

# THESAURUS: Contrastive Graph Clustering by Swapping Fused Gromov-Wasserstein Couplings

Bowen Deng<sup>1,2\*</sup>, Tong Wang<sup>1\*</sup>, Lele Fu<sup>1,2</sup>, Sheng Huang<sup>1,2</sup>, Chuan Chen<sup>1†</sup>, Tao Zhang<sup>2†</sup>,

<sup>1</sup>School of Computer Science and Engineering, Sun Yat-sen University, Guangzhou, China

<sup>2</sup>School of Systems Science and Engineering, Sun Yat-sen University, Guangzhou, China  
{dengbw3, wangt, full, huangs}@mail2.sysu.com, {chenchuan,zhangt358}@mail.sysu.edu.cn

## Abstract

Graph node clustering is a fundamental unsupervised task. Existing methods typically train an encoder through self-supervised learning and then apply K-means to the encoder output. Some methods use this clustering result directly as the final assignment, while others initialize centroids based on this initial clustering and then finetune both the encoder and these learnable centroids. However, due to their reliance on K-means, these methods inherit its drawbacks when the cluster separability of encoder output is low, facing challenges from the Uniform Effect and Cluster Assimilation. We summarize three reasons for the low cluster separability in existing methods: **(1)** lack of contextual information prevents discrimination between similar nodes from different clusters; **(2)** training tasks are not sufficiently aligned with the downstream clustering task; **(3)** the cluster information in the graph structure is not appropriately exploited. To address these issues, we propose conTrastive graph cluSTering by SwApping fUsed gRomov-wasserstein coUplingS (THESAURUS). Our method introduces semantic prototypes to provide contextual information, and employs a cross-view assignment prediction pretext task that aligns well with the downstream clustering task. Additionally, it utilizes Gromov-Wasserstein Optimal Transport (GW-OT) along with the proposed prototype graph to thoroughly exploit cluster information in the graph structure. To adapt to diverse real-world data, THESAURUS updates the prototype graph and the prototype marginal distribution in OT by using momentum. Extensive experiments demonstrate that THESAURUS achieves higher cluster separability than the prior art, effectively mitigating the Uniform Effect and Cluster Assimilation issues.

**Extended version** — <https://arxiv.org/abs/2412.11550>

## 1 Introduction

Graph node clustering (Wang et al. 2024; Liu et al. 2023b) is a fundamental unsupervised task. Recently, methods based on Graph Self-Supervised Learning (Graph SSL) (Liu et al. 2022a) have become predominant (Liu et al. 2023b). Despite their success, these methods, e.g., Dink-Net (Liu et al. 2023a), heavily rely on K-means (Lloyd 1957; MacQueen

et al. 1967) to guide the representation learning process and/or to get the final clustering results, and thus inherit the shortcomings of K-means. Clusters that contain significantly more samples than others are called majority clusters, and conversely, those with fewer samples minority clusters. When the input node representations exhibit low cluster separability, K-means results may show **(1) Uniform Effect**: samples from majority clusters being assigned to neighboring minority clusters (Xiong, Wu, and Chen 2009), and **(2) Cluster Assimilation**: minority clusters being merged into neighboring majority clusters (Lu, Cheung, and Tang 2021). In contrast, high cluster separability, an ideal clustering outcome characterized by large inter-cluster and small intra-cluster distances, can alleviate these two issues (Lu, Cheung, and Tang 2021), as demonstrated by the Dink-Net finetune effect experiment presented below.

## Uniform Effect & Cluster Assimilation in Dink-Net

The current state-of-the-art (SOTA) model, Dink-Net, is pretrained by distinguishing the original data and the randomly corrupted and shuffled data. We denote the pretrained Dink-Net as Dink-Net-NoFT. After pretraining, K-means is employed to cluster the Dink-Net-NoFT output, initializing the centroids  $\{\mathbf{c}_i\}_{i=0}^{C-1}$ . In the later finetune stage, the encoder and centroids are adjusted to enhance cluster separability by minimizing the dilation loss  $\mathcal{L}_d = \frac{-1}{(C-1)C} \sum_i \sum_{i \neq j} \|\mathbf{c}_i - \mathbf{c}_j\|_2^2$  and shrink loss  $\mathcal{L}_s = \frac{1}{BC} \sum_{i=1}^B \sum_{j=0}^{C-1} \|\mathbf{z}_i - \mathbf{c}_j\|_2^2$ , where  $B$  is the node batch size and  $\mathbf{z}_i$  is the representation. For a node  $i$ , its predicted cluster is  $\hat{y}_i = \arg \min_j \|\mathbf{z}_i - \mathbf{c}_j\|_2$ . Fig. 1 shows the finetune impact on Dink-Net.

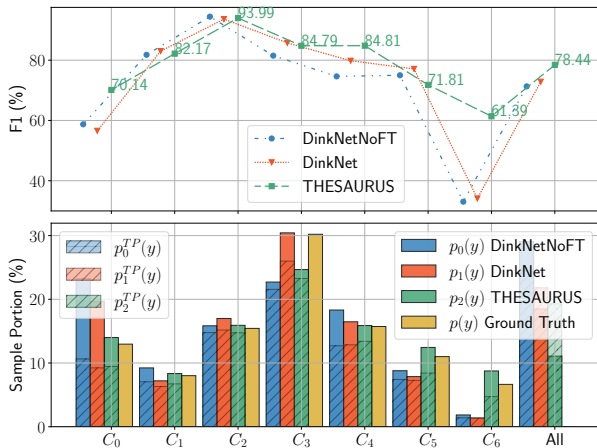
Before finetune, two phenomena are observed on the Cora dataset. **(1) Uniform Effect**: the largest cluster (cluster 3) has many external edges with the much smaller clusters 0 and 4, causing them to be close in the embedding space. This is evidenced in Fig. 1b, where many nodes from cluster 3 are misclassified into clusters 4 and 0. **(2) Cluster Assimilation**: The smallest cluster (cluster 6) has the most external edges with cluster 0, leading to its merging into cluster 0.

After finetuning towards high separability, some of the nodes from cluster 3 but misclassified into clusters 4 and 0 are returned to cluster 3, resulting in a significant improvement in the F1 score for cluster 3. However, since the cluster separability is still insufficient, the smallest cluster (clus-

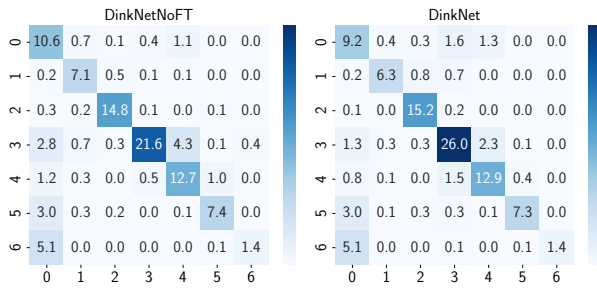
\*These authors contributed equally.

†Corresponding authors.

Copyright © 2025, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.



(a) The F1 and cluster label histograms



(b) The confusion matrices before/after finetune

Figure 1: The effect of separability-oriented finetune of Dink-Net on the Cora dataset. **(a)** The **top row** illustrates the F1 scores for each class before and after finetune, as well as the average F1 score over all classes. The **second row** shows the distribution of predicted labels from three models, along with the ground-truth labels. It also presents the distribution of predicted labels for true-positive (TP) samples, denoted as  $p_i^{TP}(y)$ ,  $i \in \{0, 1, 2\}$ . The final set of bars shows the differences between the predicted and ground-truth distributions. **(b)** displays the confusion matrices (%) of Dink-Net before and after finetune, normalized by the number of nodes.

ter 6) remains merged into cluster 0, leading to a low F1 score on class 6 and a low Macro-F1 score. This experiment underscores the critical role of high cluster separability and reveals that even the current SOTA struggles to achieve sufficient cluster separability to effectively address the challenges of Uniform Effect and Cluster Assimilation.

### Current Model Limitations and Contributions

There are three **limitations (Ls)** impeding current methods achieving high cluster separability. **L1: Lack of contextual information hinders distinguishing synonymous nodes (i.e., similar nodes from different classes).** When a graph has many inter-class edges, low-pass (learnable) graph filters, e.g., GCN (Kipf and Welling 2017), tend to generate similar embeddings for neighboring nodes from different classes (Chen et al. 2020). Distinguishing such nodes based solely on embedding distances is challenging. Just

as synonyms in a thesaurus need context to be properly understood, “childish” implies immaturity in “Stop being so childish!” while “childlike” conveys innocence in “She has a childlike wonder.” Nodes also require contextual information for accurate clustering. Current methods depend only on the embedding distance and no contextual details are provided. As a result, they fail to differentiate between closely embedded nodes from different classes, e.g., mixing nodes from clusters 0, 5, and 6 in Fig. 1. **L2: Training tasks are not well aligned with downstream clustering task.** The alignment between pretext and downstream tasks is crucial for SSL (Lee et al. 2021; Wei et al. 2021). However, current pretext tasks often lack this alignment. **(1)** SDCN (Bo et al. 2020) reconstructs attributes and adopts a dual self-supervised strategy derived from DEC (Xie, Girshick, and Farhadi 2016). DFCN (Tu et al. 2021) reconstructs both attributes and structures with a DEC-like triplet self-supervised task. The reconstruction tasks do not optimize cluster separability and fail to align closely with clustering. DEC-style tasks use K-means to cluster pretrained representations initially. However, the pretrained representations often exhibit low separability, incurring Uniform Effect and Cluster Assimilation. Since finetune only refines the initial clustering, substantial improvements in addressing these issues are not available via DEC-style tasks. **(2)** Regarding the contrastive ones, DCRN (Liu et al. 2022b) and HSAN (Liu et al. 2023d) only preserve self-correlations, not aligned with clustering. SCGC (Liu et al. 2023c) and S<sup>3</sup>GC (Devvrit et al. 2022) treat neighbors as positive pairs, but neighbors are not always of the same class, introducing clustering noise. Although the finetune task of Dink-Net aligns with the clustering objective, the unrelated pretrain task limits the representation separability. Thus, Dink-Net finetune cannot resolve all challenges, as shown in Fig. 1. **L3: The cluster information in graph structure is not appropriately extracted.** Existing methods primarily integrate structure information into embeddings via encoding with GCNs (Kipf and Welling 2017). However, over-smoothing and over-squashing (Nguyen et al. 2023) may hurt structure information. HSAN (Liu et al. 2023d) processes structure through a linear layer, yet it lacks permutation invariance, unduly emphasizing node indices over structure information. DFCN (Tu et al. 2021) and DCRN (Liu et al. 2022b) reflect the structure information through adjacency reflection loss. SCGC (Liu et al. 2023c) and S<sup>3</sup>GC (Devvrit et al. 2022) take neighbors as positive samples and maximizes their similarities. These four methods implicitly treat neighbors as belonging to the same class, misleading the clustering on adjacent nodes from different classes.

To address the above limitations, we propose a novel contrastive graph clustering method, THESAURUS. **(1)** It establishes semantic prototypes in the embedding space, each representing a semantic category. The relationships between one node and these prototypes constitute its context. **(2)** Inspired by SwAv (Caron et al. 2020), THESAURUS considers semantic prototypes as centroids and learns by predicting the node clustering assignments across different data augmentation views. **(3)** To explore the structure cluster information, we encode the relationships between prototypes

as prototype graph, and then match it with the data graph using GW-OT (Mémoli 2011; Peyré and Cuturi 2019). (4) To exploit structure and attribute information comprehensively, designs (2) and (3) are unified by our Task and Structure Alignment (TSA) module based on Fused Gromov-Wasserstein OT (FGW-OT) (Titouan et al. 2019). (5) A momentum module for prototype graph and marginal distribution is developed for data adaptability.

Our main contributions are as follows. (1) We identify that prior methods has insufficient cluster separability and face the Uniform Effect and Cluster Assimilation challenges. (2) We propose a novel graph contrastive learning framework THESAURUS, which leverages semantic prototypes to provide contextual information. We design the TSA module to align the pretext task to clustering and exploit the cluster information in graph structure. We develop a momentum strategy for the prototype graph and prototype marginal distribution for data adaptability. (3) Extensive experiments demonstrate that THESAURUS achieves high cluster separability and significantly outperforms existing methods.

## 2 Related Work

### Deep Graph Clustering

Early graph clustering methods use Autoencoder (AE) and Graph Autoencoder (GAE) (Kipf and Welling 2016) for feature extraction, followed by K-means or spectral clustering (Cao, Lu, and Xu 2016; Wang et al. 2017; Pan et al. 2018). Later methods such as DAEGC (Wang et al. 2019), SDCN (Bo et al. 2020), AGCN (Peng et al. 2021), and DFCN (Tu et al. 2021) incorporate DEC-style tasks to better align with clustering objectives. With the advent of graph contrastive learning, contrastive graph clustering gained popularity. AGE (Cui et al. 2020) uses dynamic positive and negative sample pairs constructed with pair similarities, DCRN (Liu et al. 2022b) introduces DICR loss to reduce cross-view correlations between nodes, SCGC (Liu et al. 2023c) maximizes neighbor similarity, and S<sup>3</sup>GC (Devvrit et al. 2022) optimizes a one-layer GNN with InfoNCE-style loss (van den Oord, Li, and Vinyals 2019). DinkNet (Liu et al. 2023a) maximizes differences between original and adversarial data, akin to maximizing the JSD lower bound of mutual information (Shrivastava et al. 2023), and then finetunes towards cluster separability via minimizing  $\mathcal{L}_d$  and  $\mathcal{L}_s$ .

### Optimal Transport

Optimal transport (OT) (Monge 1781; Villani 2009) is a mathematical framework for measuring distances between distributions, finding the most cost-efficient way to transform one distribution into another. It has gained prominence in machine learning for tasks like domain adaptation (Courty et al. 2017) and generative modeling (Tolstikhin et al. 2018). The optimal transport cost deduces the Wasserstein Distance, which does not require two probability distributions to have overlapping support sets. However, when distributions are defined in different or incomparable spaces, classical OT is not applicable. For example, transporting  $s \in \mathbb{R}^2$  to  $t \in \mathbb{R}^3$  lacks a meaningful cost measurement. Considering the well-defined distances within two distinct spaces  $\mathcal{S}$

and  $\mathcal{T}$ , denoted as  $D_{\mathcal{S}} : \mathcal{S} \times \mathcal{S} \rightarrow \mathbb{R}$  and  $D_{\mathcal{T}} : \mathcal{T} \times \mathcal{T} \rightarrow \mathbb{R}$ , GW-OT (Mémoli 2011) regards the cost of transporting  $s_i \in (\mathcal{S}, D_{\mathcal{S}})$  to  $t_j \in (\mathcal{T}, D_{\mathcal{T}})$  as the difference between the relative relationship between  $s_i$  and other elements  $s_k$  in  $\mathcal{S}$ , and that between  $t_j$  and other elements  $t_l$  in  $\mathcal{T}$ . Such ability to transport across incomparable spaces makes it useful in tasks such as cross-lingual alignment (Alvarez-Melis and Jaakkola 2018) and cross-domain alignment (Gong, Nie, and Xu 2022).

## 3 Methodology

In this section, we present the proposed graph contrastive learning framework in detail.

### THESAURUS Overview and Notations

The attribute graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathbf{X})$ , consisting of  $N$  nodes from  $\mathcal{V}$  and  $E$  edges from  $\mathcal{E}$ , can be summarized by the tuple  $\mathcal{G} = (\mathbf{A}, \mathbf{X})$ . Here,  $\mathbf{A} \in \mathbb{R}^{N \times N}$  is the (binary) adjacency matrix, and  $\mathbf{X} \in \mathbb{R}^{N \times d_0}$  is the node attribute matrix. For convenience, we will use these two graph notations interchangeably in the following sections.

Our framework is illustrated in Fig. 2. Initially, part of edges and feature dimensions of the original graph are masked to generate two distinct (but similar) augmented views  $\mathcal{G}_1 = (\mathbf{A}_1, \mathbf{X}_1)$  and  $\mathcal{G}_2 = (\mathbf{A}_2, \mathbf{X}_2)$ . Subsequently, a GCN encoder  $f_{\theta}$  (Kipf and Welling 2017) and an MLP projector  $f_{\omega} : \mathbb{R}^{d_1} \rightarrow \mathbb{R}^d$  are employed to map these views into representations  $\mathbf{Z}_1 = f_{\omega} \circ f_{\theta}(\mathbf{A}_1, \mathbf{X}_1) \in \mathbb{R}^{N \times d}$  and  $\mathbf{Z}_2$ , where  $f_1 \circ f_2$  denotes the function composition and the whole neural network. The cosine similarities between  $\mathbf{z}_{1,i} = [\mathbf{Z}_1]_i$  and  $S$  semantic prototypes  $\{\mathbf{s}_i \in \mathbb{R}^{1 \times d}\}_{i=1}^S$  form the context-aware representation  $\mathbf{r}_{1,i} \in \mathbb{R}^{1 \times S}$  of node  $i$  in view 1, where  $[\cdot]_i$  is the  $i$ -th row of matrix. Similarly, the context-aware vector  $\mathbf{r}_{2,i}$  in view 2 is obtained.

Let  $\mathbf{S} = [\mathbf{s}_1, \dots, \mathbf{s}_S] \in \mathbb{R}^{S \times d}$  be the learnable prototype matrix,  $\mathbf{R}_1 = \mathbf{Z}_1 \mathbf{S}^{\top} \in \mathbb{R}^{N \times S}$  be the context-aware representation matrix, and  $\mathbf{B}_1 \in \mathbb{R}^{S \times S}$  be the prototype graph in view 1. THESAURUS computes the (optimal) node-prototype assignment  $\mathbf{Q}_1 \in \mathbb{R}^{N \times S}$  for view 1 by solving the FGW-OT problem involving  $\mathbf{R}_1$ ,  $\mathbf{B}_1$ , and  $\mathbf{A}_1$ . Similarly, the assignment  $\mathbf{Q}_2$  is obtained. Once the assignments are got, we train the network  $f_{\omega} \circ f_{\theta}$  to predict  $\mathbf{Q}_2$  from  $\mathbf{Z}_1$  and vice versa. After training, the  $\mathbf{R}$  of the original graph  $\mathcal{G}$  is fed into K-means to get the final result  $\Phi \in \{0, 1\}^{N \times C}$ .

### Address L1: Context via Sematic Prototypes

To capture the subtle differences between synonymous nodes, we draw inspiration from the human ability to accurately distinguish synonyms using textual context. For instance, consider the sentences ‘‘Stop being so childish!’’ and ‘‘She has a childlike wonder.’’ In the first sentence, ‘‘stop being’’ conveys a negative connotation, while ‘‘so’’ is neutral. In the second sentence, ‘‘wonder’’ carries a positive connotation, and ‘‘She has a’’ is neutral. The word co-occurrence in these sentences indicates that ‘‘childish’’ is associated with negative and neutral semantics (prototypes), whereas ‘‘childlike’’ is linked with positive and neutral semantics (prototypes). This allows us to immediately infer that ‘‘childish’’

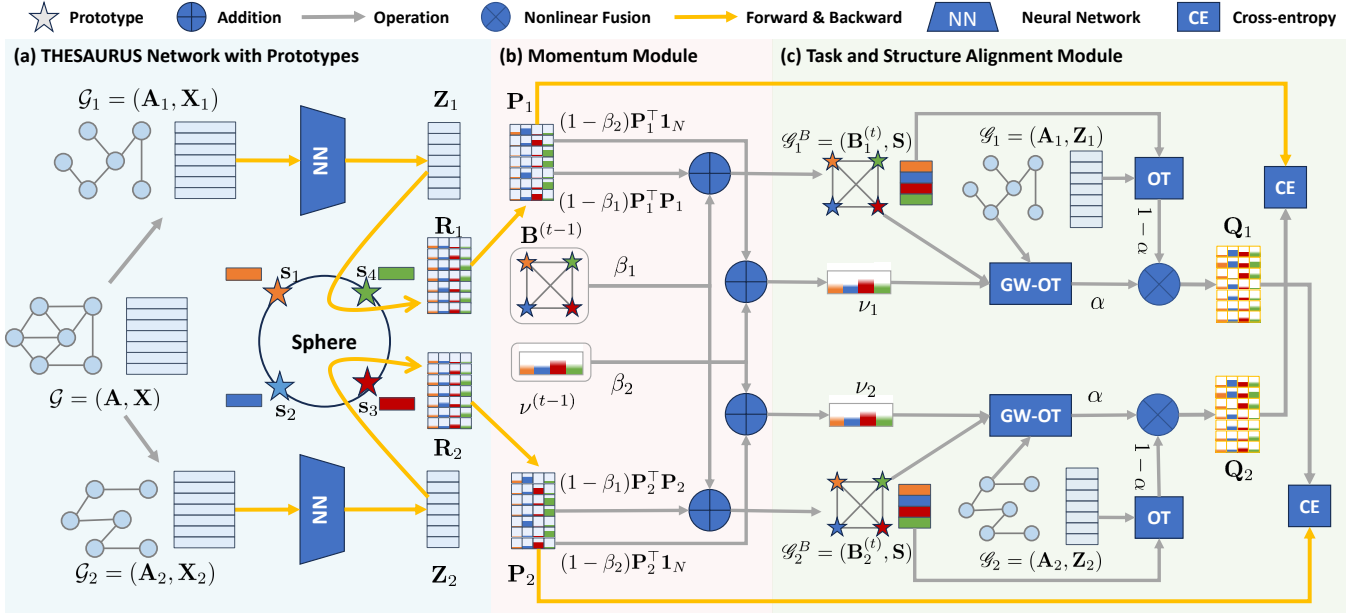


Figure 2: The illustration of our proposed THESAURUS. And the details are summarized in Algorithm 1 in the appendix.

has a negative meaning related to children, while “child-like” has a positive connotation, thus distinguishing these two synonyms.

Similarly, in the embedding space, we define  $S$  semantic prototypes, each representing a specific semantic category. The positional relationships between nodes and these prototypes constitute their contextual semantics. Unlike existing methods that measure the association between two nodes solely by the distance between their embedding vectors, THESAURUS uses that between  $S$ -dimensional context-aware representations — derived from the distances between nodes and semantic prototypes — to represent their associations. With this approach, the input for K-means is not the encoder output, but the context-aware representations  $\mathbf{R}$ .

### Task and Structure Alignment Module

To closely align with the downstream clustering task, we design the pretext task of predicting optimal clustering assignments cross views. To thoroughly exploit the cluster information in the graph structure, we align the prototype and data graph structures. These designs are integrated as the TSA module based on FGW-OT (Titouan et al. 2019).

#### Address L2: Pretext Task Aligned with Clustering

THESAURUS treats semantic prototypes as centroids providing clear semantic meanings and learns by predicting swapped clustering assignments across views. The node-prototype assignment  $\mathbf{Q}$  is derived from context-aware representations  $\mathbf{R} = \mathbf{Z}\mathbf{S}^\top$ , calculated via solving the problem

$$\min_{\pi \in \Pi} \text{Tr}(-\mathbf{R}^\top \pi) - \epsilon H(\pi), \quad (1)$$

where  $\pi \in \mathbb{R}_{\geq 0}^{N \times S}$  deduces  $\mathbf{Q}$  with row-sum normalization,  $H(\pi) = -\sum_{ij} \pi_{ij} \log \pi_{ij}$  denotes the entropy, and

$\epsilon > 0$  controls the smoothness of the unnormalized assignment (i.e., the coupling  $\pi$  in OT). Here  $\pi$  is constrained by the node marginal distribution  $\mu \in \mathbb{R}_{\geq 0}^N$  and the prototype marginal distribution  $\nu \in \{\nu \in \mathbb{R}_{\geq 0}^S \mid \sum_i \nu_i = 1\}$

$$\Pi = \{\pi \mid \pi \mathbf{1}_S = \mu, \pi^\top \mathbf{1}_N = \nu\}, \quad (2)$$

where  $\mathbf{1}_S$  is a  $S$ -dimensional all-one column vector. Problem (1) can be viewed as a relaxed and regularized K-means problem. K-means minimizes the sum of squared distances between data points and centroids, and we take the negative similarities  $-\mathbf{R}$  as “distances” (and OT costs) here. Thanks to the entropy regularization, this problem can be efficiently solved by the scalable Sinkhorn (Cuturi 2013).

Feeding  $\mathbf{R}_1$  of view 1 to the above procedure gives  $\mathbf{Q}_1$ , and similarly  $\mathbf{Q}_2$  is got from view 2. The goal is to predict the assignment  $\mathbf{Q}_2$  of view 2 from the representation  $\mathbf{Z}_1$  of view 1, and  $\mathbf{Q}_1$  from  $\mathbf{Z}_2$ . Since  $\mathbf{Z}_1$  lacks interaction with prototypes  $\mathbf{S}$ , the prediction distribution  $[\mathbf{P}_1^\top]_n$  for node  $n$  is built on  $\mathbf{R}_1 = \mathbf{Z}_1 \mathbf{S}^\top$  instead of  $\mathbf{Z}_1$

$$[\mathbf{P}_1^\top]_{n,s} = \frac{\exp\left([\mathbf{Z}_1]_n [\mathbf{S}]_s^\top / \tau\right)}{\sum_{s'} \exp\left([\mathbf{Z}_1]_n [\mathbf{S}]_{s'}^\top / \tau\right)}, \quad (3)$$

where  $\tau$  is the temperature that controls the distribution sharpness and  $\mathbf{P}_1^\top$  is abbreviated to  $\mathbf{P}_1$ . The distributions of all nodes are stacked into  $\mathbf{P}_1^\top \in \mathbb{R}_{\geq 0}^{N \times S}$ . Similarly,  $\mathbf{P}_2^\top$  is obtained. The overall training loss is then computed as the averaged cross-entropy

$$\mathcal{L} = -\frac{1}{2N} \sum_{n=1}^N \sum_{s=1}^S \left( [\mathbf{Q}_1]_{n,s} \log [\mathbf{P}_2^\top]_{n,s} + [\mathbf{Q}_2]_{n,s} \log [\mathbf{P}_1^\top]_{n,s} \right) \quad (4)$$

**Address L3: Structure Alignment via GW-OT** The above node-prototype clustering assignment is derived from

the relationships between node embeddings  $\mathbf{Z}$  and prototypes  $\mathbf{S}$ . Like prior methods, this approach does not explicitly extract structure cluster information. To fill this gap, we propose deriving the optimal assignment from the structure and using this assignment as a prediction target for  $\mathbf{Z}$ . One effective method to assign nodes to different semantic prototypes based on the structures is to perform GW-OT between  $\mathbf{A}$  and an isolated graph  $\mathbf{I}_S$  (Xu et al. 2019). In this isolated graph, each node represents a cluster with no inter-cluster edges. Such approach adheres to cluster definition but ignores inter-cluster relationships, which is unrealistic. Therefore, we replace the isolated graph with a complete prototype graph  $\mathbf{B} \in \mathbb{R}_{\geq 0}^{S \times S}$ , which is constructed as

$$\mathbf{B} = \mathbf{P}^\top \mathbf{P}. \quad (5)$$

GW-OT is formally defined in Def. 1. For an attribute graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathbf{X})$ , each node  $v_i$  contains observable attribute  $x_i = [\mathbf{X}]_{i,:} \in \Omega_x \subset \mathbb{R}^d$  and implicit structure embedding  $s_i \in \Omega_s$ . Although  $s_i$  is not known, the pairwise relationship  $\mathbf{C} \in \mathbb{R}_{\geq 0}^{N \times N}$  determined by the metric  $D_{\Omega_s} : \Omega_s \times \Omega_s \rightarrow \mathbb{R}_{\geq 0}$  on the space  $\Omega_s$  is given by the adjacency matrix  $\mathbf{A}$ , the Laplacian  $\mathbf{L}$ , or the pairwise shortest path matrix (Chowdhury and Mémoli 2019). To use OT, the probability measure on the space  $(\mathcal{V}, D_{\Omega_s})$  or equivalently  $(\mathcal{V}, \mathbf{C})$  must be defined. Denote the importance of  $N$  nodes by the histogram

$$h \in \mathcal{H}_N = \left\{ h \mid h \in \mathbb{R}_{>0}^N, \sum_{i=1}^N h_i = 1 \right\}. \quad (6)$$

Then this space has a measure  $\mu = \sum_i h_i \delta_{(s_i)}$ , where  $\delta_{(s_i)}$  denotes the Dirac delta function at  $s_i$ .

**Definition 1.** Let  $(\mathcal{V}_1, \mathbf{C}_1, \mu)$  and  $(\mathcal{V}_2, \mathbf{C}_2, \nu)$  be Metric-Measure (MM) spaces defined on  $\mathcal{G}_1 = (\mathbf{C}_1, \emptyset)$  and  $\mathcal{G}_2 = (\mathbf{C}_2, \emptyset)$ , respectively.  $\mu = \sum_i h_i^{(1)} \delta_{(s_i)}$ ,  $h^{(1)} \in \mathcal{H}_{N_1}$  and  $\nu = \sum_i h_i^{(2)} \delta_{(s_i)}$ ,  $h^{(2)} \in \mathcal{H}_{N_2}$  are the probability measures on  $\mathcal{V}_1$  and  $\mathcal{V}_2$ , separately. The Gromov-Wasserstein distance  $GW_p(\mathcal{G}_1, \mathcal{G}_2)$  between these two measures is given by

$$\inf_{\pi \in \Pi} \sum_{i,k=1}^{N_1} \sum_{j,l=1}^{N_2} \left( [\mathbf{C}_1]_{i,k} - [\mathbf{C}_2]_{j,l} \right)^p \pi_{i,j} \pi_{k,l}, \quad (7)$$

where  $\Pi = \left\{ \pi \in \mathbb{R}_{\geq 0}^{N_1 \times N_2} \mid \pi \mathbf{1}_{N_2} = h^{(1)}, \pi^\top \mathbf{1}_{N_1} = h^{(2)} \right\}$ .

We can utilize the  $GW_1$  OT between the non-attribute data graph  $(\mathbf{A}, \emptyset)$  and prototype graph  $(\mathbf{B}, \emptyset)$  to get the optimal node-prototype assignment  $\mathbf{Q}$ . It is row-normalized from the solution  $\pi$  of following OT problem

$$\min_{\pi \in \Pi} \sum_{i,k=1}^S \sum_{j,l=1}^N \left| [\mathbf{A}]_{i,k} - [\mathbf{B}]_{j,l} \right| \pi_{i,j} \pi_{k,l}, \quad (8)$$

where  $\mu$  is the uniform node marginal distribution and  $\nu$  is the current prototype marginal. After the structure-induced assignments  $\mathbf{Q}_1, \mathbf{Q}_2$  of two views  $\mathcal{G}_1, \mathcal{G}_2$  are separately got via Eq. (8), they can be used with the loss Eq. (4).

**Fused Clustering Assignment via FGW-OT** The above introduce two kinds of ‘‘clustering’’ assignments respectively acquired from the context-aware node representation  $\mathbf{R}$  and the graph structure  $\mathbf{A}$ . The assignment from  $\mathbf{R}$  focuses on attribute information, while that from  $\mathbf{A}$  emphasizes structural information. We fuse them with FGW-OT for more comprehensive graph mining. We build  $\mathcal{G} = (\mathbf{A}, \mathbf{Z})$  with the embeddings  $\mathbf{Z}$  as node attributes. And we add prototypes  $\mathbf{S} \in \mathbb{R}^{S \times d}$  as attributes to the prototype graph  $\mathbf{B}$ , resulting in  $\mathcal{G}^B = (\mathbf{B}, \mathbf{S})$ . The optimal transport between  $\mathcal{G}^B$  and  $\mathcal{G}$  encapsulate both attribute and structure cluster information, and can be achieved via FGW-OT defined below.

**Definition 2.** Let  $(\mathcal{V}_1, \mathbf{C}_1, \mu)$  and  $(\mathcal{V}_2, \mathbf{C}_2, \nu)$  be MM-spaces on  $\mathcal{G}_1 = (\mathbf{C}_1, \mathbf{X}_1)$  with measure  $\mu = \sum_i h_i^{(1)} \delta_{(s_i, x_i)}$ ,  $h^{(1)} \in \mathcal{H}_{N_1}$  and on  $\mathcal{G}_2 = (\mathbf{C}_2, \mathbf{X}_2)$  with measure  $\nu = \sum_i h_i^{(2)} \delta_{(s_i, x_i)}$ ,  $h^{(2)} \in \mathcal{H}_{N_2}$ , respectively. The Fused Gromov-Wasserstein distance  $FGW_{p,\alpha}(\mathcal{G}_1, \mathcal{G}_2)$  is

$$\inf_{\pi \in \Pi} \left\{ \sum_{i,k=1}^{N_1} \sum_{j,l=1}^{N_2} \left[ (1-\alpha) D_{\Omega_x}(x_{1,i}, x_{2,j}) + \alpha | \mathbf{C}_1(i,k) - \mathbf{C}_2(j,l) |^p \right] \pi_{i,j} \pi_{k,l} \right\}^{\frac{1}{p}} - \epsilon H(\pi) \quad (9)$$

where  $\Pi = \left\{ \pi \in \mathbb{R}_{\geq 0}^{N_1 \times N_2} \mid \pi \mathbf{1}_{N_2} = h^{(1)}, \pi^\top \mathbf{1}_{N_1} = h^{(2)} \right\}$ ;  $H(\pi) = -\sum_{ij} \pi_{ij} \log \pi_{ij}$  is the coupling entropy and  $\epsilon$  weights this regularization;  $x_{1,i}, x_{2,j} \in \Omega_x \subset \mathbb{R}^d$  are the attributes of nodes  $v_i^{(1)} \in \mathcal{V}_1$  and  $v_j^{(2)} \in \mathcal{V}_2$ , respectively.

The optimal coupling of  $FGW_{1,\alpha}(\mathcal{G}, \mathcal{G}^B)$  encapsulates the information in  $\mathbf{R}$  and  $\mathbf{A}$ , where  $\alpha$  balances the contribution of these two parts. In view 1, we can construct  $\mathcal{G}_1 = (\mathbf{A}_1, \mathbf{Z}_1)$  and  $\mathcal{G}_1^B = (\mathbf{B}_1, \mathbf{S})$ , and then the assignment  $\mathbf{Q}_1$  of  $FGW_{1,\alpha}(\mathcal{G}_1, \mathcal{G}_1^B)$  is obtained via normalizing the optimal coupling  $\pi_1$  at each row  $n \in \{1, 2, \dots, N\}$

$$[\mathbf{Q}_1]_{n,s} = \frac{[\pi_1]_{n,s}}{\sum_{s'=1}^S [\pi_1]_{n,s'}}. \quad (10)$$

Similarly,  $\mathbf{Q}_2$  is got from view 2. Finally, these fused assignments are used to guide the training with Eq. (4).

## Prototype Momentum Module

Prototype graph and marginal should adapt to the data. So momentum update is adopted here. Let  $\mathbf{B}^{(t-1)}$  be the prototype graph of last forward step  $t-1$ . The step  $t$  has

$$\mathbf{B}^{(t)} = \beta_1 \mathbf{B}^{(t-1)} + (1 - \beta_1) \mathbf{P}^\top \mathbf{P}. \quad (11)$$

To reflect the common graph homophily,  $\mathbf{B}$  is initialized as an identity matrix and  $\beta_1 \in [0, 1]$  is set to a high value.

Meanwhile, the logits  $\mathbf{P}$ , e.g. Eq. (3), are used to update prototype marginal, which is initialized as an uniform distribution  $\nu^{(0)} = \mathbf{1}_S / S$  to reflect the presumed even cluster size before data exposure.

$$\nu^{(t)} = \beta_2 \nu^{(t-1)} + (1 - \beta_2) (\mathbf{P}^\top \mathbf{1}_N) \quad (12)$$

Datasets	Metrics	K-means	DEC	SDCN	GRACE	DAEGC	DFCN	DCRN	HSAN	S <sup>3</sup> GC	SCGC	Dink-Net*	Dink-Net	Ours
Cora	ACC	33.80	46.50	35.60	73.90	70.43	36.33	61.93	77.21	74.21	73.88	75.55	78.21	<b>80.72</b>
	NMI	14.98	23.54	14.28	57.10	52.89	19.36	45.13	59.56	58.80	56.10	60.03	62.48	<b>62.92</b>
	ARI	8.60	15.13	7.78	52.70	49.63	4.67	33.15	57.93	54.43	51.79	54.56	61.48	<b>63.61</b>
	F1	30.26	39.23	24.37	72.50	68.27	26.16	49.50	75.13	72.10	70.81	71.31	72.85	<b>78.44</b>
Citeseer	ACC	39.32	46.51	65.96	63.10	64.54	69.50	69.86	71.05	68.81	71.02	69.34	69.91	<b>71.99</b>
	NMI	16.94	23.54	38.71	39.91	36.41	43.90	44.86	45.62	44.11	45.25	44.36	45.29	<b>47.37</b>
	ARI	13.43	15.13	40.17	37.70	37.78	45.50	45.64	48.22	44.80	46.29	45.65	46.29	<b>48.99</b>
	F1	36.08	39.23	63.62	60.30	62.24	64.30	64.83	64.52	64.30	64.80	65.54	65.79	<b>66.42</b>
Pubmed	ACC	59.83	60.14	64.20	63.72	68.73	68.89	69.87	OOM	71.31	45.12	67.32	67.51	<b>79.64</b>
	NMI	31.05	22.14	22.87	30.86	28.26	31.43	32.20		33.35	7.04	32.49	33.01	<b>41.43</b>
	ARI	28.10	19.55	22.30	27.61	29.84	30.64	31.41		34.52	7.04	29.95	30.44	<b>48.25</b>
	F1	58.88	61.45	65.01	62.85	68.23	68.10	68.94		70.33	44.54	67.12	67.35	<b>79.00</b>
A-Photo	ACC	27.22	47.22	53.44	67.66	75.96	76.82	79.94	77.02	75.15	77.48	77.19	80.71	<b>84.42</b>
	NMI	13.23	37.35	44.85	53.46	65.25	66.23	73.70	67.54	59.78	67.67	68.94	70.50	<b>74.99</b>
	ARI	5.50	18.59	31.21	42.74	58.12	58.28	63.69	58.05	56.13	58.48	60.20	66.54	<b>72.01</b>
	F1	23.96	46.71	50.66	60.32	69.87	71.25	73.82	72.60	72.85	72.22	71.23	73.09	<b>76.49</b>
CoraFull	ACC	26.27	31.92	26.67	32.38	34.35	37.51	38.80	OOM	36.46	41.89	38.51	39.45	<b>42.94</b>
	NMI	34.68	41.31	37.83	50.42	49.16	51.30	51.91		52.82	53.21	53.39	54.41	<b>55.83</b>
	ARI	9.35	16.89	22.60	20.64	22.60	24.46	25.25		24.78	24.23	26.53	27.45	<b>30.07</b>
	F1	22.57	27.77	22.14	27.82	26.96	31.22	31.68		29.78	32.98	30.90	31.95	<b>37.01</b>

Table 1: Clustering performance (%). The best result is in bold. Dink-Net\* denotes Dink-Net-NoFT.

## 4 Experiments

### Experimental Setup

To evaluate the performance of THESAURUS, we run the proposed method on nine attribute graph datasets, including Cora, Citeseer, Pubmed, Amazon-Photo (A-Photo), Cora-Full, ACM, DBLP, UAT, and Wiki. The baselines are K-means, DEC, GRACE (Zhu et al. 2020), SDCN, DFCN, DCRN, S<sup>3</sup>GC, SCGC, HSAN, and Dink-Net.

Our evaluation protocol follows that of the previous SOTA Dink-Net (Liu et al. 2023a). Besides Normalized Mutual Information (NMI) and Adjusted Rand Index (ARI), the metrics include Accuracy (ACC) and the Macro-F1 score (F1), computed after mapping the clusters to the ground-truth classes with linear assignment (Lovasz 1986; Crouse 2016). The F1 score, defined as the harmonic mean of precision and recall, balances the effects of false positives and false negatives. Meanwhile, ARI quantifies the number of true positive and true negative pairs and normalizes these values to ensure that the assessment is not influenced by variations in cluster sizes. Therefore, F1 and ARI are more effective than ACC and NMI on imbalanced data.

### Overall Performance

Part of the results are summarized in Table 1, with OOM indicating out-of-memory failures on one RTX 4090 GPU. And the rest results are presented in Tables 4 and 5 in the appendix. These results demonstrate that the proposed THESAURUS significantly outperforms existing methods across all datasets. And several key observations can be made.

**The contextual information from semantic prototypes is important.** Existing methods do not achieve sufficient cluster separability and are affected by the Uniform Effect and Cluster Assimilation. This results in suboptimal performance on clusters of varying sizes, particularly minority

clusters, adversely impacting F1 and ARI. In contrast, THESAURUS, utilizing semantic prototype contexts, achieves better distinction between synonymous nodes, leading to higher cluster separability. This effectively mitigates Uniform Effect and Cluster Assimilation, evidenced by THESAURUS’s F1 and ARI significantly surpassing existing methods. **Pretext task aligned with clustering is vital.** Contrastive clustering methods like S<sup>3</sup>GC and HSAN outperform DEC-style methods such as SDCN and DFCN, likely due to the implicit but insufficient alignment between InfoNCE and (spectral) clustering (HaoChen et al. 2021; Tan et al. 2023). DinkNet, which explicitly optimize towards clustering during finetune, is better aligned with clustering tasks, thus surpassing other baselines. Furthermore, THESAURUS representations are learned towards high cluster separability from start to finish, and thus far outperform all other methods. **The cluster information in structure matters.** Methods utilizing no graph structures, such as K-means and DEC, are unsuitable for graph clustering. Methods like HSAN and SCGC, which inject structure information with supervision signals, generally surpass those that only leverage structure through graph filters, such as SDCN and DAEGC. Furthermore, THESAURUS, which exhaustively exploits structural cluster information, surpasses all other methods.

### Class-wise Performance & Visualization

While Dink-Net’s finetune step mitigates the Uniform Effect in the majority cluster (cluster 3), it fails to address Cluster Assimilation in the minority cluster (cluster 6), as indicated by the experiments presented in Section Introduction. This section evaluates the class-wise performance of Dink-Net and THESAURUS, demonstrating that THESAURUS achieves high cluster separability and addresses the failure cases of Dink-Net. Fig. 3 visualizes the final representations

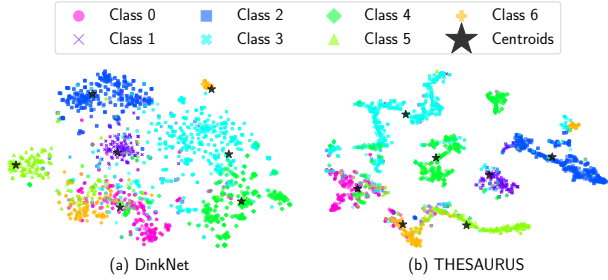


Figure 3: The t-SNE visualization on Cora

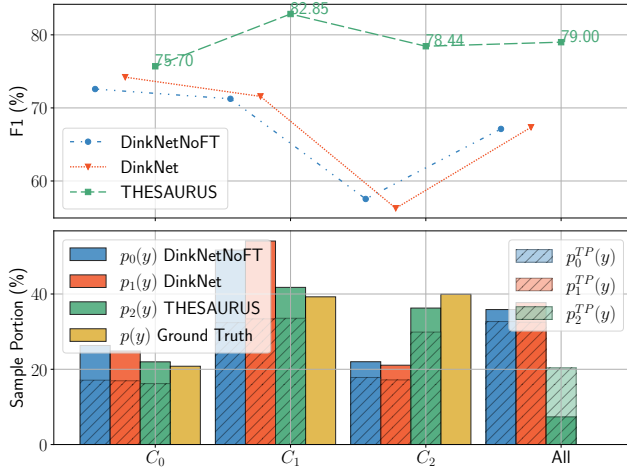


Figure 4: Dink-Net and THESAURUS on Pubmed. The **top** figure illustrates the F1 scores for each category, as well as the Macro-F1. The **bottom** shows the distribution of labels predicted by Dink-Net and THESAURUS, along with the ground-truth labels. It also presents the distribution of predicted labels for true-positive (TP) samples, denoted as  $p_i^{TP}(y)$ ,  $i \in \{0, 1, 2\}$ . The final set of bars shows the differences between the predicted and ground-truth distributions.

to cluster and the centroids on Cora using t-SNE (van der Maaten and Hinton 2008). Fig. 4 displays the class-wise performance on Pubmed, in a style like Fig. 1a.

Overall, the predicted label distribution of THESAURUS exhibits a lower deviation from the ground-truth distribution than that of DinkNet, as shown in Fig. 1a and Fig. 4. This suggests fewer mis-clustered nodes and corresponds with the higher cluster separability observed in THESAURUS’s representation space in Fig. 3, resulting in a Macro-F1 score up to 5.5 percentage points higher than Dink-Net on Cora and 11.65 points higher on Pubmed.

**Uniform Effect** Fig. 3a (right part) shows that Dink-Net does not separate clusters 3 and 4 as well as THESAURUS in Fig. 3b (top-left). So although Dink-Net reduces the Uniform Effect in cluster 3, the F1 score for cluster 4 is significantly lower than THESAURUS, as shown by Fig. 1a. Besides, Dink-Net’s slightly higher F1 score for cluster 3 comes at the cost of many false positives, which severely compromises the performance of other clusters (e.g., clus-

Datasets	Metrics	w/o $\mathbf{B}$	w/o $\mathbf{A}$	fixed $\nu$ & $\mathbf{B}$	Ours
Citeseer	ACC	70.27	71.54	71.24	<b>71.99</b>
	NMI	46.41	46.74	46.37	<b>47.37</b>
	ARI	47.82	47.92	47.57	<b>48.99</b>
	F1	65.02	65.75	65.68	<b>66.42</b>
Pubmed	ACC	75.73	78.35	75.84	<b>79.64</b>
	NMI	35.44	39.19	35.50	<b>41.43</b>
	ARI	39.25	45.47	40.62	<b>48.25</b>
	F1	75.27	77.81	75.39	<b>79.00</b>

Table 2: THESAURUS ablation. The best is in bold.

ters 4 and 0). In contrast, THESAURUS doesn’t favor the majority cluster and performs well overall.

**Cluster Assimilation** Fig. 3a (bottom-left) shows Dink-Net mixing clusters 0, 5, and 6, leading to a low F1 score of about 34% for cluster 6 due to samples being merged into other clusters. Fig. 3b shows that THESAURUS effectively separates these clusters, significantly mitigating Cluster Assimilation, as evidenced by a 27.29 percentage point higher F1 score for cluster 6 compared to Dink-Net (see Fig. 1a).

### Ablation Study

To validate the effectiveness of the complete prototype graph, we set the prototype graph  $\mathbf{B}$  as isolated graph  $\mathbf{I}_S$  (w/o  $\mathbf{B}$ ). To test the impact of structural information extraction, we replace  $\mathbf{A}_1$  and  $\mathbf{A}_2$  in FGW-OT with  $\mathbf{I}_N$  (w/o  $\mathbf{A}$ ). Additionally, to assess the effectiveness of the proposed momentum module, we set the momentum to 1 (fixed  $\nu$  &  $\mathbf{B}$ ). The results in Table 2 indicate that the complete prototype graph outperforms the isolated one, structural information extraction is effective, and momentum updates for the prototype graph and marginal is useful.

## 5 Conclusion

This work identifies challenges in prior deep graph clustering methods, particularly the Uniform Effect and Cluster Assimilation issues, which arise due to low cluster separability in the learned embedding space. To address these challenges, we propose a novel contrastive graph learning framework, THESAURUS. Our method 1) utilizes semantic prototypes to provide contextual information crucial for distinguishing similar nodes from different classes, 2) leverages a pretext task well-aligned with the downstream clustering task for better feature transferability, 3) takes GW-OT to thoroughly exploit the cluster information in graph structure, and 4) employs a momentum module for data adaptability. To achieve comprehensive information mining, cross-view alignments used in the pretext task are acquired via FGW-OT. This approach integrates designs 2 and 3 organically, providing an unified TSA module for better alignment. Experimental results strongly validate the effectiveness and superiority of THESAURUS compared to existing methods, showcasing its enhanced capability in achieving more accurate and efficient alignments.

## Ethics Statement

This research utilizes publicly available datasets and comparison methods, all of which are based on open-source code. No human participants or private data are involved in this study. All datasets used have been anonymized, and ethical guidelines regarding data usage have been strictly followed. We ensure that the methods used are transparent.

## Acknowledgements

The research is supported by the National Key Research and Development Program of China (2023YFB2703700), the National Natural Science Foundation of China (62176269), and the Guangzhou Science and Technology Program (2023A04J0314).

## References

- Alvarez-Melis, D.; and Jaakkola, T. 2018. Gromov-Wasserstein Alignment of Word Embedding Spaces. In Riloff, E.; Chiang, D.; Hockenmaier, J.; and Tsujii, J., eds., *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, 1881–1890. Brussels, Belgium: Association for Computational Linguistics.
- Bo, D.; Wang, X.; Shi, C.; Zhu, M.; Lu, E.; and Cui, P. 2020. Structural Deep Clustering Network. In *Proceedings of the Web Conference 2020*, Www '20, 1400–1410. New York, NY, USA: Association for Computing Machinery. ISBN 978-1-4503-7023-3.
- Cao, S.; Lu, W.; and Xu, Q. 2016. Deep Neural Networks for Learning Graph Representations. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 30.
- Caron, M.; Misra, I.; Mairal, J.; Goyal, P.; Bojanowski, P.; and Joulin, A. 2020. Unsupervised Learning of Visual Features by Contrasting Cluster Assignments. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, NIPS '20, 9912–9924. Red Hook, NY, USA: Curran Associates Inc. ISBN 978-1-7138-2954-6.
- Chen, D.; Lin, Y.; Li, W.; Li, P.; Zhou, J.; and Sun, X. 2020. Measuring and Relieving the Over-Smoothing Problem for Graph Neural Networks from the Topological View. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, 3438–3445.
- Chowdhury, S.; and Mémoli, F. 2019. The Gromov-Wasserstein Distance between Networks and Stable Network Invariants. *Information and Inference: A Journal of the IMA*, 8(4): 757–787.
- Courty, N.; Flamary, R.; Habrard, A.; and Rakotomamonjy, A. 2017. Joint Distribution Optimal Transportation for Domain Adaptation. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, NIPS'17, 3733–3742. Red Hook, NY, USA: Curran Associates Inc. ISBN 978-1-5108-6096-4.
- Crouse, D. F. 2016. On Implementing 2D Rectangular Assignment Algorithms. *IEEE Transactions on Aerospace and Electronic Systems*, 52(4): 1679–1696.
- Cui, G.; Zhou, J.; Yang, C.; and Liu, Z. 2020. Adaptive Graph Encoder for Attributed Graph Embedding. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, Kdd '20, 976–985. New York, NY, USA: Association for Computing Machinery. ISBN 978-1-4503-7998-4.
- Cuturi, M. 2013. Sinkhorn Distances: Lightspeed Computation of Optimal Transport. In Burges, C.; Bottou, L.; Welling, M.; Ghahramani, Z.; and Weinberger, K., eds., *Advances in Neural Information Processing Systems*, volume 26. Curran Associates, Inc.
- Devvrit, F.; Sinha, A.; Dhillon, I.; and Jain, P. 2022. S3GC: Scalable Self-Supervised Graph Clustering. In Koyejo, S.; Mohamed, S.; Agarwal, A.; Belgrave, D.; Cho, K.; and Oh, A., eds., *Advances in Neural Information Processing Systems*, volume 35, 3248–3261. Curran Associates, Inc.
- Gong, F.; Nie, Y.; and Xu, H. 2022. Gromov-Wasserstein Multi-modal Alignment and Clustering. In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management*, CIKM '22, 603–613. New York, NY, USA: Association for Computing Machinery. ISBN 978-1-4503-9236-5.
- HaoChen, J. Z.; Wei, C.; Gaidon, A.; and Ma, T. 2021. Provable Guarantees for Self-Supervised Deep Learning with Spectral Contrastive Loss. In Ranzato, M.; Beygelzimer, A.; Dauphin, Y.; Liang, P.; and Vaughan, J. W., eds., *Advances in Neural Information Processing Systems*, volume 34, 5000–5011. Curran Associates, Inc.
- Kipf, T. N.; and Welling, M. 2016. Variational Graph Auto-Encoders. In *NIPS Workshop on Bayesian Deep Learning*.
- Kipf, T. N.; and Welling, M. 2017. Semi-Supervised Classification with Graph Convolutional Networks. In *International Conference on Learning Representations*.
- Lee, J. D.; Lei, Q.; Saunshi, N.; and ZHUO, JIACHENG. 2021. Predicting What You Already Know Helps: Provable Self-Supervised Learning. In *Advances in Neural Information Processing Systems*, volume 34, 309–323. Curran Associates, Inc.
- Liu, Y.; Jin, M.; Pan, S.; Zhou, C.; Zheng, Y.; Xia, F.; and Yu, P. 2022a. Graph Self-Supervised Learning: A Survey. *IEEE Transactions on Knowledge and Data Engineering*, 1–1.
- Liu, Y.; Liang, K.; Xia, J.; Zhou, S.; Yang, X.; Liu, X.; and Li, S. Z. 2023a. Dink-Net: Neural Clustering on Large Graphs. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of ICML'23, 21794–21812. JMLR.org.
- Liu, Y.; Tu, W.; Zhou, S.; Liu, X.; Song, L.; Yang, X.; and Zhu, E. 2022b. Deep Graph Clustering via Dual Correlation Reduction. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, 7603–7611.
- Liu, Y.; Xia, J.; Zhou, S.; Yang, X.; Liang, K.; Fan, C.; Zhuang, Y.; Li, S. Z.; Liu, X.; and He, K. 2023b. A Survey of Deep Graph Clustering: Taxonomy, Challenge, Application, and Open Resource. arXiv:2211.12875.
- Liu, Y.; Yang, X.; Zhou, S.; Liu, X.; Wang, S.; Liang, K.; Tu, W.; and Li, L. 2023c. Simple Contrastive Graph Clustering.

- IEEE Transactions on Neural Networks and Learning Systems*, 1–12.
- Liu, Y.; Yang, X.; Zhou, S.; Liu, X.; Wang, Z.; Liang, K.; Tu, W.; Li, L.; Duan, J.; and Chen, C. 2023d. Hard Sample Aware Network for Contrastive Deep Graph Clustering. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, 8914–8922.
- Lloyd, SP. 1957. Least Squares Quantization in PCM. Technical Report RR-5497, Bell Lab, September 1957.
- Lovasz, L. 1986. *Matching Theory (North-Holland Mathematics Studies)*. GBR: Elsevier Science Ltd. ISBN 0-444-87916-1.
- Lu, Y.; Cheung, Y.-M.; and Tang, Y. Y. 2021. Self-Adaptive Multiprototype-Based Competitive Learning Approach: A k-Means-Type Algorithm for Imbalanced Data Clustering. *IEEE Transactions on Cybernetics*, 51(3): 1598–1612.
- MacQueen, J.; et al. 1967. Some Methods for Classification and Analysis of Multivariate Observations. In *Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability*, volume 1, 281–297. Oakland, CA, USA.
- Mémoli, F. 2011. Gromov–Wasserstein Distances and the Metric Approach to Object Matching. *Foundations of Computational Mathematics*, 11(4): 417–487.
- Monge, G. 1781. Mémoire Sur La Théorie Des Déblais et Des Remblais. *Mem. Math. Phys. Acad. Royale Sci.*, 666–704.
- Nguyen, K.; Nong, H.; Nguyen, V.; Ho, N.; Osher, S.; and Nguyen, T. 2023. Revisiting Over-Smoothing and over-Squashing Using Ollivier-Ricci Curvature. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *ICML’23*, 25956–25979. Honolulu, Hawaii, USA: JMLR.org.
- Pan, S.; Hu, R.; Long, G.; Jiang, J.; Yao, L.; and Zhang, C. 2018. Adversarially Regularized Graph Autoencoder for Graph Embedding. In *Proceedings of the 27th International Joint Conference on Artificial Intelligence, IJCAI’18*, 2609–2615. Stockholm, Sweden: AAAI Press. ISBN 978-0-9992411-2-7.
- Peng, Z.; Liu, H.; Jia, Y.; and Hou, J. 2021. Attention-Driven Graph Clustering Network. In *Proceedings of the 29th ACM International Conference on Multimedia, Mm ’21*, 935–943. New York, NY, USA: Association for Computing Machinery. ISBN 978-1-4503-8651-7.
- Peyré, G.; and Cuturi, M. 2019. Computational Optimal Transport: With Applications to Data Science. *Foundations and Trends® in Machine Learning*, 11(5-6): 355–607.
- Shrivastava, A.; Selvaraju, R. R.; Naik, N.; and Ordonez, V. 2023. CLIP-Lite: Information Efficient Visual Representation Learning with Language Supervision. In *Proceedings of The 26th International Conference on Artificial Intelligence and Statistics*, 8433–8447. PMLR.
- Tan, Z.; Zhang, Y.; Yang, J.; and Yuan, Y. 2023. Contrastive Learning Is Spectral Clustering on Similarity Graph. In *The Twelfth International Conference on Learning Representations*.
- Titouan, V.; Courty, N.; Tavenard, R.; and Flamary, R. 2019. Optimal Transport for Structured Data with Application on Graphs. In *International Conference on Machine Learning*, 6275–6284. PMLR.
- Tolstikhin, I.; Bousquet, O.; Gelly, S.; and Schoelkopf, B. 2018. Wasserstein Auto-Encoders. In *International Conference on Learning Representations*.
- Tu, W.; Zhou, S.; Liu, X.; Guo, X.; Cai, Z.; Zhu, E.; and Cheng, J. 2021. Deep Fusion Clustering Network. In *Proceedings of the Thirty-Fifth AAAI Conference on Artificial Intelligence*, 9978–9987.
- van den Oord, A.; Li, Y.; and Vinyals, O. 2019. Representation Learning with Contrastive Predictive Coding. arXiv:1807.03748.
- van der Maaten, L.; and Hinton, G. 2008. Visualizing Data Using T-SNE. *Journal of Machine Learning Research*, 9(86): 2579–2605.
- Villani, C. 2009. *Optimal Transport: Old and New*, volume 338 of *Grundlehren Der Mathematischen Wissenschaften*. Berlin, Heidelberg: Springer Berlin Heidelberg. ISBN 978-3-540-71049-3 978-3-540-71050-9.
- Wang, C.; Pan, S.; Hu, R.; Long, G.; Jiang, J.; and Zhang, C. 2019. Attributed Graph Clustering: A Deep Attentional Embedding Approach. In *Proceedings of the 28th International Joint Conference on Artificial Intelligence, IJCAI’19*, 3670–3676. AAAI Press. ISBN 978-0-9992411-4-1.
- Wang, C.; Pan, S.; Long, G.; Zhu, X.; and Jiang, J. 2017. Mgae: Marginalized Graph Autoencoder for Graph Clustering. In *Proceedings of the 2017 ACM on Conference on Information and Knowledge Management*, 889–898.
- Wang, S.; Yang, J.; Yao, J.; Bai, Y.; and Zhu, W. 2024. An Overview of Advanced Deep Graph Node Clustering. *IEEE Transactions on Computational Social Systems*, 11(1): 1302–1314.
- Wei, F.; Gao, Y.; Wu, Z.; Hu, H.; and Lin, S. 2021. Aligning Pretraining for Detection via Object-Level Contrastive Learning. In *Advances in Neural Information Processing Systems*.
- Xie, J.; Girshick, R.; and Farhadi, A. 2016. Unsupervised Deep Embedding for Clustering Analysis. In Balcan, M. F.; and Weinberger, K. Q., eds., *Proceedings of the 33rd International Conference on Machine Learning*, volume 48 of *Proceedings of Machine Learning Research*, 478–487. New York, New York, USA: PMLR.
- Xiong, H.; Wu, J.; and Chen, J. 2009. K-Means Clustering Versus Validation Measures: A Data-Distribution Perspective. *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)*, 39(2): 318–331.
- Xu, H.; Luo, D.; Zha, H.; and Duke, L. C. 2019. Gromov-Wasserstein Learning for Graph Matching and Node Embedding. In Chaudhuri, K.; and Salakhutdinov, R., eds., *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, 6932–6941. PMLR.
- Zhu, Y.; Xu, Y.; Yu, F.; Liu, Q.; Wu, S.; and Wang, L. 2020. Deep Graph Contrastive Representation Learning. arXiv:2006.04131.