 ± 0.28	55.10 ± 0.31	8.33 ± 0.74
Hyper2vec	35.47 ± 0.41	<u>43.11 ± 0.48</u>	OOT	85.74 ± 0.00	90.70 ± 0.00	41.21 ± 0.09	46.67 ± 0.27	55.62 ± 0.19	6.57 ± 2.44
LBSN	21.01 ± 0.09	19.05 ± 0.14	29.11 ± 0.37	81.30 ± 0.00	93.22 ± 0.00	30.60 ± 0.09	40.09 ± 0.26	43.50 ± 0.39	8.62 ± 2.05
TriCL	45.11 ± 0.56	OOM	OOM	89.91 ± 0.00	97.64 ± 0.00	42.37 ± 0.10	54.98 ± 0.26	61.94 ± 0.25	4.66 ± 2.13
<b>VilLain<sub>B</sub><sup>128</sup></b>	47.27 ± 0.54	35.39 ± 0.76	50.34 ± 0.53	89.65 ± 0.00	99.21 ± 0.00	43.33 ± 0.10	53.57 ± 0.29	<u>62.50 ± 0.35</u>	3.87 ± 1.05
<b>VilLain<sub>B</sub><sup>256</sup></b>	53.78 ± 0.58	42.36 ± 0.62	<u>50.61 ± 0.52</u>	<u>90.03 ± 0.00</u>	<u>99.22 ± 0.00</u>	<u>44.35 ± 0.09</u>	<u>56.11 ± 0.30</u>	61.51 ± 0.35	<u>2.50 ± 1.00</u>
<b>VilLain</b>	<b>60.15 ± 0.55</b>	<b>67.23 ± 0.72</b>	<b>53.64 ± 0.47</b>	<b>91.26 ± 0.00</b>	<b>98.99 ± 0.00</b>	<b>46.37 ± 0.10</b>	<b>57.96 ± 0.27</b>	<b>64.43 ± 0.36</b>	<b>1.37 ± 0.99</b>

**Table 20: Effects of  $k$ s in node classification (NCS), hyperedge prediction (HP), and node clustering (NCT) node retrieval (NR). VilLain benefits from the long-range propagation during training.**

Method	DBLP	Trivago	Amazon	Primary	High	Citeseer	Cora	Pubmed	Rank	
NCS	$k = 1$	74.25 ± 1.05	78.14 ± 1.89	52.16 ± 2.36	<b>94.67 ± 2.47</b>	<b>99.51 ± 0.39</b>	60.48 ± 3.17	74.96 ± 1.04	78.21 ± 2.18	3.00 ± 1.22
	$k = 2$	75.76 ± 1.41	78.44 ± 1.84	55.09 ± 2.55	93.43 ± 3.84	<u>99.29 ± 0.39</u>	60.17 ± 3.69	<b>75.15 ± 1.30</b>	<u>78.97 ± 1.38</u>	2.62 ± 0.85
	$k = 4$	77.16 ± 1.26	79.43 ± 1.63	<u>57.95 ± 2.47</u>	<u>93.66 ± 3.93</u>	99.19 ± 0.41	61.53 ± 3.17	75.03 ± 1.38	78.82 ± 1.47	<u>2.25 ± 0.43</u>
	$k = 8$	<b>78.22 ± 1.20</b>	<b>80.24 ± 1.89</b>	<b>59.12 ± 2.58</b>	92.47 ± 4.00	98.78 ± 0.69	<b>62.05 ± 3.52</b>	74.24 ± 1.49	<b>79.22 ± 1.73</b>	<b>2.12 ± 1.45</b>
HP	$k = 1$	80.72 ± 0.68	94.60 ± 0.56	93.81 ± 0.48	<b>83.51 ± 0.64</b>	87.54 ± 0.63	81.51 ± 1.28	77.23 ± 0.91	81.80 ± 0.79	3.50 ± 1.00
	$k = 2$	81.22 ± 0.61	94.82 ± 0.48	94.79 ± 0.44	<u>83.36 ± 0.72</u>	<u>87.69 ± 0.55</u>	82.02 ± 1.22	78.80 ± 0.90	82.29 ± 0.59	2.75 ± 0.43
	$k = 4$	<u>81.61 ± 0.52</u>	<u>95.12 ± 0.37</u>	<u>94.91 ± 0.36</u>	83.19 ± 0.56	<b>87.79 ± 0.68</b>	<u>82.08 ± 1.42</u>	<u>78.95 ± 0.79</u>	<u>82.79 ± 0.79</u>	<u>2.00 ± 0.50</u>
	$k = 8$	<b>81.97 ± 0.72</b>	<b>95.24 ± 0.49</b>	<b>95.27 ± 0.26</b>	82.99 ± 0.56	87.39 ± 0.47	<b>82.62 ± 1.13</b>	<b>79.13 ± 0.82</b>	<b>82.94 ± 0.49</b>	<b>1.75 ± 1.29</b>
NCT	$k = 1$	43.36 ± 1.66	65.28 ± 0.28	32.66 ± 0.22	<b>90.40 ± 1.86</b>	<b>98.72 ± 0.00</b>	32.29 ± 2.06	44.14 ± 2.38	<b>33.96 ± 0.24</b>	2.87 ± 1.45
	$k = 2$	45.50 ± 1.10	67.06 ± 0.52	33.13 ± 0.44	<u>87.70 ± 1.65</u>	<b>98.72 ± 0.00</b>	34.16 ± 1.78	47.38 ± 2.13	<u>33.95 ± 0.27</u>	2.50 ± 0.70
	$k = 4$	<u>46.58 ± 0.62</u>	<u>69.35 ± 0.32</u>	<u>35.24 ± 0.48</u>	85.67 ± 1.88	<b>98.72 ± 0.00</b>	<u>34.53 ± 0.45</u>	<u>50.38 ± 2.25</u>	32.73 ± 0.00	<u>2.12 ± 0.59</u>
	$k = 8$	<b>48.35 ± 0.06</b>	<b>71.44 ± 0.28</b>	<b>36.97 ± 0.07</b>	84.61 ± 1.19	<b>98.72 ± 0.00</b>	<b>35.16 ± 0.42</b>	<b>50.52 ± 0.94</b>	32.52 ± 0.00	<b>1.75 ± 1.29</b>
NR	$k = 1$	55.55 ± 0.63	57.90 ± 0.74	48.44 ± 0.52	91.07 ± 0.00	<b>99.30 ± 0.00</b>	43.78 ± 0.11	53.84 ± 0.27	62.94 ± 0.33	3.50 ± 1.00
	$k = 2$	58.42 ± 0.61	63.12 ± 0.70	51.61 ± 0.42	<b>91.35 ± 0.00</b>	<u>99.10 ± 0.00</u>	45.07 ± 0.11	55.98 ± 0.27	63.60 ± 0.35	2.62 ± 0.69
	$k = 4$	<u>60.15 ± 0.55</u>	<u>67.23 ± 0.72</u>	<u>53.64 ± 0.47</u>	91.26 ± 0.00	98.99 ± 0.00	<u>46.37 ± 0.10</u>	<b>57.96 ± 0.27</b>	<b>64.43 ± 0.36</b>	<b>1.87 ± 0.59</b>
	$k = 8$	<b>62.75 ± 0.48</b>	<b>69.08 ± 0.58</b>	<b>56.01 ± 0.64</b>	90.69 ± 0.00	98.70 ± 0.00	<b>47.68 ± 0.10</b>	<u>57.72 ± 0.25</u>	<u>63.64 ± 0.37</u>	<u>2.00 ± 1.22</u>

**Table 21: Effects of  $k'$ 's in node classification (NCS), hyperedge prediction (HP), and node clustering (NCT) node retrieval (NR). Villain benefits from the long-range propagation at inference (i.e., embedding generation).**

Method	DBLP	Trivago	Amazon	Primary	High	Citeseer	Cora	Pubmed	Rank	
NCS	$k' = 1$	64.71 ± 1.98	60.60 ± 1.54	48.46 ± 3.45	<b>96.74 ± 0.79</b>	<b>99.58 ± 0.20</b>	60.62 ± 2.94	74.68 ± 1.57	77.94 ± 1.84	5.62 ± 2.86
	$k' = 2$	65.29 ± 1.91	61.59 ± 1.62	49.42 ± 2.76	<u>96.36 ± 2.19</u>	<u>99.57 ± 0.21</u>	60.44 ± 3.03	74.70 ± 1.58	78.18 ± 1.75	5.50 ± 2.29
	$k' = 4$	66.64 ± 2.19	63.25 ± 1.63	50.52 ± 2.45	96.33 ± 1.99	99.39 ± 0.21	60.54 ± 3.28	74.77 ± 1.75	78.29 ± 2.14	5.00 ± 1.58
	$k' = 8$	67.88 ± 1.79	65.08 ± 1.55	53.22 ± 3.74	93.91 ± 2.57	99.26 ± 0.38	61.29 ± 3.30	<b>75.06 ± 1.44</b>	78.75 ± 1.91	4.50 ± 1.41
	$k' = 16$	70.83 ± 1.70	68.28 ± 1.38	54.65 ± 3.63	94.49 ± 3.05	98.86 ± 0.79	61.62 ± 3.28	<u>74.86 ± 1.34</u>	79.12 ± 1.44	<u>3.75 ± 0.82</u>
	$k' = 32$	73.20 ± 1.60	72.31 ± 1.65	55.80 ± 2.41	92.57 ± 3.78	98.59 ± 1.42	61.96 ± 3.50	<u>74.68 ± 1.30</u>	<u>79.22 ± 1.69</u>	4.00 ± 1.65
	$k' = 64$	<u>76.47 ± 1.30</u>	<u>76.77 ± 1.71</u>	<u>56.42 ± 2.59</u>	94.21 ± 3.34	98.50 ± 2.31	<u>62.42 ± 2.93</u>	74.25 ± 1.99	78.98 ± 1.68	4.00 ± 2.29
	<b>77.62 ± 1.26</b>	<b>80.63 ± 1.36</b>	<b>57.46 ± 2.04</b>	88.68 ± 4.90	98.09 ± 1.90	<b>63.67 ± 3.26</b>	74.41 ± 1.76	<b>79.37 ± 1.92</b>	<b>3.50 ± 3.24</b>	
HP	$k' = 1$	77.77 ± 0.57	91.46 ± 0.57	93.11 ± 0.56	82.05 ± 0.92	87.45 ± 0.65	82.14 ± 0.97	<u>79.63 ± 0.84</u>	82.09 ± 0.70	6.75 ± 1.92
	$k' = 2$	77.86 ± 0.57	91.53 ± 0.61	93.39 ± 0.75	82.30 ± 0.56	<b>87.74 ± 0.61</b>	82.39 ± 1.59	79.29 ± 0.87	82.27 ± 0.56	5.25 ± 1.98
	$k' = 4$	78.58 ± 0.86	92.51 ± 0.82	93.77 ± 0.77	<u>82.89 ± 0.84</u>	87.53 ± 0.55	81.79 ± 1.24	78.28 ± 0.79	82.12 ± 0.82	5.87 ± 1.83
	$k' = 8$	79.06 ± 0.83	92.31 ± 0.61	94.07 ± 0.51	82.86 ± 0.79	87.69 ± 0.66	82.23 ± 1.19	78.35 ± 0.93	82.46 ± 0.65	5.00 ± 1.11
	$k' = 16$	79.87 ± 0.63	92.97 ± 0.44	94.59 ± 0.39	<b>83.06 ± 0.71</b>	87.55 ± 0.62	82.16 ± 0.97	78.87 ± 0.89	82.49 ± 0.75	4.12 ± 1.45
	$k' = 32$	80.39 ± 0.53	93.80 ± 0.60	<u>95.15 ± 0.38</u>	<u>82.89 ± 0.84</u>	87.28 ± 0.45	82.74 ± 1.15	79.29 ± 0.84	<u>82.92 ± 0.65</u>	3.25 ± 1.56
	$k' = 64$	<u>81.06 ± 0.47</u>	<u>94.58 ± 0.61</u>	<u>95.15 ± 0.42</u>	82.78 ± 0.85	<u>87.70 ± 0.44</u>	<u>83.03 ± 0.92</u>	79.51 ± 0.78	82.75 ± 0.55	<b>2.62 ± 0.99</b>
	<b>81.89 ± 0.87</b>	<b>95.15 ± 0.49</b>	<b>95.21 ± 0.36</b>	81.91 ± 0.77	87.04 ± 0.70	<b>83.42 ± 1.19</b>	<b>80.23 ± 0.85</b>	<b>83.00 ± 0.55</b>	<u>2.75 ± 3.03</u>	
NCT	$k' = 1$	29.62 ± 1.44	48.38 ± 0.50	31.74 ± 0.05	<u>93.08 ± 0.05</u>	<b>98.72 ± 0.00</b>	33.04 ± 1.07	48.93 ± 2.07	32.36 ± 0.00	5.25 ± 2.72
	$k' = 2$	30.30 ± 1.57	49.05 ± 0.33	31.74 ± 0.40	<b>93.09 ± 0.01</b>	<b>98.72 ± 0.00</b>	33.79 ± 0.53	<u>49.78 ± 1.39</u>	32.30 ± 0.00	4.87 ± 2.75
	$k' = 4$	30.26 ± 1.88	50.24 ± 0.32	31.81 ± 0.45	91.95 ± 1.25	<b>98.72 ± 0.00</b>	33.95 ± 0.78	48.40 ± 1.85	32.31 ± 0.48	4.87 ± 1.76
	$k' = 8$	30.39 ± 1.45	51.75 ± 0.75	33.62 ± 0.03	86.48 ± 1.53	<b>98.72 ± 0.00</b>	34.78 ± 0.54	<b>50.74 ± 1.88</b>	<u>32.77 ± 0.01</u>	<b>3.50 ± 1.73</b>
	$k' = 16$	31.90 ± 1.66	54.59 ± 0.27	34.29 ± 0.05	84.45 ± 1.26	<b>98.72 ± 0.00</b>	34.97 ± 0.68	48.65 ± 1.49	32.61 ± 0.00	<u>3.62 ± 1.11</u>
	$k' = 32$	35.63 ± 1.87	58.75 ± 0.66	34.05 ± 0.46	83.87 ± 0.32	98.42 ± 0.19	35.54 ± 1.93	47.49 ± 0.47	<b>32.98 ± 0.00</b>	4.25 ± 2.10
	$k' = 64$	<u>44.85 ± 1.29</u>	<u>64.91 ± 0.36</u>	<u>35.40 ± 0.58</u>	84.87 ± 0.32	97.65 ± 0.00	<b>37.41 ± 0.76</b>	45.56 ± 1.05	32.65 ± 0.00	3.87 ± 2.52
	<b>47.62 ± 0.05</b>	<b>70.88 ± 0.30</b>	<b>35.66 ± 0.61</b>	83.92 ± 0.38	98.09 ± 0.00	<u>36.68 ± 0.22</u>	44.35 ± 0.87	32.12 ± 0.00	4.37 ± 3.15	
NR	$k' = 1$	44.77 ± 0.54	31.08 ± 0.79	51.00 ± 0.42	91.94 ± 0.00	<b>99.41 ± 0.00</b>	45.47 ± 0.11	57.46 ± 0.26	63.04 ± 0.33	5.50 ± 2.39
	$k' = 2$	45.34 ± 0.55	30.95 ± 0.82	50.25 ± 0.41	<b>91.95 ± 0.00</b>	<u>99.35 ± 0.00</u>	45.62 ± 0.11	57.51 ± 0.26	63.15 ± 0.34	5.25 ± 2.33
	$k' = 4$	46.06 ± 0.60	33.43 ± 0.83	48.84 ± 0.45	91.72 ± 0.00	99.22 ± 0.00	45.05 ± 0.10	56.59 ± 0.27	63.21 ± 0.34	5.75 ± 2.04
	$k' = 8$	48.00 ± 0.57	36.56 ± 0.77	51.32 ± 0.40	91.41 ± 0.00	99.06 ± 0.00	46.25 ± 0.10	57.79 ± 0.26	<u>63.46 ± 0.35</u>	4.25 ± 0.96
	$k' = 16$	50.65 ± 0.58	40.19 ± 0.65	51.19 ± 0.39	91.01 ± 0.00	98.79 ± 0.00	46.77 ± 0.10	58.04 ± 0.26	<b>63.50 ± 0.36</b>	4.00 ± 1.22
	$k' = 32$	54.32 ± 0.59	46.72 ± 0.76	<b>54.67 ± 0.37</b>	90.89 ± 0.00	98.49 ± 0.00	<u>47.09 ± 0.10</u>	58.23 ± 0.25	63.31 ± 0.37	<b>3.37 ± 1.65</b>
	$k' = 64$	<u>58.34 ± 0.58</u>	<u>60.01 ± 0.77</u>	<u>53.54 ± 0.48</u>	90.88 ± 0.00	98.33 ± 0.00	<b>47.10 ± 0.10</b>	<b>58.29 ± 0.26</b>	62.97 ± 0.36	<u>3.62 ± 2.64</u>
	<b>61.23 ± 0.53</b>	<b>70.75 ± 0.79</b>	52.60 ± 0.42	90.88 ± 0.00	98.29 ± 0.00	46.80 ± 0.10	<u>58.28 ± 0.26</u>	62.76 ± 0.36	4.12 ± 2.84	