

---

# DIGRAC: Digraph Clustering Based on Flow Imbalance

---

Anonymous Author(s)

Anonymous Affiliation

Anonymous Email

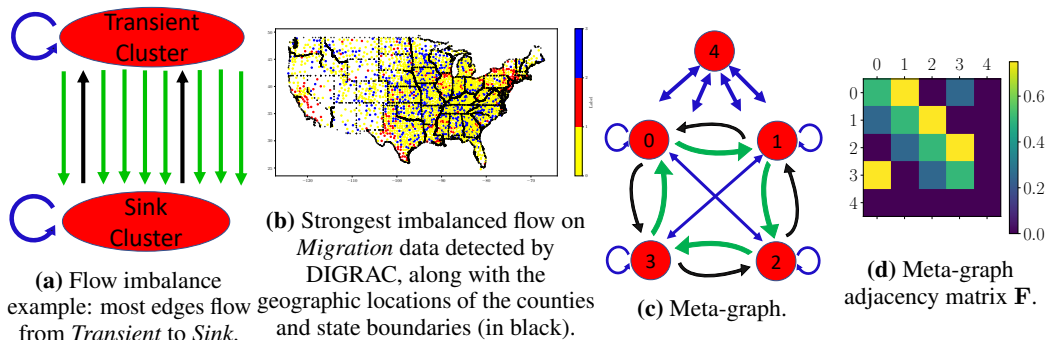
## Abstract

Node clustering is a powerful tool in the analysis of networks. We introduce a graph neural network framework to obtain node embeddings for directed networks in a self-supervised manner, including a novel probabilistic imbalance loss, which can be used for network clustering. Here, we propose *directed flow imbalance* measures, which are tightly related to directionality, to reveal clusters in the network even when there is no density difference between clusters. In contrast to standard approaches in the literature, in this paper, directionality is not treated as a nuisance, but rather contains the main signal. DIGRAC optimizes directed flow imbalance for clustering without requiring label supervision, unlike existing graph neural network methods, and can naturally incorporate node features, unlike existing spectral methods. Extensive experimental results on synthetic data, in the form of directed stochastic block models, and real-world data at different scales, demonstrate that our method, based on flow imbalance, attains state-of-the-art results on directed graph clustering when compared against 10 state-of-the-art methods from the literature, for a wide range of noise and sparsity levels, graph structures and topologies, and even outperforms supervised methods.

## 1 Introduction

Revealing an underlying community structure of *directed* networks (*digraphs*) is an important problem in many applications, see for example [1] and [2], such as detecting influential social groups [3] and analyzing migration patterns [4]. While most existing methods that could be applied to directed clustering use local edge densities as main signal and directionality (i.e., edge orientation) as additional signal, we argue that even in the absence of any edge density differences, directionality can play a vital role in directed clustering as it can reveal latent properties of network flows. The underlying intuition is that homogeneous clusters of nodes form *meta-nodes* in a *meta-graph*, with the meta-graph directing the flow between clusters; directed core-periphery structure is such an example [5]. Loosely speaking, a meta-node is a collection of nodes, and a meta-graph is a graph on such meta-nodes, with weighted edges collecting the overall sum of edge weights between the meta-nodes. Fig. 1(a) is an example of flow imbalance between two clusters, here on an unweighted network for simplicity: while 80% of the edges flow from the *Transient* cluster to the *Sink* cluster, only 20% flow in the other direction. As a real-world example, Fig. 1(b) shows the strongest flow imbalances between clusters detected by our method in a network of US migration flow [4]; most edges flow from the red cluster (label 1) to the blue one (label 2). Figures 1(c-d) show examples on a synthetic meta-graph. We could also think of a social network in which a set of fake accounts  $\mathcal{A}$  have been created, and these target another subset  $\mathcal{B}$  of real accounts by sending them messages. Most likely, there would be many more messages from  $\mathcal{A}$  to  $\mathcal{B}$  than from  $\mathcal{B}$  to  $\mathcal{A}$ , hinting that  $\mathcal{A}$  is most likely comprised of fake accounts.

Thus, instead of finding relatively dense groups of nodes in digraphs with a relatively small amount of flow between the groups, as in [6–11], our main goal is to recover clusters with *strongly imbalanced flow* among them, in the spirit of [12, 13], where directionality is the main signal. This task is not addressed by most methods for node clustering in digraphs, including community detection methods. Those methods that do lay emphasis on directionality are usually spectral methods, for which incorporating features is non-trivial, or graph neural network (GNN) methods that require labeling information. An exception is the network community detection method InfoMap [14] which uses directed random walks; however, it still relies on some edge density information within clusters,



**Figure 1:** Visualization of cut flow imbalance and meta-graph: (a) 80% of edges flow from *Transient* to *Sink*, while 20% of edges flow in the opposite direction; (b) top pair imbalanced flow on *Migration* data [4]: most edges flow from red (1) to blue (2); (c) & (d) are for a Directed Stochastic Block Model with a cycle meta-graph with ambient nodes, for a total of 5 clusters. Most edges flow in direction  $0 \rightarrow 1 \rightarrow 2 \rightarrow 3 \rightarrow 0$ , while few flow in the opposite direction. Cluster 4 is the ambient cluster. In (a) and (c), blue lines indicate flows with random, equally likely directions; these flows do not exist in the meta-graph adjacency matrix  $\mathbf{F}$ . For (d), the lighter the color, the stronger the flow.

45 as a walk is more likely to happen when the density is higher. [15] and [16] employ Markov chains,  
 46 but we pick InfoMap as a representative for methods based on information theory/Markov chains.

47 Here we introduce DIGRAC, a GNN framework to obtain node embeddings for clustering digraphs  
 48 (allowing weighted edges and self-loops but no multiple edges). In a self-supervised manner, a  
 49 novel *probabilistic imbalance loss* is proposed to act on the digraph induced by all training nodes.  
 50 The global imbalance score, one minus whom is the self-supervised loss function, is aggregated  
 51 from pairwise normalized cut imbalances. The method is end-to-end in combining embedding  
 52 generation and clustering without an intermediate step. To the best of our knowledge, this is the first  
 53 GNN method which derives node embeddings for digraphs that directly maximizes flow imbalance  
 54 between pairs of clusters. With an emphasis on the use of a direction-based flow imbalance objective,  
 55 experimental results on synthetic data and real-world data at different scales demonstrate that our  
 56 method can achieve leading performance for a wide range of network densities and topologies.

57 DIGRAC’s main novelty is the ability to cluster based on direction-based flow imbalance, instead  
 58 of using classical criteria such as maximizing relative densities within clusters. Compared with prior  
 59 methods that focus on directionality, DIGRAC can easily consider node features and also does not  
 60 require known node clustering labels. DIGRAC complements existing approaches in various aspects:  
 61 (1) Our results show that DIGRAC complements classical community detection by detecting alterna-  
 62 tive patterns in the data, such as meta-graph structures, which are otherwise not detectable by existing  
 63 methods. This aspect of detecting novel structures in directed graphs has also been emphasized  
 64 in [12]. (2) DIGRAC complements existing spectral methods, through the possibility of including  
 65 exogenous information, in the form of node-level features or labels, thus borrowing their strength.  
 66 (3) DIGRAC complements existing GNN methods by introducing an imbalance-based objective.  
 67 (4) DIGRAC introduces imbalance measures for evaluation when ground-truth is unavailable.

68 DIGRAC’s applicability extends beyond settings where the input data is a digraph: with time series  
 69 data as input, the digraph construction mechanism can accommodate any procedure that encodes a  
 70 pairwise directional association between the corresponding time series, such as lead-lag relationships  
 71 and Granger causality [17], with applications such as in the analysis of information flow in brain  
 72 networks [18], biology [19], finance [20, 21] and earth sciences [22]. DIGRAC could also facilitate  
 73 tasks in ranking and anomaly detection, as it allows one to extrapolate from *local* pairwise (directed)  
 74 interactions to a *global* structure inference, in the high-dimensional low signal-to-noise ratio regime.

75 **Main contributions.** Our main contributions are as follows. •(1) We propose a GNN framework for  
 76 self-supervised end-to-end node clustering on (possibly attributed and weighted) digraphs explicitly  
 77 taking into account the directed flow imbalance. •(2) We propose a family of probabilistic global  
 78 imbalance scores to serve as the self-supervised loss function and evaluation objective, including one  
 79 based on hypothesis testing for directionality signal. To the best of our knowledge, this is the first  
 80 method directly maximizing flow imbalance for node clustering in digraphs using GNNs. •(3) We  
 81 extend our method to the semi-supervised setting when label information is available.

## 82 2 Related work

83 Directed clustering has been explored by non-GNN methods. [23] performs directed clustering that  
 84 hinges on symmetrizations of the adjacency matrix, but is not scalable as it requires large matrix mul-  
 85 tiplications. [24] proposes a spectral co-clustering algorithm for asymmetry discovery that relies on  
 86 in-degree and out-degree. Whenever direction is the sole information, such as in a complete network  
 87 with lead-lag structure derived from time series [20], a purely degree-based method cannot detect the  
 88 clusters. While [25] produces two partitions of the node set, one based on out-degree and one based on  
 89 in-degree, our partition simultaneously takes both directions into account. The *directed graph Lapla-*  
 90 *cians* introduced by [2] are only applicable to strongly connected digraphs, which is rarely the case  
 91 in sparse networks arising in applications. InfoMap by [14] assumes that there is a “map” underlying  
 92 the network, similar to a meta-graph in DIGRAC. InfoMap aims to minimize the expected description  
 93 length of a random walk and is recommended for networks where edges encode patterns of movement  
 94 among nodes. While related to DIGRAC, InfoMap still relies on some amount of density-based signal  
 95 being present within each of the modules. [12] seeks to uncover clusters characterized by a strongly  
 96 imbalanced flow circulating among them, based on eigenvectors of the Hermitian matrix  $(\mathbf{A} - \mathbf{A}^T) \cdot i$ ,  
 97 where  $\mathbf{A}$  is the (normalized) adjacency matrix and  $i$  the imaginary unit. [12] is a purely spectral-based  
 98 method and is not able to naturally incorporate any available node features or label information; in  
 99 contrast, DIGRAC is a GNN-based method that is naturally able to account for such information.  
 100 Moreover, [12] is not driven by an optimization function, but only proposes evaluation metrics that cap-  
 101 ture the imbalance of the pairs of clusters. In contrast, inspired by [12], in DIGRAC a family of novel  
 102 imbalance loss functions is proposed, with a probabilistic interpretation, rendering DIGRAC a fully  
 103 trainable end-to-end pipeline. Furthermore, the rich class of imbalance evaluation and training objec-  
 104 tives/losses proposed in this paper go far beyond the evaluation metrics considered in [12]. [13] uncov-  
 105 ers higher-order structural information among clusters in digraphs, while maximizing the imbalance  
 106 of the edge directions, but its definition of the flow ratio restricts the underlying meta-graph to a path.

107 GNNs have been applied to digraph node classification, which is similar to digraph clustering  
 108 but requires known clustering labels. [26] uses first and second-order proximity, constructs three  
 109 Laplacians, but the method is space and speed-inefficient. [27] simplifies [26], builds a directed  
 110 Laplacian based on PageRank, and aggregates information dependent on higher-order proximity.  
 111 Building on [12, 28], [29] constructs a Hermitian matrix that encodes undirected geometric structure  
 112 in the magnitude of its entries, and directional information in their phase. [30] introduces a digraph  
 113 data augmentation method called *Laplacian perturbation* and conducts digraph contrastive learning.  
 114 [31] proposes a spectral-based graph convolution network for digraphs, yet is restricted to strongly  
 115 connected digraphs that are usually not realistic. [32] utilizes convolution-like anisotropic filters  
 116 based on local subgraph structures (motifs) for semi-supervised node classification tasks in digraphs,  
 117 but relies on pre-defined structures and fails to handle complex networks.

118 In particular, [26, 27, 29, 30, 32] all require known labels, which are not generally available for  
 119 real-world data. [2, 12, 13, 23, 24] could not trivially incorporate node attributes or node labels. In  
 120 contrast, we propose an efficient GNN-based method that maximizes a probabilistic flow imbalance  
 121 objective, in a self-supervised manner, and which can naturally analyze attributed weighted digraphs.

122 To avoid potential misunderstanding, we briefly mention several related works that we are aware of,  
 123 but do not compare against in our experiments in the main text. While DIGRAC addresses the task  
 124 of partitioning the nodes into disjoint sets, [33] locates a certain community within a network. In  
 125 particular, [33] proposes a local algorithm while this paper proposes a global one. OSLOM by [34]  
 126 is very flexible but based on a density heuristic and hence a comparison to DIGRAC on networks  
 127 without density signal would not be fair to begin with. [35] introduces directionality in the Louvain  
 128 algorithm. This algorithm optimizes a modularity-type function that compares the number of edges  
 129 within communities to the expected number of edges under a specified model. It is thus an approach  
 130 that aims to find denser-than-expected groups of vertices. When all groups have the same density, as  
 131 in our synthetic data sets, and the only structure lies in the directionality of the edges, this method  
 132 simply cannot be expected to perform well. The Leiden algorithm in [36] also builds on the Louvain  
 133 method, again optimizing a modularity-type function that compares the number of edges within  
 134 communities to the expected number of edges under a specified model. It is a powerful method for  
 135 that task, but cannot be fairly compared to DIGRAC which is tailored to find imbalances.

136 We also do not compare DIGRAC against graph pooling methods [37], which are inspired by pooling  
 137 in CNNs and developed to discard information which is superfluous for the task at hand, as a partition

of the nodes which can be interpreted as clustering is only a byproduct. Moreover, graph pooling methods are usually developed only for undirected networks. While graph matching as in [38–40] and [41] can be viewed as a clustering method of networks, matching the graph of interest to a disconnected graph by connecting each node in the observed graph with an isolated node of the disconnected graph, this approach is not developed for directed networks. The underlying idea of these papers is complementary to the meta-graph idea which underpins DIGRAC; in the meta-graph, the components are connected, and estimating the directionality of these connections is the main focus. Hence this work addresses a very different task. We emphasize that these are all excellent methods, but they address different objectives and tasks. As confirmed by our experiments in Appendix (App.) E, comparing these methods to DIGRAC is not appropriate. DIGRAC is tailored to detect an imbalance signal in directed networks, and such a signal cannot be present in an undirected network. As it is based on imbalance, DIGRAC will not be able to detect a signal in an undirected network, thus rendering it not applicable to undirected networks.

### 3 The DIGRAC framework

**Problem definition.** Denote a (possibly weighted) digraph with node attributes as  $\mathcal{G} = (\mathcal{V}, \mathcal{E}, w, \mathbf{X})$ , with  $\mathcal{V}$  the set of nodes,  $\mathcal{E}$  the set of directed edges or links, and  $w \in [0, \infty)^{|\mathcal{E}|}$  the set of edge weights.  $\mathcal{G}$  may have self-loops, but no multiple edges. The number of nodes is  $n = |\mathcal{V}|$ , and  $\mathbf{X} \in \mathbb{R}^{n \times d_{\text{in}}}$  is a matrix whose rows encode the nodes’ attributes. Such a network can be represented by the attribute matrix  $\mathbf{X}$  and the adjacency matrix  $\mathbf{A} = (A_{ij})_{i,j \in \mathcal{V}}$ , with  $A_{ij} = 0$  if no edge exists from  $v_i$  to  $v_j$ ; if there is an edge  $e$  from  $v_i$  to  $v_j$ , we set  $A_{ij} = w_e$ , the edge weight.

Digraphs often lend themselves to interpreting weighted directed edges as flows, with a meta-graph on clusters of vertices describing the overall flow directions; see Fig. 1. A clustering is a partition of the set of nodes into  $K$  disjoint sets (clusters)  $\mathcal{V} = \mathcal{C}_0 \cup \mathcal{C}_1 \cup \dots \cup \mathcal{C}_{K-1}$  (ideally,  $K \geq 2$ ). Intuitively, nodes within a cluster should be similar to each other with respect to flow directions, while nodes across clusters should be dissimilar. In a self-supervised setting, only the number of clusters  $K$  is given. In a semi-supervised setting, for each of the  $K$  clusters, a fraction set  $\mathcal{V}^{\text{seed}} \subseteq \mathcal{V}^{\text{train}} \subset \mathcal{V}$  of the set  $\mathcal{V}^{\text{train}}$  of all training nodes is selected to serve as the set of seed nodes, for which the cluster membership labels are known before training. The goal of semi-supervised clustering is to assign each node  $v \in \mathcal{V}$  to a cluster containing some known seed nodes, without knowledge of the underlying flow meta-graph. The corresponding self-supervised clustering task does not use seed nodes.

#### 3.1 Self-supervised loss for clustering

Our self-supervised loss function is inspired by [12], aiming to cluster the nodes by maximizing a normalized form of cut imbalance across clusters. We first define probabilistic versions of cuts, imbalance flows, and probabilistic volumes. For  $K$  clusters, the *assignment probability matrix*  $\mathbf{P} \in \mathbb{R}^{n \times K}$  has as row  $i$  the probability vector  $\mathbf{P}_{(i,:)} \in \mathbb{R}^K$  with entries denoting the probabilities of each node to belong to each cluster; its  $k^{\text{th}}$  column is denoted by  $\mathbf{P}_{(:,k)}$ .

•  $\forall k, l \in \{0, \dots, K-1\}$  where  $K \geq 2$ , the **probabilistic cut** from cluster  $\mathcal{C}_k$  to  $\mathcal{C}_l$  is defined as

$$W(\mathcal{C}_k, \mathcal{C}_l) = \sum_{i,j} \mathbf{A}_{i,j} \cdot \mathbf{P}_{i,k} \cdot \mathbf{P}_{j,l} = (\mathbf{P}_{(:,k)})^T \mathbf{A} \mathbf{P}_{(:,l)}.$$

• The **imbalance flow** between  $\mathcal{C}_k$  and  $\mathcal{C}_l$  is defined as  $|W(\mathcal{C}_k, \mathcal{C}_l) - W(\mathcal{C}_l, \mathcal{C}_k)|$ .

For interpretability and ease of comparison, we normalize the imbalance flows to obtain an imbalance score with values in  $[0, 1]$  as follows (we defer additional details to App. B.2).

• The **probabilistic volume** for cluster  $\mathcal{C}_k$  is defined as

$$\begin{aligned} \text{VOL}(\mathcal{C}_k) &= \text{VOL}^{\text{(out)}}(\mathcal{C}_k) + \text{VOL}^{\text{(in)}}(\mathcal{C}_k) \\ &= \sum_{i,j} (\mathbf{A}_{j,i} + \mathbf{A}_{i,j}) \cdot \mathbf{P}_{j,k} \end{aligned}$$

Then  $\text{VOL}(\mathcal{C}_k) \geq W(\mathcal{C}_k, \mathcal{C}_l)$  for all  $l = 1, \dots, K-1$  and

$$\min(\text{VOL}(\mathcal{C}_k), \text{VOL}(\mathcal{C}_l)) \geq |W(\mathcal{C}_k, \mathcal{C}_l) - W(\mathcal{C}_l, \mathcal{C}_k)|. \quad (1)$$

The imbalance term, which is used in most of our experiments, denoted  $\text{CI}^{\text{vol\_sum}}$ , is defined as

$$\text{CI}^{\text{vol\_sum}}(k, l) = 2 \frac{|W(\mathcal{C}_k, \mathcal{C}_l) - W(\mathcal{C}_l, \mathcal{C}_k)|}{\text{VOL}(\mathcal{C}_k) + \text{VOL}(\mathcal{C}_l)} \in [0, 1]. \quad (2)$$

183 In particular for  $K = n$ , every node is a single cluster, and  $\text{CI}^{\text{vol\_sum}}(k, l) = 1$ , but then the partition is  
 184 not informative. The aim is to find a partition that maximizes the imbalance flow under the constraint  
 185 that the partition has at least two sets, to capture groups of nodes that could be viewed as representing  
 186 clusters in the meta-graph. The normalization by the volumes penalizes partitions that put most nodes  
 187 into a single cluster. The range  $[0, 1]$  follows from Eq. (1). Other variants are discussed in App. B.3.

188 To obtain a **global probabilistic imbalance score**, based on  $\text{CI}^{\text{vol\_sum}}$  from Eq. (2), we average over  
 189 pairwise imbalance scores of different pairs of clusters. Since the scores discussed are symmetric and  
 190 the cut difference before taking absolute value is skew-symmetric, we only need to consider the pairs  
 191 in the set  $\mathcal{T} = \{(C_k, C_l) : 0 \leq k < l \leq K - 1, k, l \in \mathbb{Z}\}$ .

192 A naive approach, which we call the **“naive”** variant, considers all possible  $\binom{K}{2}$  pairwise cut  
 193 imbalance values. However, due to potentially high noise levels in certain data sets, one may only be  
 194 interested in pairs that are not just noise but exhibit true signals. To this end, we introduce a **“std”**  
 195 variant, which only considers pairwise cut imbalance values that are 3 standard deviations away  
 196 from the observed purely noisy imbalance values; the standard deviation is calculated under the null  
 197 hypothesis that the between-cluster relationship has no direction preference, i.e.  $\mathbf{F}_{k,l} = \mathbf{F}_{l,k}$  (entries  
 198 of the meta-graph adjacency matrix  $\mathbf{F}$  to be introduced later in this section), as follows.

199 Suppose two clusters  $C_k$  and  $C_l$  have only noisy links between them, with no edge in the meta-graph  
 200  $\mathbf{F}$ , i.e.  $\mathbf{F}_{kl} = 0$ . Assume also that the underlying network is fixed in terms of the number of nodes  
 201 and locations of edges; the only randomness stems from the direction the edges. Then we can provide  
 202 the following theoretical guarantee.

203 **Proposition 1.** *Suppose that  $C_k$  and  $C_l$  are two clusters of  $n_k$  and  $n_l$  nodes, respectively, with  $m(k, l)$   
 204 edges between them, edge weights  $w_{ij} = w_{ji} \in [0, 1]$  and edge direction drawn independently  
 205 at random with equal probability  $\frac{1}{2}$  for each direction. We assume that the edge weights satisfy  
 206  $\max_e |w_e| (\sum_e w_e^2)^{-\frac{1}{2}} = o(m(k, l))$ . Then  $W(C_k, C_l) - W(C_l, C_k)$  is approximately normally  
 207 distributed with mean 0 and variance  $\|w\|^2$  as  $m(k, l) \rightarrow \infty$ .*

208 A consequence of Proposition 1, which is proved in App. B.1, is that under its assumptions, ap-  
 209 proximately 99.7 % of the observations fall within 3 standard deviations from 0. While Proposition  
 210 1 makes many assumptions and ignores reciprocal edges, the resulting threshold is still a useful  
 211 guideline for restricting attention to pairwise imbalance values which are very likely to capture a true  
 212 signal. In particular, we use it as motivation for our “std” variant to pick cluster pairs from  $\mathcal{T}$  that  
 213 satisfy  $(W(C_k, C_l) - W(C_l, C_k))^2 > 9(W(C_k, C_l) + W(C_l, C_k))$ .

214 As we are mainly concerned about the top pairs (i.e., those exhibiting the largest imbalance flow),  
 215 another option is the **“sort”** variant, which selects the largest  $\beta$  pairwise cut imbalance values,  
 216 where  $\beta$  is half of the number of nonzero entries in the off-diagonal entries of the meta-graph  
 217 adjacency matrix  $\mathbf{F}$ , if the meta-graph is known or can be approximated. For example, for a “cycle”  
 218 meta-graph with three clusters and no ambient nodes,  $\beta = 3$ . When the meta-graph is a “path”  
 219 with three clusters and ambient nodes, then  $\beta = 1$ . When considering the “sort” variant, with  
 220  $\mathcal{T}(\beta) = \{(C_k, C_l) \in \mathcal{T} : \text{CI}^{\text{vol\_sum}}(k, l) \text{ is among the top } \beta \text{ values}\}$ , where  $1 \leq \beta \leq \binom{K}{2}$ , we set

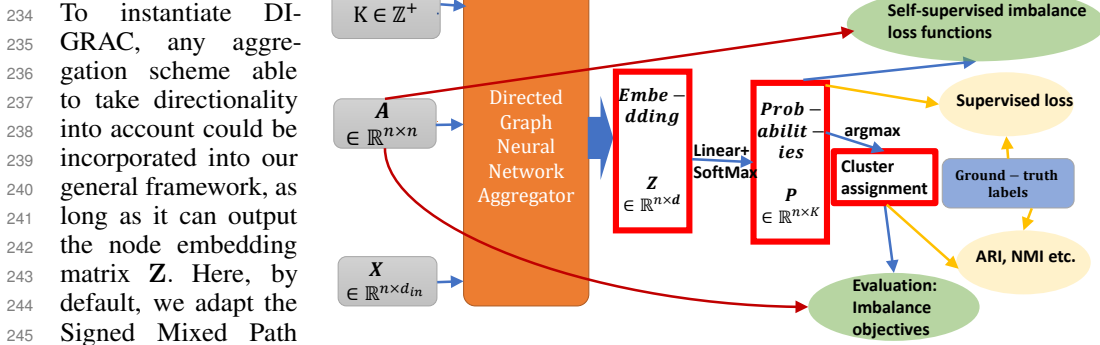
$$\mathcal{O}_{\text{vol\_sum}}^{\text{sort}} = \frac{1}{\beta} \sum_{(C_k, C_l) \in \mathcal{T}(\beta)} \text{CI}^{\text{vol\_sum}}(k, l), \quad \text{and} \quad \mathcal{L}_{\text{vol\_sum}}^{\text{sort}} = 1 - \mathcal{O}_{\text{vol\_sum}}^{\text{sort}}, \quad (3)$$

221 as the corresponding loss function. Definitions of meta-graph structures are discussed in Section 4.1.  
 222 For the other variants, the corresponding scores and loss functions are defined analogously. We apply  
 223 the “std” variant when we have no prior knowledge on the meta-graph structure during training, and  
 224 the “sort” variant when we have information of the number of pairs to count.  
 225

226 When using the “std” variant for training, for the initial 50 epochs, we apply the “sort” variant with  
 227  $\beta = 3$  for a reasonable starting clustering probability matrix for training, as otherwise during the initial  
 228 training epochs possibly no pairs could be picked out. During the epochs actually utilizing this “std”  
 229 variant, if no pairs could be picked out, we temporarily switch to the “naive” variant for that epoch.

230 Regarding complexity, the objective mainly contains matrix-vector multiplications and element-wise  
 231 matrix divisions, which are at most quadratic in the number of nodes, but usually faster with our  
 232 sparsity-aware implementation.

### 233 3.2 Instantiation of DIGRAC



**Figure 2:** DIGRAC overview: from feature matrix  $X$ , adjacency matrix  $A$  and number of clusters  $K$ , we first apply a directed GNN aggregator to obtain the node embedding matrix  $Z$ , then apply a linear layer followed by a unit softmax function to get the probability matrix  $P$ . Applying  $\text{argmax}$  on each row of  $P$  yields node cluster assignments. Green circles involves our proposed imbalance objective, while the yellow circles can only be used when ground-truth labels are provided.

## 4 Experiments

In our synthetic experiments, when by design ground truth is available, performance is assessed by the Adjusted Rand Index (ARI) [43]. Normalized Mutual Information (NMI) results give almost the same ranking for the best-performing methods as the ARI, with an average Kendall tau value of 83.8% and standard deviation 24.9%, for pairwise ranking comparison, on the methods compared in our experiments. We do not focus on NMI in the main text due to its shortcomings [44], see also App. C.9.

Clustering tasks will have different ground truths, depending on the pattern they are trying to detect. Many network clustering methods focus on detecting relatively dense clusters, and try to optimize classical network clustering measures, such as directed modularity or partition density. Ground truth for these clustering algorithms then relates to relatively densely connected subgroups in the data. DIGRAC is a novel method that addresses a novel task, namely that of detecting flow imbalances. To the best of our knowledge, real-world data sets with ground-truth flow imbalances are not available to date, and hence we introduce normalized imbalance scores to evaluate clustering performance based on flow imbalance. As ARI and NMI require ground-truth labels, they thus cannot be applied to the available real-world data sets. To address this shortcoming, for the real-world data sets, in Table 1, we include three performance measures which we introduce in the paper, and the appendix contains an additional 11 performance measures. Implementation details are provided in App. C. Anonymized codes and preprocessed data are available at <https://anonymous.4open.science/r/DIGRAC>.

We compare DIGRAC against the most recent related methods from the literature for clustering digraphs. The 10 methods are • (1) InfoMap [14], • (2) Bibliometric and • (3) Degree-discounted introduced in [23], • (4) DI\_SIM [24], • (5) Herm and • (6) Herm\_sym introduced in [12], • (7) MagNet [29], • (8) DGCN [26], • (9) DiGCN [27], and • (10) DiGCL [30]. The abbreviations of these methods, when reported in the numerical experiments, are InfoMap, Bi\_sym, DD\_sym, DISG\_LR, Herm, Herm\_sym, MagNet, DGCN, DiGCN, DiGCL, respectively. DGCN is the least efficient method in terms of speed and space complexity, followed by DiGCN which involves the so-called *inception blocks*. We use the same hyperparameter settings stated in these papers. Methods (7), (8), (9), (10) are GNN methods which are trained with 80% nodes under label supervision, while all the other methods are trained without label supervision. DIGRAC further restricts itself to be trained on the subgraph induced by only the training nodes. All methods are designed for directed graphs, and all except Infomap require  $K$  to be known. Runtime comparison is provided in App. C.2, illustrating that DIGRAC is among the fastest among competing GNNs. Implementation details for competitors are provided in App. C.7.

**Table 1:** Performance comparison on real-world data sets. The best is marked in **bold red** and the second best is marked in underline blue. The objectives are defined in Section 3.1.

Metric	Data set	InfoMap	Bi_sym	DD_sym	DISG_LR	Herm	Herm_rw	DIGRAC
$\mathcal{O}_{\text{vol\_sum}}^{\text{sort}}$	<i>Telegram</i>	0.04±0.00	<u>0.21±0.0</u>	<u>0.21±0.0</u>	<u>0.21±0.01</u>	0.2±0.01	0.14±0.0	<b>0.32±0.01</b>
	<i>Blog</i>	0.07±0.00	<u>0.07±0.0</u>	0.0±0.0	0.05±0.0	<u>0.37±0.0</u>	0.0±0.0	<b>0.44±0.0</b>
	<i>Migration</i>	N/A	0.03±0.00	0.01±0.00	0.02±0.00	<u>0.04±0.00</u>	0.02±0.00	<b>0.05±0.00</b>
	<i>WikiTalk</i>	N/A	N/A	N/A	<u>0.18±0.03</u>	0.15±0.02	0.0±0.0	<b>0.24±0.05</b>
	<i>Lead-Lag</i>	N/A	<u>0.07±0.01</u>	<u>0.07±0.01</u>	<u>0.07±0.01</u>	<u>0.07±0.02</u>	<u>0.07±0.02</u>	<b>0.15±0.03</b>
$\mathcal{O}_{\text{vol\_sum}}^{\text{std}}$	<i>Telegram</i>	0.01±0.00	0.26±0.00	0.26±0.00	0.26±0.01	0.25±0.02	<b>0.35±0.00</b>	<u>0.28±0.01</u>
	<i>Blog</i>	0.00±0.00	0.07±0.00	0.00±0.00	0.05±0.00	<u>0.37±0.00</u>	0.00±0.00	<b>0.44±0.00</b>
	<i>Migration</i>	N/A	0.01±0.00	0.01±0.00	0.01±0.00	<u>0.02±0.00</u>	<u>0.02±0.00</u>	<b>0.04±0.01</b>
	<i>WikiTalk</i>	N/A	N/A	N/A	<b>0.17±0.04</b>	0.06±0.01	0.01±0.00	<u>0.14±0.02</u>
	<i>Lead-Lag</i>	N/A	<u>0.04±0.01</u>	<u>0.04±0.01</u>	<u>0.04±0.01</u>	<u>0.04±0.01</u>	<u>0.04±0.01</u>	<b>0.12±0.03</b>
$\mathcal{O}_{\text{vol\_sum}}^{\text{naive}}$	<i>Telegram</i>	0.01±0.00	<u>0.26±0.0</u>	<u>0.26±0.0</u>	<u>0.26±0.01</u>	0.25±0.02	0.23±0.0	<b>0.27±0.01</b>
	<i>Blog</i>	0.00±0.00	<u>0.07±0.0</u>	0.0±0.0	0.05±0.0	<u>0.37±0.0</u>	0.0±0.0	<b>0.44±0.0</b>
	<i>Migration</i>	N/A	0.01±0.00	0.01±0.00	0.01±0.00	<u>0.02±0.00</u>	0.01±0.00	<b>0.04±0.01</b>
	<i>WikiTalk</i>	N/A	N/A	N/A	<u>0.1±0.02</u>	0.04±0.0	0.0±0.0	<b>0.12±0.01</b>
	<i>Lead-Lag</i>	N/A	<u>0.30±0.06</u>	0.28±0.06	0.27±0.06	0.29±0.05	0.29±0.05	<b>0.32±0.11</b>

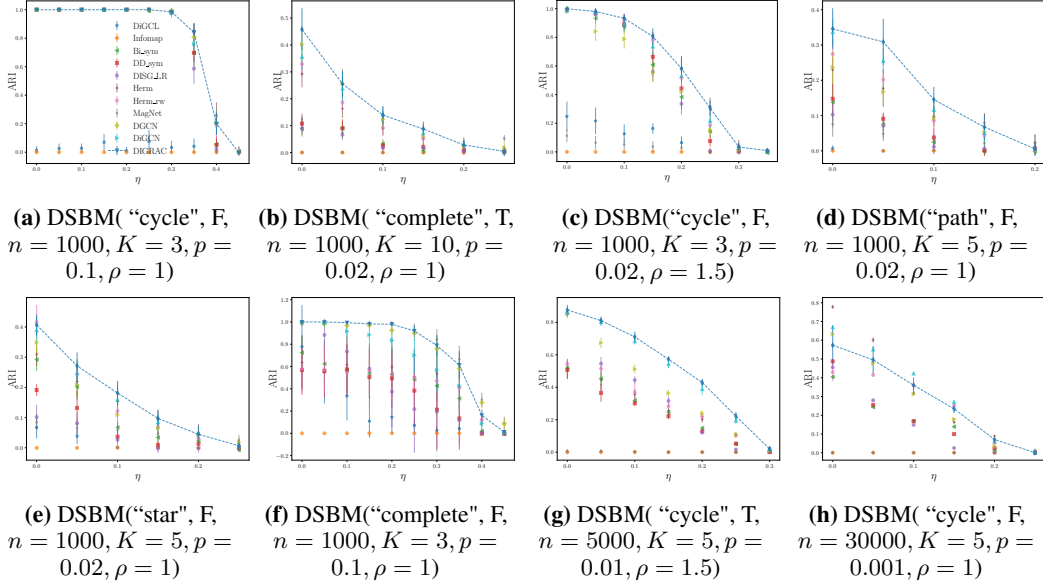
## 291 4.1 Data sets

292 **Synthetic data: Directed Stochastic Block Models** A standard directed stochastic blockmodel  
 293 (DSBM) is often used to represent a network cluster structure, see for example [1]. Its parameters are  
 294 the number  $K$  of clusters and the edge probabilities; given the cluster assignment of the nodes, the  
 295 edge indicators are independent. The DSBMs used in our experiments also depend on a meta-graph  
 296 adjacency matrix  $\mathbf{F} = (\mathbf{F}_{k,l})_{k,l=0,\dots,K-1}$  and a *filled* version of it,  $\tilde{\mathbf{F}} = (\tilde{\mathbf{F}}_{k,l})_{k,l=0,\dots,K-1}$ , and  
 297 on a noise level parameters  $\eta \leq 0.5$ . The meta-graph adjacency matrix  $\mathbf{F}$  is generated from the  
 298 given meta-graph structure, called  $\mathcal{M}$ . To include an ambient background, the filled meta-graph  
 299 adjacency matrix  $\tilde{\mathbf{F}}$  replaces every zero in  $\mathbf{F}$  that is not part of the imbalance structure by 0.5. The  
 300 filled meta-graph thus creates a number of *ambient nodes* which correspond to entries which are not  
 301 part of  $\mathcal{M}$  and thus are not part of a meaningful cluster; this set of *ambient nodes* is also called the  
 302 *ambient cluster*. First, we provide examples of structures of  $\mathbf{F}$  without any ambient nodes, where  $\mathbb{1}$   
 303 denotes the indicator function.

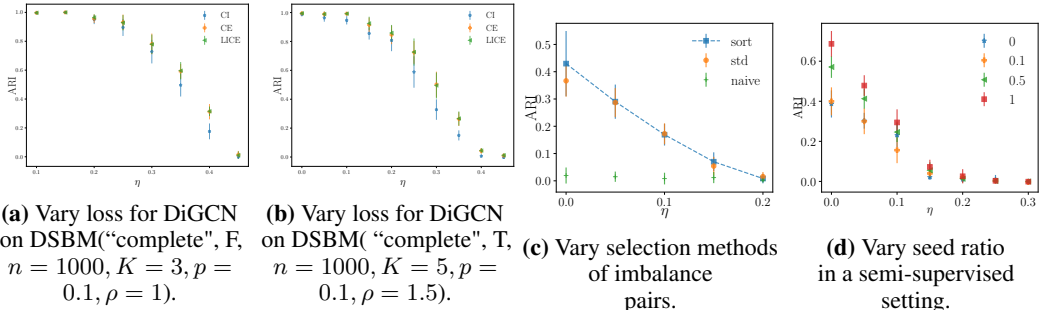
- 304 •(1) “cycle”:  $\mathbf{F}_{k,l} = (1 - \eta)\mathbb{1}(l = ((k + 1) \bmod K)) + \eta\mathbb{1}(l = ((k - 1) \bmod K)) + \frac{1}{2}\mathbb{1}(l = k)$ .
- 305 •(2) “path”:  $\mathbf{F}_{k,l} = (1 - \eta)\mathbb{1}(l = k + 1) + \eta\mathbb{1}(l = k - 1) + \frac{1}{2}\mathbb{1}(l = k)$ .
- 306 •(3) “complete”: assign diagonal entries  $\frac{1}{2}$ . For each pair  $(k, l)$  with  $k < l$ , let  $\mathbf{F}_{k,l}$  be  $\eta$  and  $1 - \eta$   
 307 with equal probability, then assign  $\mathbf{F}_{l,k} = 1 - \mathbf{F}_{k,l}$ .
- 308 •(4) “star”, following [45]: select the center node as  $\omega = \lfloor \frac{K-1}{2} \rfloor$  and set  $\mathbf{F}_{k,l} = (1 - \eta)\mathbb{1}(k =$   
 309  $\omega, l \text{ odd}) + \eta\mathbb{1}(k = \omega, l \text{ even}) + (1 - \eta)\mathbb{1}(l = \omega, k \text{ odd}) + \eta\mathbb{1}(l = \omega, k \text{ even})$ .

310 When ambient nodes are present, the construction involves two steps, with the first step the same as the  
 311 above, but with the following changes: For “cycle” meta-graph structure,  $\mathbf{F}_{k,l} = (1 - \eta)\mathbb{1}(l = ((k + 1)$   
 312  $\bmod (K - 1))) + \eta\mathbb{1}(l = ((k - 1) \bmod (K - 1))) + 0.5\mathbb{1}(l = k)$ . The second step is to assign  
 313 0 (0.5, resp.) to the last row and the last column of  $\mathbf{F}$  ( $\tilde{\mathbf{F}}$ , resp.). Figures 1(c-d) display a “cycle”  
 314 meta-graph structure with ambient nodes (in cluster 4). The majority of edges flow in the form  
 315  $0 \rightarrow 1 \rightarrow 2 \rightarrow 3 \rightarrow 0$ , while few flow from the opposite direction. Fig. 1(d) illustrates the  
 316 meta-graph adjacency matrix corresponding to this  $\tilde{\mathbf{F}}$ .

317 In our experiments, we choose the number of clusters, the (approximate) ratio,  $\rho$ , between the largest  
 318 and the smallest cluster size, and the number,  $n$ , of nodes. To tackle the hardest clustering task and also  
 319 focus on directionality, all pairs of nodes within a cluster and all pairs of nodes between clusters have  
 320 the same edge probability,  $p$ . Note that for  $\mathcal{M} = \text{“cycle”}$ , even the expected in-degree and out-degree  
 321 of all nodes are identical. Our DSBM, which we denote by  $\text{DSBM}(\mathcal{M}, \mathbb{1}(\text{ambient}), n, K, p, \rho, \eta)$ ,  
 322 is built similarly to [12] but with possibly unequal cluster sizes, with more details in App. C.3. For  
 323 each node  $v_i \in \mathcal{C}_k$ , and each node  $v_j \in \mathcal{C}_l$ , independently sample an edge from node  $v_i$  to node  $v_j$   
 324 with probability  $p \cdot \tilde{\mathbf{F}}_{k,l}$ . The parameter settings in our experiments are  $p \in \{0.001, 0.01, 0.02, 0.1\}$ ,  
 325  $\rho \in \{1, 1.5\}$ ,  $K \in \{3, 5, 10\}$ ,  $\mathbb{1}(\text{ambient}) \in \{\text{T}, \text{F}\}$  (True and False),  $n \in \{1000, 5000, 30000\}$ , and  
 326 we also vary the direction flip probability  $\eta$  from 0 to 0.45, with a 0.05 step size.



**Figure 3:** Test ARI comparison on synthetic data. Dashed lines highlight DIGRAC’s performance. Error bars are given by one standard error.



**Figure 4:** Ablation study. (c-d) are on DSBM("cycle", F,  $n = 1000, K = 5, p = 0.02, \rho = 1$ ).

327 **Real-world data** We perform experiments on five real-world digraph data sets with size ranging from  
 328 245 to over 2 million nodes: *Telegram* [3], *Blog* [46], *Migration* [4], *WikiTalk* [47], and *Lead-Lag*  
 329 [20], with details in App. C.3. We set the number of clusters  $K$  to be 4, 2, 10, 10, 10, respectively,  
 330 and values of  $\beta$  to be 5, 1, 9, 10, 3, respectively. Note that *Lead-Lag* comprises of 19 separate  
 331 networks constructed from yearly financial time series, rendering a total of 23 real-world networks.

## 332 4.2 Experimental results

333 **Training set-up** As training setup, we use 10% of all nodes from each cluster as test nodes, 10% as  
 334 validation nodes to select the model, and the remaining 80% as training nodes. In each setting, unless  
 335 otherwise stated, we carry out 10 experiments with different data splits. Error bars are given by one  
 336 standard error. When no node attributes are given, the matrix  $\mathbf{X}$  for DIGRAC is taken as the stacked  
 337 eigenvectors corresponding to the largest  $K$  eigenvalues of the random-walk symmetrized Hermitian  
 338 matrix used in the comparison method *Herm\_rw*. The imbalance loss function acts on the subgraph  
 339 induced by the training nodes. To further clarify the training setup, DIGRAC uses 0% of the labels in  
 340 training. As DIGRAC is a self-supervised method, in principle, we could use all nodes for training.  
 341 However for a fair comparison with other GNN methods we use only 80% of the nodes for training.  
 342 For supervised methods our split of 80% - 10% - 10% is a standard split. For the non-GNN methods,  
 343 all nodes are used for training. The default loss function for DIGRAC is  $\mathcal{L}_{\text{vol\_sum}}^{\text{sort}}$ .

344 **Results on synthetic data** Fig. 3 compares the numerical performance of DIGRAC with other  
 345 methods on synthetic data. For this Fig. we generate 5 DSBM networks under each parameter setting  
 346 and use 10 different data splits for each network, then average over the 50 runs. Error bars are given  
 347 by one standard error. App. C provides additional implementation details.



348 We conclude that DIGRAC compares favorably against state-of-the-art methods, on a wide range  
 349 of network densities and noise levels, on different network sizes, and with different underlying  
 350 meta-graph structures, with and without ambient nodes. Being a self-supervised method, DIGRAC  
 351 even attains comparable or better performance than fully-supervised GNN competitors.

352 **Results on real-world data** For our real-world data sets, the node in- and out-degrees may not  
 353 be identical across clusters. Moreover, as these data sets do not contain node attributes, DIGRAC  
 354 considers the eigenvectors corresponding to the largest  $K$  eigenvalues of the Hermitian matrix from  
 355 [12] to construct an input feature matrix. Table 1 reveals that DIGRAC provides competitive global  
 356 imbalance scores in three objectives discussed and across all real-world data sets, and outperforms all  
 357 other methods in 13 out of 15 instances, while attains the second-best performance for the remaining  
 358 two instances. The N/A entries for *WikiTalk* are caused by memory error, and the N/A entries  
 359 for InfoMap on *Migration* and *Lead-Lag* are due to its prediction of only one single cluster. For  
 360 *Migration*, as detailed in Fig. 1(b) and App. D.4, DIGRAC is able to uncover nontrivial migration  
 361 patterns, such as migration from California to Arizona, as discovered by [4]. *Lead-Lag* results in each  
 362 year are averaged over ten runs, while the mean and standard deviation values are calculated with  
 363 respect to the 19 years. The experiments indicate that edge directionality contains an important signal  
 364 that DIGRAC is able to capture. As App. D.2 illustrates, DIGRAC is able to provide comparable or  
 365 higher pairwise imbalance scores for the leading pairs. The fitted meta-graph plots in App. D.3 reveal  
 366 that DIGRAC is able to recover a directed flow imbalance between clusters in all of the selected data  
 367 sets. A comprehensive numerical comparison in App. D reveals similar conclusions.

### 368 4.3 Ablation study

369 Figures 4(a-b) compare the performance of DiGCN replacing the loss function by  $\mathcal{L}_{\text{vol\_sum}}^{\text{sort}}$  from  
 370 Eq. (3), indicated by “CI” (self-supervised loss only), or “LICE” (sum of supervised and self-  
 371 supervised loss), on two synthetic models. We find that replacing the supervised loss function with  
 372  $\mathcal{L}_{\text{vol\_sum}}^{\text{sort}}$  leads to comparable results, and that adding  $\mathcal{L}_{\text{vol\_sum}}^{\text{sort}}$  to the loss could be beneficial, indicating  
 373 that the imbalance objectives are more general than only applicable to DIMPA. Fig. 4(c) compares the  
 374 test ARI performance using three variants of loss functions on the same digraph. The current choice  
 375 “*sort*” performs best among these variants, indicating a benefit in only considering top pairs of individ-  
 376 ual imbalance scores. The “*std*” variant is comparable with the “*sort*” variant, but the “*sort*” variant  
 377 performs the best with prior knowledge on the network structure. More details on loss functions, compar-  
 378 ison with other variants, and evaluation on additional metrics are discussed in App. B, with similar  
 379 conclusions. As illustrated in Fig. 4(d), again on the same digraph, we also experiment on adding  
 380 seeds, with the seed ratio defined as the ratio of the number of seed nodes to the number of training  
 381 nodes. A supervised loss, following [42], is then applied to these seeds; App. C.5 contains additional  
 382 details. In conclusion, seed nodes with a supervised loss function enhance performance, and we infer  
 383 that our model can further boost its performance when additional label information is available.

## 384 5 Conclusion, limitations and outlook

385 DIGRAC provides an end-to-end pipeline to create node embeddings and perform directed clustering,  
 386 with or without available additional node features or cluster labels. We illustrate DIGRAC on  
 387 publicly available data without any personally identifiable information. DIGRAC could potentially  
 388 have societal impact, for example, in detecting clusters of fake accounts in social networks. While  
 389 we do not envision our work to have any negative societal impact, vigilance is of course required.

390 Current limitations that could be addressed by future work include detecting the number of clusters  
 391 [10, 48], instead of specifying it a-priori, as this is typically not available in real-world applications.  
 392 The relatively small sizes of the networks used in the paper (the largest has 2 million nodes) also  
 393 opens future direction in adapting our pipeline to extremely large networks, possibly combined  
 394 with sampling methods or mini-batch [49], rendering DIGRAC applicable to large scale industrial  
 395 applications. We also intent to further explore the effect of normalization terms in our objectives,  
 396 and to design more powerful objectives that could explicitly account for varying edge density.

397 Another future direction pertains to additional experiments in the semi-supervised setting, when there  
 398 exist seed nodes with known cluster labels, or when additional information is available in the form of  
 399 *must-link* and *cannot-link* constraints, popular in the *constrained clustering* literature [50, 51]. Further  
 400 research directions will also address the performance in the sparse regime, where spectral methods  
 401 are known to underperform, and various regularizations have been proven to be effective theoretically  
 402 and empirically; e.g., see regularization in the sparse regime for the undirected settings [52–54].

## References

- 403
- 404 [1] Fragkiskos D Malliaros and Michalis Vazirgiannis. Clustering and community detection in  
405 directed networks: A survey. *Physics reports*, 533(4):95–142, 2013. 1, 7
- 406 [2] William R. Palmer and Tian Zheng. Spectral clustering for directed networks. In Rosa M.  
407 Benito, Chantal Cherifi, Hocine Cherifi, Esteban Moro, Luis Mateus Rocha, and Marta Sales-  
408 Pardo, editors, *Complex Networks & Their Applications IX*, pages 87–99, Cham, 2021. Springer  
409 International Publishing. ISBN 978-3-030-65347-7. 1, 3
- 410 [3] Alexandre Bovet and Peter Grindrod. The Activity of the Far Right on Tele-  
411 gram. [https://www.researchgate.net/publication/346968575\\_The\\_Activity\\_](https://www.researchgate.net/publication/346968575_The_Activity_of_the_Far_Right_on_Telegram_v21)  
412 [of\\_the\\_Far\\_Right\\_on\\_Telegram\\_v21](https://www.researchgate.net/publication/346968575_The_Activity_of_the_Far_Right_on_Telegram_v21), 2020. 1, 8, 23, 34
- 413 [4] Marc J Perry. *State-to-state Migration Flows, 1995 to 2000*. US Department of Commerce,  
414 Economics and Statistics Administration, US . . . , 2003. 1, 2, 8, 9, 23, 35, 36
- 415 [5] Andrew Elliott, Angus Chiu, Marya Bazzi, Gesine Reinert, and Mihai Cucuringu. Core-  
416 periphery structure in directed networks. *Proceedings of the Royal Society A*, 476(2241):  
417 20190783, 2020. 1, 23, 34
- 418 [6] Michelle Girvan and Mark EJ Newman. Community structure in social and biological networks.  
419 *Proceedings of the National Academy of Sciences*, 99(12):7821–7826, 2002. 1
- 420 [7] Mark EJ Newman. Modularity and community structure in networks. *Proceedings of the*  
421 *National Academy of Sciences*, 103(23):8577–8582, 2006.
- 422 [8] Jure Leskovec, Kevin J Lang, Anirban Dasgupta, and Michael W Mahoney. Statistical properties  
423 of community structure in large social and information networks. In *Proceedings of the 17th*  
424 *International Conference on World Wide Web*, pages 695–704, 2008.
- 425 [9] Elizabeth A Leicht and Mark EJ Newman. Community structure in directed networks. *Physical*  
426 *Review Letters*, 100(11):118703, 2008.
- 427 [10] Yilin Chen and Jack W Baker. Community detection in spatial correlation graphs: Application  
428 to non-stationary ground motion modeling. *Computers and Geosciences*, 2021. 9
- 429 [11] Caiyan Jia, Yafang Li, Matthew B Carson, Xiaoyang Wang, and Jian Yu. Node attribute-  
430 enhanced community detection in complex networks. *Scientific Reports*, 7(1):1–15, 2017.  
431 1
- 432 [12] Mihai Cucuringu, Huan Li, He Sun, and Luca Zanetti. Hermitian matrices for clustering directed  
433 graphs: insights and applications. In *International Conference on Artificial Intelligence and*  
434 *Statistics*, pages 983–992. PMLR, 2020. 1, 2, 3, 4, 6, 7, 9, 22, 23, 24
- 435 [13] Steinar Laenen and He Sun. Higher-order spectral clustering of directed graphs. *Advances in*  
436 *Neural Information Processing Systems*, 2020. 1, 3
- 437 [14] Martin Rosvall and Carl T Bergstrom. Maps of random walks on complex networks reveal  
438 community structure. *Proceedings of the national academy of sciences*, 105(4):1118–1123,  
439 2008. 1, 3, 6, 40
- 440 [15] Kun Deng, Prashant G Mehta, and Sean P Meyn. Optimal kullback-leibler aggregation via  
441 spectral theory of markov chains. *IEEE Transactions on Automatic Control*, 56(12):2793–2808,  
442 2011. 2, 40
- 443 [16] Bernhard C Geiger, Tatjana Petrov, Gernot Kubin, and Heinz Koepl. Optimal kullback-leibler  
444 aggregation via information bottleneck. *IEEE Transactions on Automatic Control*, 60(4):  
445 1010–1022, 2014. 2, 40
- 446 [17] Ali Shojaie and Emily B Fox. Granger causality: A review and recent advances. *arXiv preprint*  
447 *arXiv:2105.02675*, 2021. 2
- 448 [18] Mukeshwar Dhamala, Govindan Rangarajan, and Mingzhou Ding. Analyzing information  
449 flow in brain networks with nonparametric Granger causality. *NeuroImage*, 41(2):354–362,  
450 2008. ISSN 1053-8119. doi: <https://doi.org/10.1016/j.neuroimage.2008.02.020>. URL <https://www.sciencedirect.com/science/article/pii/S1053811908001328>. 2
- 451
- 452 [19] Jakob Runge, Peer Nowack, Marlene Kretschmer, Seth Flaxman, and Dino Sejdinovic. Detecting  
453 and quantifying causal associations in large nonlinear time series datasets. *Science advances*, 5  
454 (11), 2019. 2, 41

- 455 [20] Stefanos Bennett, Mihai Cucuringu, and Gesine Reinert. Detection and clustering of lead-lag  
456 networks for multivariate time series with an application to financial markets. *7th SIGKDD*  
457 *Workshop on Mining and Learning from Time Series (MiLeTS)*, 2021. 2, 3, 8, 23
- 458 [21] Stefanos Bennett, Mihai Cucuringu, and Gesine Reinert. Lead-lag detection and network  
459 clustering for multivariate time series with an application to the us equity market, 2022. 2
- 460 [22] A Harzallah and R Sadourny. Observed lead-lag relationships between indian summer monsoon  
461 and some meteorological variables. *Climate dynamics*, 13(9):635–648, 1997. 2
- 462 [23] Venu Satuluri and Srinivasan Parthasarathy. Symmetrizations for clustering directed graphs. In  
463 *Proceedings of the 14th International Conference on Extending Database Technology*, pages  
464 343–354, 2011. 3, 6, 16
- 465 [24] Karl Rohe, Tai Qin, and Bin Yu. Co-clustering directed graphs to discover asymmetries  
466 and directional communities. *Proceedings of the National Academy of Sciences*, 113(45):  
467 12679–12684, 2016. 3, 6
- 468 [25] Jingnan Zhang, Xin He, and Junhui Wang. Directed community detection with network  
469 embedding. *Journal of the American Statistical Association*, pages 1–11, 2021. 3
- 470 [26] Zekun Tong, Yuxuan Liang, Changsheng Sun, David S Rosenblum, and Andrew Lim. Directed  
471 graph convolutional network. *arXiv preprint arXiv:2004.13970*, 2020. 3, 6
- 472 [27] Zekun Tong, Yuxuan Liang, Changsheng Sun, Xinke Li, David Rosenblum, and Andrew Lim.  
473 Digraph inception convolutional networks. *Advances in Neural Information Processing Systems*,  
474 33, 2020. 3, 6, 26
- 475 [28] Bojan Mohar. A new kind of Hermitian matrices for digraphs. *Linear Algebra and its Applica-*  
476 *tions*, 584:343–352, 2020. 3
- 477 [29] Xitong Zhang, Yixuan He, Nathan Brugnone, Michael Perlmutter, and Matthew Hirn. Magnet:  
478 A neural network for directed graphs. *arXiv preprint arXiv:2102.11391*, 2021. 3, 6, 26
- 479 [30] Zekun Tong, Yuxuan Liang, Henghui Ding, Yongxing Dai, Xinke Li, and Changhu Wang.  
480 Directed graph contrastive learning. *Advances in Neural Information Processing Systems*, 34,  
481 2021. 3, 6
- 482 [31] Yi Ma, Jianye Hao, Yaodong Yang, Han Li, Junqi Jin, and Guangyong Chen. Spectral-based  
483 graph convolutional network for directed graphs. *arXiv preprint arXiv:1907.08990*, 2019. 3
- 484 [32] Federico Monti, Karl Otness, and Michael M. Bronstein. Motifnet: A motif-based graph  
485 convolutional network for directed graphs. In *2018 IEEE Data Science Workshop (DSW)*, pages  
486 225–228, 2018. doi: 10.1109/DSW.2018.8439897. 3
- 487 [33] Yuan Yao, Yicheng Pan, Shaojiang Wang, Hongan Wang, and Waqas Nazeer. Flow-based  
488 clustering on directed graphs: A structural entropy minimization approach. *IEEE Access*, 8:  
489 152579–152591, 2019. 3
- 490 [34] Andrea Lancichinetti, Filippo Radicchi, José J Ramasco, and Santo Fortunato. Finding statisti-  
491 cally significant communities in networks. *PloS one*, 6(4):e18961, 2011. 3, 37
- 492 [35] Nicolas Dugué and Anthony Perez. *Directed Louvain: maximizing modularity in directed*  
493 *networks*. PhD thesis, Université d’Orléans, 2015. 3, 37
- 494 [36] Vincent A Traag, Ludo Waltman, and Nees Jan Van Eck. From louvain to leiden: guaranteeing  
495 well-connected communities. *Scientific reports*, 9(1):1–12, 2019. 3, 37
- 496 [37] Filippo Maria Bianchi, Daniele Grattarola, and Cesare Alippi. Spectral clustering with graph  
497 neural networks for graph pooling. In *International Conference on Machine Learning*, pages  
498 874–883. PMLR, 2020. 3
- 499 [38] Hongteng Xu, Dixin Luo, and Lawrence Carin. Scalable gromov-wasserstein learning for graph  
500 partitioning and matching. *Advances in neural information processing systems*, 32:3052–3062,  
501 2019. 4
- 502 [39] Hongteng Xu, Dixin Luo, Hongyuan Zha, and Lawrence Carin Duke. Gromov-wasserstein  
503 learning for graph matching and node embedding. In *International conference on machine*  
504 *learning*, pages 6932–6941. PMLR, 2019.
- 505 [40] Samir Chowdhury and Tom Needham. Generalized spectral clustering via gromov-wasserstein  
506 learning. In *International Conference on Artificial Intelligence and Statistics*, pages 712–720.  
507 PMLR, 2021. 4

- 508 [41] Hongteng Xu. Gromov-wasserstein factorization models for graph clustering. In *Proceedings*  
509 *of the AAAI conference on artificial intelligence*, volume 34, pages 6478–6485, 2020. 4
- 510 [42] Yixuan He, Gesine Reinert, Songchao Wang, and Mihai Cucuringu. SSSNET: Semi-Supervised  
511 Signed Network Clustering. In *Proceedings of the 2022 SIAM International Conference on*  
512 *Data Mining (SDM)*, pages 244–252. SIAM, 2022. 6, 9, 15, 25, 26
- 513 [43] Lawrence Hubert and Phipps Arabie. Comparing partitions. *Journal of Classification*, 2(1):  
514 193–218, 1985. 6
- 515 [44] Xin Liu, Hui-Min Cheng, and Zhong-Yuan Zhang. Evaluation of community detection methods.  
516 *IEEE Transactions on Knowledge and Data Engineering*, 32(9):1736–1746, 2019. 6, 26
- 517 [45] Andrew Elliott, Paul Reidy Milton Martinez Luaces, Mihai Cucuringu, and Gesine Reinert.  
518 Anomaly detection in networks using spectral methods and network comparison approaches.  
519 *arXiv preprint arXiv:1901.00402*, 2019. 7
- 520 [46] Lada A Adamic and Natalie Glance. The political blogosphere and the 2004 US election:  
521 divided they blog. In *Proceedings of the 3rd International Workshop on Link Discovery*, pages  
522 36–43, 2005. 8, 23
- 523 [47] Jure Leskovec, Daniel Huttenlocher, and Jon Kleinberg. Signed networks in social media.  
524 In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, pages  
525 1361–1370, 2010. 8, 23
- 526 [48] Maria A Riolo, George T Cantwell, Gesine Reinert, and Mark EJ Newman. Efficient method for  
527 estimating the number of communities in a network. *Physical Review E*, 96(3):032310, 2017. 9
- 528 [49] William L Hamilton, Rex Ying, and Jure Leskovec. Inductive representation learning on large  
529 graphs. In *Proceedings of the 31st International Conference on Neural Information Processing*  
530 *Systems*, pages 1025–1035, 2017. 9, 16
- 531 [50] Sugato Basu, Ian Davidson, and Kiri Wagstaff. *Constrained Clustering: Advances in Algorithms,*  
532 *Theory, and Applications*. CRC Press, 2008. 9
- 533 [51] Mihai Cucuringu, Ioannis Koutis, Sanjay Chawla, Gary Miller, and Richard Peng. Scalable  
534 Constrained Clustering: A Generalized Spectral Method. *Artificial Intelligence and Statistics*  
535 *Conference (AISTATS) 2016*, 2016. 9
- 536 [52] Kamalika Chaudhuri, Fan Chung, and Alexander Tsiatas. Spectral clustering of graphs with  
537 general degrees in the extended planted partition model. In *25th Annual Conference on Learning*  
538 *Theory*, volume 23 of *Proceedings of Machine Learning Research*, pages 35.1–35.23, Edinburgh,  
539 Scotland, 2012. JMLR Workshop and Conference Proceedings. 9
- 540 [53] Arash A. Amini, Aiyu Chen, Peter J. Bickel, and Elizaveta Levina. Pseudo-likelihood methods  
541 for community detection in large sparse networks. *The Annals of Statistics*, 41(4):2097–2122,  
542 2013.
- 543 [54] Mihai Cucuringu, Apoorv Vikram Singh, Déborah Sulem, and Hemant Tyagi. Regularized  
544 spectral methods for clustering signed networks. *arXiv:2011.01737*, 2020. 9
- 545 [55] Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional  
546 networks. *arXiv preprint arXiv:1609.02907*, 2016. 15
- 547 [56] Adam P Harrison and Dileepan Joseph. High performance rearrangement and multiplication  
548 routines for sparse tensor arithmetic. *SIAM Journal on Scientific Computing*, 40(2):C258–C281,  
549 2018. 15
- 550 [57] Gero Greiner and Riko Jacob. The I/O Complexity of Sparse Matrix Dense Matrix Multiplica-  
551 tion", booktitle="LATIN 2010: Theoretical Informatics. pages 143–156, Berlin, Heidelberg,  
552 2010. Springer Berlin Heidelberg. ISBN 978-3-642-12200-2. 15
- 553 [58] Elan Sopher Markowitz, Keshav Balasubramanian, Mehrnoosh Mirtaheri, Sami Abu-El-Haija,  
554 Bryan Perozzi, Greg Ver Steeg, and Aram Galstyan. Graph traversal with tensor functionals: A  
555 meta-algorithm for scalable learning. In *International Conference on Learning Representations*,  
556 2021. URL <https://openreview.net/forum?id=6DOZ8XNNfGN>. 16
- 557 [59] Matthias Fey, Jan E Lenssen, Frank Weichert, and Jure Leskovec. Gnnautoscale: Scalable and  
558 expressive graph neural networks via historical embeddings. *arXiv preprint arXiv:2106.05609*,  
559 2021. 16

- 560 [60] Louis HY Chen, Larry Goldstein, and Qi-Man Shao. *Normal approximation by Stein's method*.  
561 Springer Science & Business Media, 2010. 17
- 562 [61] Ryan L Phillips and Rita Ormsby. Industry classification schemes: An analysis and review.  
563 *Journal of Business & Finance Librarianship*, 21(1):1–25, 2016. 23
- 564 [62] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*  
565 *arXiv:1412.6980*, 2014. 26
- 566 [63] Pan Zhang. Evaluating accuracy of community detection using the relative normalized mutual  
567 information. *Journal of Statistical Mechanics: Theory and Experiment*, 2015(11):P11006, 2015.  
568 26
- 569 [64] Mihai Cucuringu, Vincent Blondel, and Paul Van Dooren. Extracting spatial information from  
570 networks with low-order eigenvectors. *Phys. Rev. E*, 87:032803, Mar 2013. doi: 10.1103/  
571 PhysRevE.87.032803. URL <http://link.aps.org/doi/10.1103/PhysRevE.87.032803>.  
572 35, 36
- 573 [65] Stijn Van Dongen. Graph clustering via a discrete uncoupling process. *SIAM Journal on Matrix*  
574 *Analysis and Applications*, 30(1):121–141, 2008. 40

575	<b>Contents</b>	
576	<b>1 Introduction</b>	<b>1</b>
577	<b>2 Related work</b>	<b>3</b>
578	<b>3 The DIGRAC framework</b>	<b>4</b>
579	3.1 Self-supervised loss for clustering . . . . .	4
580	3.2 Instantiation of DIGRAC . . . . .	5
581	<b>4 Experiments</b>	<b>6</b>
582	4.1 Data sets . . . . .	7
583	4.2 Experimental results . . . . .	8
584	4.3 Ablation study . . . . .	9
585	<b>5 Conclusion, limitations and outlook</b>	<b>9</b>
586	<b>Appendix</b>	<b>15</b>
587	<b>A Directed Mixed Path Aggregation (DIMPA)</b>	<b>15</b>
588	<b>B Loss and objectives</b>	<b>16</b>
589	B.1 Proof of Proposition 1 . . . . .	16
590	B.2 Additional details on probabilistic cut and volume . . . . .	17
591	B.3 Variants of normalization . . . . .	18
592	B.4 Selection of the loss function . . . . .	19
593	<b>C Implementation details</b>	<b>19</b>
594	C.1 Code . . . . .	19
595	C.2 Hardware . . . . .	19
596	C.3 Data . . . . .	22
597	C.3.1 Data splits and preprocessing . . . . .	22
598	C.3.2 Synthetic data . . . . .	22
599	C.3.3 Real-world data . . . . .	23
600	C.4 Hyperparameter selection for DIMPA . . . . .	24
601	C.5 Use of seed nodes in a semi-supervised manner . . . . .	25
602	C.5.1 Supervised loss . . . . .	25
603	C.5.2 Overall objective function . . . . .	26
604	C.6 Training . . . . .	26
605	C.7 Implementation details for the comparison methods . . . . .	26
606	C.8 Enlarged synthetic result figures . . . . .	26
607	C.9 NMI results example and reasons against using NMI . . . . .	26

608	<b>D Additional results on real-world data</b>	<b>27</b>
609	D.1 Extended result tables . . . . .	27
610	D.2 Ranked pairwise imbalance scores . . . . .	32
611	D.3 Predicted meta-graph flow matrix plots . . . . .	34
612	D.4 Migration plots . . . . .	35
613	D.5 Coping with outliers . . . . .	35
614	<b>E Discussion of related methods that are not compared against in the main text</b>	<b>36</b>

## 615 A Directed Mixed Path Aggregation (DIMPA)

616 To instantiate DIGRAC, we can employ any digraph aggregator that could generate the probability  
617 matrix  $\mathbf{P}$ . In this paper, we devise a simple yet effective directed mixed path aggregation scheme,  
618 to obtain the probability assignment matrix  $\mathbf{P}$  and feed it to the loss function, as a special case of  
619 the successful *SSSNET* method introduced by [42]. Thus, in order to build node embeddings, we  
620 capture local network information by taking a weighted average of information from neighbors  
621 within  $h$  hops. To this end, we row-normalize the adjacency matrix,  $\mathbf{A}$ , to obtain  $\overline{\mathbf{A}}^s$ . Similar  
622 to the regularization discussed in [55], we add a weighted self-loop to each node and normalize  
623 by setting  $\overline{\mathbf{A}}^s = (\tilde{\mathbf{D}}^s)^{-1} \tilde{\mathbf{A}}^s$ , where  $\tilde{\mathbf{A}}^s = \mathbf{A} + \tau \mathbf{I}$ , with  $\tilde{\mathbf{D}}^s$  the diagonal matrix with entries  
624  $\tilde{\mathbf{D}}^s(i, i) = \sum_j \tilde{\mathbf{A}}^s(i, j)$ , and  $\tau$  is a small value; we take  $\tau = 0.5$ ; see Section C.4 for details.

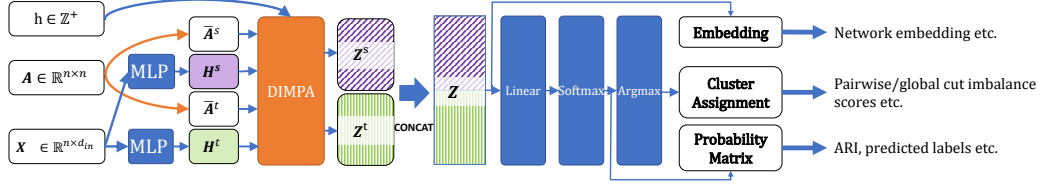
625 The  $h$ -hop **source** matrix is given by  $(\overline{\mathbf{A}}^s)^h$ . We denote the set of *up-to- $h$ -hop* source neighborhood  
626 matrices as  $\mathcal{A}^{s,h} = \{\mathbf{I}, \overline{\mathbf{A}}^s, \dots, (\overline{\mathbf{A}}^s)^h\}$ . Similarly, for aggregating information when each node is  
627 viewed as a **target** node of a link, we carry out the same procedure for  $\mathbf{A}^T$  which is the transpose of  
628  $\mathbf{A}$ . We denote the set of *up-to- $h$ -hop* target neighborhood matrices as  $\mathcal{A}^{t,h} = \{\mathbf{I}, \overline{\mathbf{A}}^t, \dots, (\overline{\mathbf{A}}^t)^h\}$ ,  
629 where  $\overline{\mathbf{A}}^t$  is the row-normalized target adjacency matrix calculated from  $\mathbf{A}^T$ . As convention, the  
630 superscript  $s$  stands for *source* and the superscript  $t$  stands for *target*.

631 Next, we define two feature mapping functions for source and target embeddings, respectively.  
632 Assume that for each node in  $\mathcal{V}$ , a vector of features is available, and summarize these features in the  
633 input feature matrix  $\mathbf{X}$ . The source embedding is given by

$$634 \mathbf{Z}^s = \left( \sum_{\mathbf{M} \in \mathcal{A}^{s,h}} \omega_{\mathbf{M}}^s \cdot \mathbf{M} \right) \cdot \mathbf{H}^s \in \mathbb{R}^{n \times d}, \quad (4)$$

635 where for each  $\mathbf{M}$ ,  $\omega_{\mathbf{M}}^s$  is a learnable scalar,  $d$  is the dimension of this embedding, and  $\mathbf{H}^s =$   
636  $\mathbf{MLP}^{(s,l)}(\mathbf{X})$ . Here, the hyperparameter  $l$  controls the number of layers in the multilayer perceptron  
637 (MLP) with ReLU activation; we fix  $l = 2$  throughout. Each layer of the MLP has the same  
638 number  $d$  of hidden units. The target embedding  $\mathbf{Z}^t$  is defined similarly, with  $s$  replaced by  $t$  in  
639 Eq. (4). Different parameters for the MLPs for different embeddings are possible. After these two  
640 decoupled aggregations, we concatenate the embeddings to obtain the final node embedding as a  
641  $n \times (2d)$  matrix  $\mathbf{Z} = \text{CONCAT}(\mathbf{Z}^s, \mathbf{Z}^t)$ . The embedding vector  $\mathbf{z}_i$  for a node  $v_i$  is the  $i^{\text{th}}$  row of  $\mathbf{Z}$ ,  
642  $\mathbf{z}_i := (\mathbf{Z})_{(i,:)} \in \mathbb{R}^{2d}$ .

643 After obtaining the embedding matrix  $\mathbf{Z}$ , we apply a linear layer (an affine transformation) to  $\mathbf{Z}$ , so  
644 that the resulting matrix has  $K$  columns. Next, we apply the unit *softmax* function to the rows and  
645 obtain the assignment probability matrix  $\mathbf{P}$ . Fig. 5 gives an overview of this implementation.  
646 To avoid computationally expensive and space unfriendly matrix operations, as described in Eq. 4,  
647 DIGRAC uses an efficient sparsity-aware implementation, described in Algorithm 1, without explicitly  
648 calculating the sets of powers  $\mathcal{A}^{s,h}$  and  $\mathcal{A}^{t,h}$ . We omit the subscript  $\mathcal{V}$  for ease of notation. The  
649 algorithm is efficient in the sense that it takes sparse matrices as input, and never explicitly computes a  
650 multiplication of two  $n \times n$  matrices. Therefore, for input feature dimension  $d_{\text{in}}$  and hidden dimension  
651  $d$ , if  $d' = \max(d_{\text{in}}, d) \ll n$ , time and space complexity of DIMPA, and implicitly DIGRAC, is  
652  $\mathcal{O}(|\mathcal{E}|d'h^2 + 2nd'K)$  and  $\mathcal{O}(2|\mathcal{E}| + 4nd' + nK)$ , respectively [56, 57].



**Figure 5:** DIGRAC with DIMPA as aggregator overview: from feature matrix  $\mathbf{X}$  and adjacency matrix  $\mathbf{A}$ , we first compute the row-normalized adjacency matrices  $\bar{\mathbf{A}}^s$  and  $\bar{\mathbf{A}}^t$ . Then, we apply two separate MLPs on  $\mathbf{X}$ , to obtain hidden representations  $\mathbf{H}^s$  and  $\mathbf{H}^t$ . Next, we compute their decoupled embeddings using Eq. (4), and its equivalent for target embeddings. The concatenated decoupled embeddings are the final embeddings. For node clustering tasks, we add a linear layer followed by a unit *softmax* to obtain the probability matrix  $\mathbf{P}$ . Applying *argmax* on each row of  $\mathbf{P}$  yields node cluster assignments.

653 While it is a current shortcoming of DIGRAC that it does not scale well to very large networks, this  
 654 limitation is shared by all the GNN competitors compared against in the paper, and some of the  
 655 spectral methods. DIGRAC scales well in the sense that when the underlying network is sparse,  
 656 the sparsity is preserved throughout the pipeline. In contrast, Bi\_sym and DD\_sym [23] construct  
 657 derived dense matrices for manipulation, rendering the methods no longer scalable. These methods  
 658 resulted in N/A values in Table 1 in the main text. For large-scale networks, DIMPA is amenable  
 659 to a minibatch version using neighborhood sampling, similar to the minibatch forward propagation  
 660 algorithm in [49, 58]. We are also aware of a framework [59] for scaling up graph neural networks  
 661 automatically, where theoretical guarantees are provided, and ideas there will be exploited in future.  
 662 We expect that the theoretical guarantees could be adapted to our situation.

---

**Algorithm 1:** Weighted Multi-Hop Neighbor Aggregation (DIMPA).

---

**Input** : (Sparse) row-normalized adjacency matrices  $\bar{\mathbf{A}}^s, \bar{\mathbf{A}}^t$ ; initial hidden representations  
 $\mathbf{H}^s, \mathbf{H}^t$ ; hop  $h (h \geq 2)$ ; lists of scalar weights  
 $\Omega^s = (\omega_M^s, \mathbf{M} \in \mathcal{A}^{s,h}), \Omega^t = (\omega_M^t, \mathbf{M} \in \mathcal{A}^{t,h})$ .  
**Output** : Vector representations  $\mathbf{z}_i$  for all  $v_i \in \mathcal{V}$  given by  $\mathbf{Z}$ .  
 $\tilde{\mathbf{X}}^s \leftarrow \bar{\mathbf{A}}^s \mathbf{H}^s; \quad \tilde{\mathbf{X}}^t \leftarrow \bar{\mathbf{A}}^t \mathbf{H}^t;$   
 $\mathbf{Z}^s \leftarrow \Omega^s[0] \cdot \mathbf{H}^s + \Omega^s[1] \cdot \tilde{\mathbf{X}}^s; \quad \mathbf{Z}^t \leftarrow \Omega^t[0] \cdot \mathbf{H}^t + \Omega^t[1] \cdot \tilde{\mathbf{X}}^t;$   
**for**  $i \leftarrow 2$  **to**  $h$  **do**  
 $\quad \tilde{\mathbf{X}}^s \leftarrow \bar{\mathbf{A}}^s \tilde{\mathbf{X}}^s; \quad \tilde{\mathbf{X}}^t \leftarrow \bar{\mathbf{A}}^t \tilde{\mathbf{X}}^t;$   
 $\quad \mathbf{Z}^s \leftarrow \mathbf{Z}^s + \Omega^s[i] \cdot \tilde{\mathbf{X}}^s; \quad \mathbf{Z}^t \leftarrow \mathbf{Z}^t + \Omega^t[i] \cdot \tilde{\mathbf{X}}^t;$   
**end**  
 $\mathbf{Z} = \text{CONCAT}(\mathbf{Z}^s, \mathbf{Z}^t);$

---

## 663 B Loss and objectives

### 664 B.1 Proof of Proposition 1

Moreover we clarify that we make an assumption on the limiting behavior of the weights, namely that

$$\frac{\max_e |w_e|}{\sqrt{\sum_e w_e^2}} = o(m(k, l))$$

665 where  $m(k, l)$  is the number of edges. This is a natural assumption: In the case that all weights are  
 666 equal in absolute value, this assumption is satisfied as then  $\frac{\max_e |w_e|}{\sqrt{\sum_e w_e^2}} = \frac{1}{\sqrt{m(k, l)}}$ . The assumption is  
 667 generally satisfied when there is not too much variability in the weights. If for example all but one  
 668 weight pair was equal to 0, then the assumption would be violated, and also a normal approximation  
 669 would not hold as there would only be two non-zero observations.

670 **Proposition 2.** Suppose that  $C_k$  and  $C_l$  are two clusters of  $n_k$  and  $n_l$  nodes, respectively, with  $m(k, l)$   
 671 edges between them, with symmetric edge weights  $w_{ij} = w_{ji} \in [0, 1]$  and with edge direction drawn



672 independently at random with equal probability  $\frac{1}{2}$  for each direction. We assume that the edge weights  
 673 satisfy  $\frac{\max_e |w_e|}{\sqrt{\sum_e w_e^2}} = o(m(k, l))$ . Then  $W(\mathcal{C}_k, \mathcal{C}_l) - W(\mathcal{C}_l, \mathcal{C}_k)$  is approximately normally distributed  
 674 with mean 0 and variance  $\|w\|^2$  as  $m(k, l) \rightarrow \infty$ .

**Proof.** For each edge between the two clusters  $\mathcal{C}_k$  and  $\mathcal{C}_l$ , the edge direction is random, i.e. the edge is from  $\mathcal{C}_k$  to  $\mathcal{C}_l$  with probability 0.5, and  $\mathcal{C}_l$  to  $\mathcal{C}_k$  with probability 0.5 also. Let  $\mathcal{E}^{k,l}$  denote the set of  $m(k, l) > 0$  edges between  $\mathcal{C}_k$  and  $\mathcal{C}_l$ . For every edge  $e \in \mathcal{E}^{k,l}$ , the edge direction is encoded by a Rademacher random variable  $X_e$  with  $X_e = 1$  if the edge is from  $\mathcal{C}_k$  to  $\mathcal{C}_l$ , and  $X_e = -1$  otherwise. Then  $(X_e + 1)/2 \sim \text{Ber}(0.5)$  is a Bernoulli(0.5) random variable with mean  $2 \times 0.5 - 1 = 0$  and variance  $2^2 \times 0.5 \times (1 - 0.5) = 1$ . We have the representation

$$W(\mathcal{C}_k, \mathcal{C}_l) - W(\mathcal{C}_l, \mathcal{C}_k) = \sum_{e \in \mathcal{E}^{k,l}} X_e w_e$$

675 as the sum of  $m(k, l)$  independent bounded random variables with finite third moments. Moreover,  
 676  $W(\mathcal{C}_k, \mathcal{C}_l) - W(\mathcal{C}_l, \mathcal{C}_k)$  has mean 0 and variance  $\|w\|^2$ . The assertion now follows from a version of  
 677 the Central Limit Theorem, Theorem 3.4 in [60]; we repeat the relevant part here:

**Theorem 1** (Extract from Theorem 3.4 in [60]). *Let  $\xi_1, \dots, \xi_n$  be independent random variables with zero means satisfying  $\sum_{i=1}^n \text{Var}(\xi_i) = 1$  and assume that there is a  $\delta > 0$  such that  $|\xi_i| \leq \delta$  for  $1 \leq i \leq n$ . Let  $\Phi$  denote the cumulative distribution function of the standard normal distribution. Then*

$$\sup_{z \in \mathbb{R}} \left| \mathbb{P} \left( \sum_{i=1}^n \xi_i \leq z \right) - \Phi(z) \right| \leq 3.3\delta.$$

678 We apply this theorem with  $n$  replaced by  $m(k, l)$ , the number of edges, and take  $\xi_e = \frac{X_e w_e}{\sqrt{\sum_e w_e^2}}$ .  
 679 Then  $\xi_e$  has mean zero and, using an enumeration of the edges,  $\sum_{e=1}^{m(k,l)} \text{Var}(\xi_e) = 1$ . Moreover,  
 680  $|\xi_e| \leq \frac{\max_e |w_e|}{\sqrt{\sum_e w_e^2}} =: \delta$  holds for all  $e \in \{1, \dots, m(k, l)\}$  and hence the theorem applies for the limit  
 681  $m(k, l) \rightarrow \infty$ . The stated result follows from using that if  $Z/\sigma$  has the standard normal distribution  
 682 then  $Z$  has the mean zero normal distribution with variance  $\sigma^2$ .

683

□

## 684 B.2 Additional details on probabilistic cut and volume

Recall that the **probabilistic cut** from cluster  $\mathcal{C}_k$  to  $\mathcal{C}_l$  is defined as

$$W(\mathcal{C}_k, \mathcal{C}_l) = \sum_{i,j \in \{1, \dots, n\}} \mathbf{A}_{i,j} \cdot \mathbf{P}_{i,k} \cdot \mathbf{P}_{j,l} = (\mathbf{P}_{(:,k)})^T \mathbf{A} \mathbf{P}_{(:,l)},$$

where  $\mathbf{P}_{(:,k)}$ ,  $\mathbf{P}_{(:,l)}$  denote the  $k^{\text{th}}$  and  $l^{\text{th}}$  columns of the assignment probability matrix  $\mathbf{P}$ , respectively. The **imbalance flow** between clusters  $\mathcal{C}_k$  and  $\mathcal{C}_l$  is defined as

$$|W(\mathcal{C}_k, \mathcal{C}_l) - W(\mathcal{C}_l, \mathcal{C}_k)|,$$

685 for  $k, l \in \{0, \dots, K - 1\}$ . The loss functions proposed in the main paper can be understood in terms  
 686 of a probabilistic notion of degrees, as follows. We define the probabilistic out-degree of node  $v_i$  with  
 687 respect to cluster  $k$  by  $\tilde{d}_{i,k}^{(\text{out})} = \sum_{j=1}^n \mathbf{A}_{i,j} \cdot \mathbf{P}_{j,k} = (\mathbf{A} \mathbf{P}_{(:,k)})_i$ , where subscript  $i$  refers to the  $i^{\text{th}}$   
 688 entry of the vector  $\mathbf{A} \mathbf{P}_{(:,k)}$ . Similarly, we define the probabilistic in-degree of node  $v_i$  with respect to  
 689 cluster  $k$  by  $\tilde{d}_{i,k}^{(\text{in})} = (\mathbf{A}^T \mathbf{P}_{(:,k)})_i$ , where  $\mathbf{A}^T$  is the transpose of  $\mathbf{A}$ . The **probabilistic degree** of node  
 690  $v_i$  with respect to cluster  $k$  is  $\tilde{d}_{i,k} = \tilde{d}_{i,k}^{(\text{in})} + \tilde{d}_{i,k}^{(\text{out})} = ((\mathbf{A}^T + \mathbf{A}) \mathbf{P}_{(:,k)})_i = \sum_{j=1}^n (\mathbf{A}_{i,j} + \mathbf{A}_{j,i}) \cdot \mathbf{P}_{j,k}$ .

691 For comparisons and ease of interpretation, it is advantageous to normalize the imbalance flow  
 692 between clusters; for this purpose, we introduce the probabilistic volume of a cluster, as follows.  
 693 The **probabilistic out-volume** for cluster  $\mathcal{C}_k$  is defined as  $VOL^{(\text{out})}(\mathcal{C}_k) = \sum_{i,j} \mathbf{A}_{j,i} \cdot \mathbf{P}_{j,k}$ , and  
 694 the **probabilistic in-volume** for cluster  $\mathcal{C}_k$  is defined as  $VOL^{(\text{in})}(\mathcal{C}_k) = (\mathbf{A}^T \mathbf{P}_{(:,k)})_i$ , where  $\mathbf{A}^T$  is the  
 695 transpose of  $\mathbf{A}$ . These volumes can be viewed as sum of probabilistic out-degrees and in-degrees,  
 696 respectively; for example,  $VOL^{(\text{in})}(\mathcal{C}_k) = \sum_{i=1}^n \tilde{d}_{i,k}^{(\text{in})}$ . Then, it holds true that

$$VOL^{(\text{out})}(\mathcal{C}_k) = \sum_{i,j} \mathbf{A}_{i,j} \cdot \mathbf{P}_{i,k} \geq \sum_{i,j} \mathbf{A}_{i,j} \cdot \mathbf{P}_{i,k} \cdot \mathbf{P}_{j,l} = W(\mathcal{C}_k, \mathcal{C}_l), \quad (5)$$

697 since entries in  $\mathbf{P}$  are probabilities, which are in  $[0, 1]$ , and all entries of  $\mathbf{A}$  are nonnegative. Similarly,  
 698  $VOL^{(in)}(\mathcal{C}_k) \geq W(\mathcal{C}_l, \mathcal{C}_k)$ .

The **probabilistic volume** for cluster  $\mathcal{C}_k$  is defined as

$$VOL(\mathcal{C}_k) = VOL^{(out)}(\mathcal{C}_k) + VOL^{(in)}(\mathcal{C}_k) = \sum_{i,j} (\mathbf{A}_{i,j} + \mathbf{A}_{j,i}) \cdot \mathbf{P}_{j,k}.$$

699 Then, it holds true that  $VOL(\mathcal{C}_k) \geq W(\mathcal{C}_k, \mathcal{C}_l)$  for all  $l \in \{0, \dots, K-1\}$  and

$$\min(VOL(\mathcal{C}_k), VOL(\mathcal{C}_l)) \geq \max(W(\mathcal{C}_k, \mathcal{C}_l), W(\mathcal{C}_l, \mathcal{C}_k)) \geq |W(\mathcal{C}_k, \mathcal{C}_l) - W(\mathcal{C}_l, \mathcal{C}_k)|. \quad (6)$$

700 When there exists a strong imbalance, then  $|W(\mathcal{C}_k, \mathcal{C}_l) - W(\mathcal{C}_l, \mathcal{C}_k)| \approx \max(W(\mathcal{C}_k, \mathcal{C}_l), W(\mathcal{C}_l, \mathcal{C}_k))$ .  
 701 As an extreme case, if  $\mathbf{P}_{j,l} = 1$  for all nonnegative terms in the summations in Eq. (5), and  
 702  $VOL^{(in)}(\mathcal{C}_k) = 0$ , then  $|W(\mathcal{C}_k, \mathcal{C}_l) - W(\mathcal{C}_l, \mathcal{C}_k)| = VOL(\mathcal{C}_k)$ .

### 703 B.3 Variants of normalization

704 Recall that the imbalance term involved in most of our experiments, named  $CI^{\text{vol\_sum}}$ , is defined as

$$CI^{\text{vol\_sum}}(k, l) = 2 \frac{|W(\mathcal{C}_k, \mathcal{C}_l) - W(\mathcal{C}_l, \mathcal{C}_k)|}{VOL(\mathcal{C}_k) + VOL(\mathcal{C}_l)} \in [0, 1]. \quad (7)$$

705 An alternative, which does not take volumes into account, is given by

$$CI^{\text{plain}}(k, l) = \frac{|W(\mathcal{C}_k, \mathcal{C}_l) - W(\mathcal{C}_l, \mathcal{C}_k)|}{W(\mathcal{C}_k, \mathcal{C}_l) + W(\mathcal{C}_l, \mathcal{C}_k)} = 2 \left| \frac{W(\mathcal{C}_k, \mathcal{C}_l)}{W(\mathcal{C}_k, \mathcal{C}_l) + W(\mathcal{C}_l, \mathcal{C}_k)} - \frac{1}{2} \right| \in [0, 1]. \quad (8)$$

706 We call this cut flow imbalance  $CI^{\text{plain}}$  as it does not penalize extremely unbalanced cluster sizes.

707 To achieve balanced cluster sizes and still constrain each imbalance term to be in  $[0, 1]$ , one solution is  
 708 to multiply the imbalance flow value by the minimum of  $VOL(\mathcal{C}_k)$  and  $VOL(\mathcal{C}_l)$ , and then divide by  
 709  $\max_{(k', l') \in \mathcal{T}} (\min(VOL(\mathcal{C}_{k'}), VOL(\mathcal{C}_{l'})))$ , where  $\mathcal{T} = \{(\mathcal{C}_k, \mathcal{C}_l) : 0 \leq k < l \leq K-1, k, l \in \mathbb{Z}\}$ .

710 The reason for using  $\mathcal{T}$  is that  $CI^{\text{plain}}(k, l)$  is symmetric with respect to  $k$  and  $l$ , and  $CI^{\text{plain}}(k, l) = 0$   
 711 whenever  $k = l$ . Note that the maximum of the minimum here equals the second largest volume  
 712 among clusters. We then obtain  $CI^{\text{vol\_min}}$  as

$$CI^{\text{vol\_min}}(k, l) = CI^{\text{plain}}(k, l) \times \frac{\min(VOL(\mathcal{C}_k), VOL(\mathcal{C}_l))}{\max_{(k', l') \in \mathcal{T}} (\min(VOL(\mathcal{C}_{k'}), VOL(\mathcal{C}_{l'})))}. \quad (9)$$

713 Another potential choice, denoted  $CI^{\text{vol\_max}}$ , whose normalization follows from the same reasoning  
 714 as  $CI^{\text{vol\_sum}}$ , is given by

$$CI^{\text{vol\_max}}(k, l) = \frac{|W(\mathcal{C}_k, \mathcal{C}_l) - W(\mathcal{C}_l, \mathcal{C}_k)|}{\max(VOL(\mathcal{C}_k), VOL(\mathcal{C}_l))} \in [0, 1]. \quad (10)$$

715 Note that the current  $CI^{\text{vol\_sum}}(k, l)$  term can be reformulated as

$$CI^{\text{vol\_sum}}(k, l) = 2 \frac{|W(\mathcal{C}_k, \mathcal{C}_l) - W(\mathcal{C}_l, \mathcal{C}_k)|}{VOL(\mathcal{C}_k) + VOL(\mathcal{C}_l)} = 2 \frac{W(\mathcal{C}_k, \mathcal{C}_l) + W(\mathcal{C}_l, \mathcal{C}_k)}{VOL(\mathcal{C}_k) + VOL(\mathcal{C}_l)} \times CI^{\text{plain}}(k, l), \quad (11)$$

716 with the first term in the decomposition corresponding to the relative ratio of inter- and intra-cluster  
 717 edge density. For our synthetic data, this term is constant as we have constant edge density across the  
 718 graph. However, for certain real-world data sets, one could also maximize this first term by increasing  
 719 the inter-cluster density while decreasing the intra-cluster density, which seems to be a side effect.  
 720 However, in our experiments, we also evaluate our results with different metrics, including objectives  
 721 without any normalization, and conclude that this side effect does not create any issues in our data  
 722 sets.

**Table 2:** Naming conventions for objectives and loss functions

Selection variant / $CI$	$CI^{\text{vol\_sum}}$	$CI^{\text{vol\_min}}$	$CI^{\text{vol\_max}}$	$CI^{\text{plain}}$
sort	$\mathcal{O}_{\text{vol\_sum}}^{\text{sort}}, \mathcal{L}_{\text{vol\_sum}}^{\text{sort}}$	$\mathcal{O}_{\text{vol\_min}}^{\text{sort}}, \mathcal{L}_{\text{vol\_min}}^{\text{sort}}$	$\mathcal{O}_{\text{vol\_max}}^{\text{sort}}, \mathcal{L}_{\text{vol\_max}}^{\text{sort}}$	$\mathcal{O}_{\text{plain}}^{\text{sort}}, \mathcal{L}_{\text{plain}}^{\text{sort}}$
std	$\mathcal{O}_{\text{vol\_sum}}^{\text{std}}, \mathcal{L}_{\text{vol\_sum}}^{\text{std}}$	$\mathcal{O}_{\text{vol\_min}}^{\text{std}}, \mathcal{L}_{\text{vol\_min}}^{\text{std}}$	$\mathcal{O}_{\text{vol\_max}}^{\text{std}}, \mathcal{L}_{\text{vol\_max}}^{\text{std}}$	$\mathcal{O}_{\text{plain}}^{\text{std}}, \mathcal{L}_{\text{plain}}^{\text{std}}$
naive	$\mathcal{O}_{\text{vol\_sum}}^{\text{naive}}, \mathcal{L}_{\text{vol\_sum}}^{\text{naive}}$	$\mathcal{O}_{\text{vol\_min}}^{\text{naive}}, \mathcal{L}_{\text{vol\_min}}^{\text{naive}}$	$\mathcal{O}_{\text{vol\_max}}^{\text{naive}}, \mathcal{L}_{\text{vol\_max}}^{\text{naive}}$	$\mathcal{O}_{\text{plain}}^{\text{naive}}, \mathcal{L}_{\text{plain}}^{\text{naive}}$

#### 723 B.4 Selection of the loss function

724 Table 2 provides naming conventions of all the twelve pairs of variants of objectives and loss functions  
 725 used in this paper. We select the loss functions for DIGRAC based on two representative models,  
 726 and compare the performance of different loss functions. We use DIMPA (introduced in A) as an  
 727 instantiation of DIGRAC’s aggregator, for which  $d = 32$ , hidden units,  $h = 2$  hops, and no seed  
 728 nodes. Figures 6(a) and 7 compare twelve choices of loss combinations on a DSBM with  $n = 1000$   
 729 nodes,  $K = 5$  blocks,  $\rho = 1, p = 0.02$  without ambient nodes, with a complete meta-graph structure.  
 730 The subscript indicates the choice of pairwise imbalance, and the superscript indicates the variant  
 731 for selecting pairs. Figures 6(b) and 8 are based on a DSBM with  $n = 1000$  nodes,  $K = 5$  blocks,  
 732  $\rho = 1, p = 0.02$  without ambient nodes, with a cycle meta-graph structure. For these figures, dash  
 733 lines highlight the “*sort*” variant as well as the “*std*” variant based on  $CI^{\text{vol\_sum}}$ , which have been  
 734 introduced in the main text.

735 We also plot the imbalance evolution curves for the above two synthetic models when  $\eta = 0.05$ , for  
 736 all the loss variants, in Figure 9.

737 These figures indicate that the “*sort*” variant generally provides the best test ARI performance and  
 738 the best overall global imbalance scores, among which using normalizations  $CI^{\text{vol\_sum}}$  and  $CI^{\text{vol\_max}}$   
 739 perform the best. The “*std*” variant is comparable with the “*sort*” variant in many instances, but is less  
 740 stable in performance. We observe, however, from Figure 9, that the “*std*” variants normally converge  
 741 much faster. Taking the above into account, if we have prior knowledge on the network structure,  
 742 or when we could conduct some prior analysis on the value  $\beta$  to take, the “*sort*” variant should be  
 743 the variant of choice. Further, from Figure 9, we observe that normalization in the loss function  
 744 helps avoid the degenerate situation that the loss does not decrease. Such degeneracy can occur in  
 745 the “*plain*” variants, raising issues about the practical usefulness of these variants. We observe that  
 746  $\mathcal{L}_{\text{vol\_min}}^{\text{sort}}$  appears to behave worse than  $\mathcal{L}_{\text{vol\_sum}}^{\text{sort}}$  and  $\mathcal{L}_{\text{vol\_max}}^{\text{sort}}$ , even when using the “*sort*” variant to  
 747 select pairwise imbalance scores. One possible explanation is that  $\mathcal{L}_{\text{vol\_min}}^{\text{sort}}$  does not penalize extreme  
 748 volume sizes, and that it takes minimum as well as maximum which, as functions of the data, are not  
 749 as smooth as taking a summation. Throughout our experiments in the main text, we hence use the  
 750 loss function  $\mathcal{L}_{\text{vol\_sum}}^{\text{sort}}$ .

## 751 C Implementation details

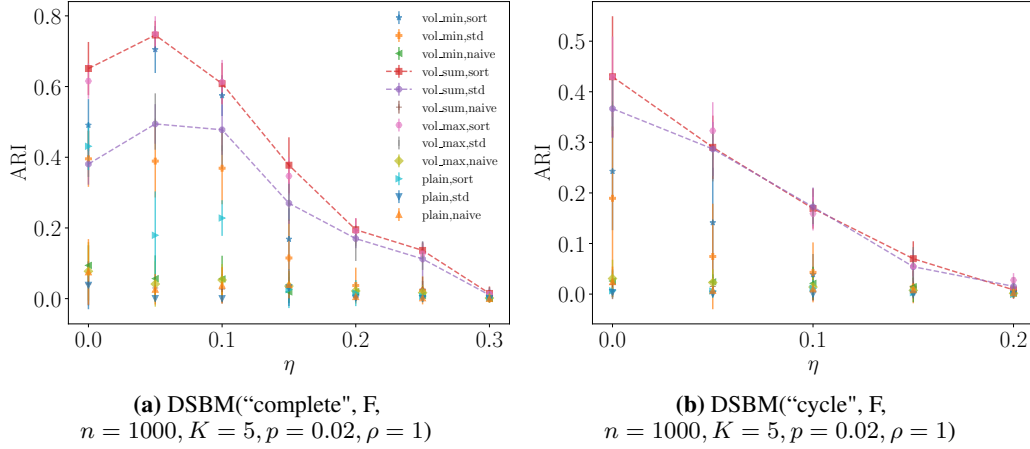
### 752 C.1 Code

753 To fully reproduce our results, anonymized code and preprocessed data are available at <https://anonymous.4open.science/r/DIGRAC>.  
 754

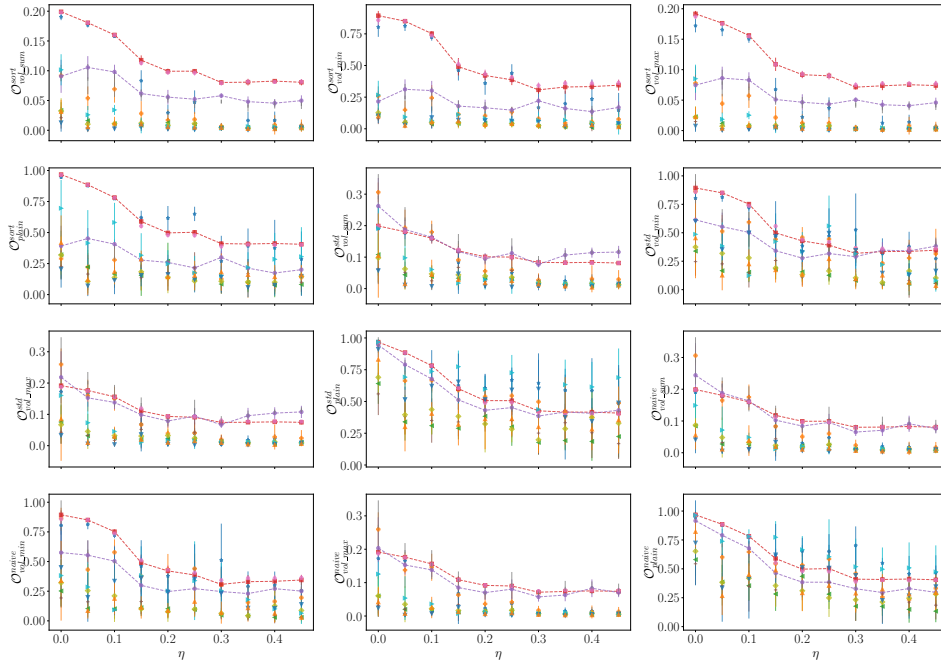
### 755 C.2 Hardware

756 Experiments were conducted on a compute node with 8 Nvidia RTX 8000, 48 Intel Xeon Silver  
 757 4116 CPUs and 1000GB RAM, a compute node with 4 NVIDIA GeForce RTX 2080, 32 Intel Xeon  
 758 E5-2690 v3 CPUs and 64GB RAM, a compute node with 2 NVIDIA Tesla K80, 16 Intel Xeon  
 759 E5-2690 CPUs and 252GB RAM, and an Intel 2.90GHz i7-10700 processor with 8 cores and 16  
 760 threads.

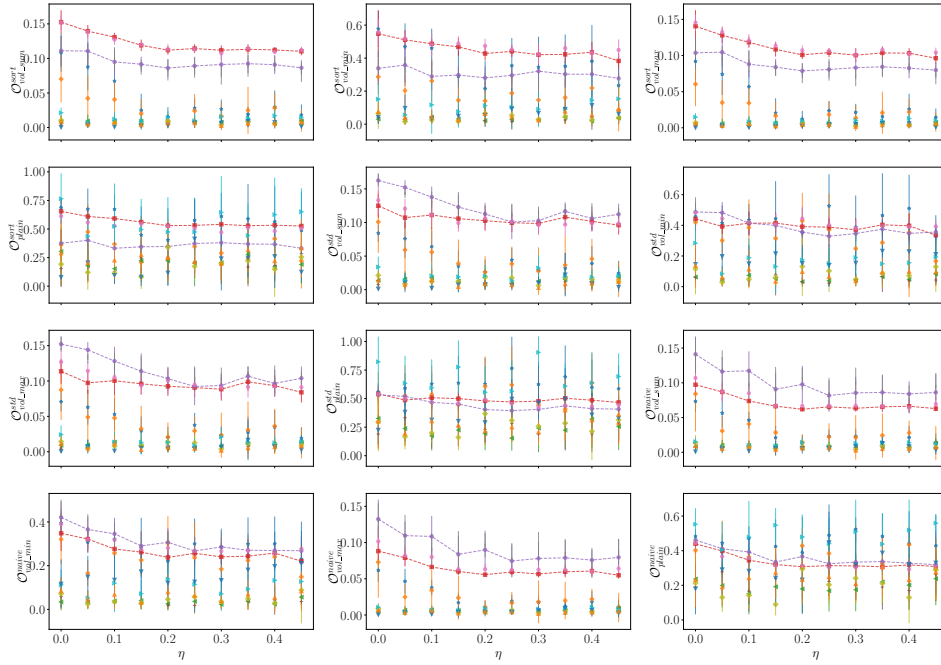
761 With this setup, all experiments for spectral methods, MagNet, DiGCL, and DIGRAC can be com-  
 762 pleted within two days, including repeated experiments, to obtain averages over multiple runs. DGCN,  
 763 DiGCN, and MagNet have much longer run time (especially DGCN, which is space-consuming,  
 764 and we cannot run many experiments in parallel), with a total of three days for them to finish. The  
 765 slow speed stems from the competitor methods; some of the other GNN methods take a long time to



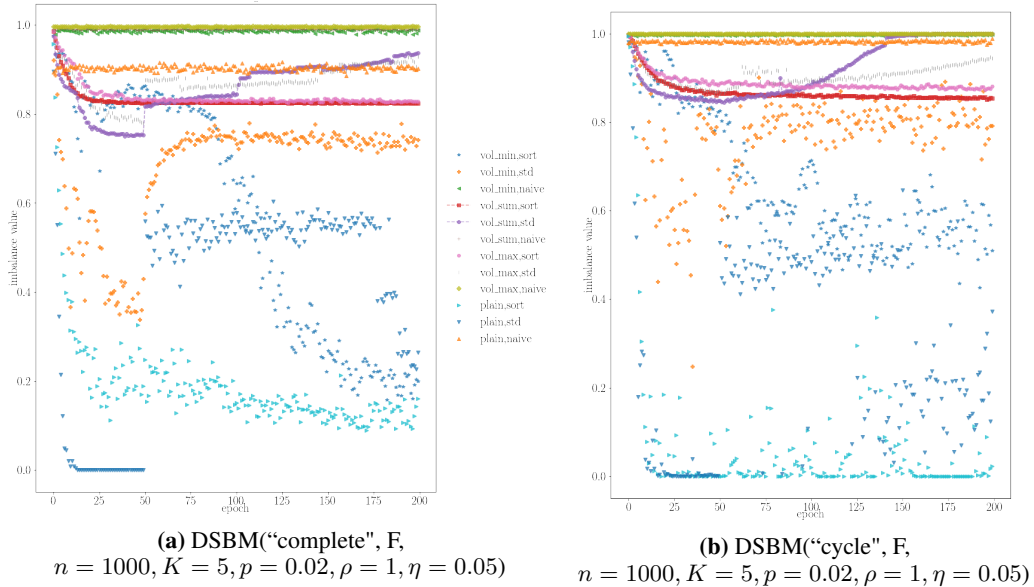
**Figure 6:** ARI comparison of loss functions on DSBM with 1000 nodes, 5 blocks,  $\rho = 1, p = 0.02$  without ambient nodes, of cycle (left) and complete (right) meta-graph structures, respectively. The first component of the legend is the choice of pairwise imbalance, and the second component is the variant of selecting pairs. The naming conventions for the abbreviations in the legend are provided in Table 2.



**Figure 7:** Imbalance scores comparison of loss functions on DSBM with 1000 nodes, 5 blocks,  $\rho = 1, p = 0.02$  without ambient nodes, of the **complete meta-graph** structure. The legend is the same as Fig. 6(a).



**Figure 8:** Imbalance scores comparison of loss functions on DSBM with 1000 nodes, 5 blocks,  $\rho = 1, p = 0.02$  without ambient nodes, of the **cyclic meta-graph** structure. The legend is the same as Fig. 6(a).



**Figure 9:** Imbalance loss evolution comparison of loss functions on DSBM with 1000 nodes, 5 blocks,  $\rho = 1, p = 0.02, \eta = 0.05$  without ambient nodes, of cycle (left) and complete (right) meta-graph structures, respectively. The first component of the legend is the choice of pairwise imbalance, and the second component is the variant of selecting pairs. The naming conventions for the abbreviations in the legend are provided in Table 2.

run. Table 1 in the main text shows N/A values for Bi\_sym and for DD\_sym exactly for this reason. Empirically, DIGRAC is among the fastest among all GNN methods to which it is compared. In detail, Table 3 reports the average runtime for all GNN methods on a variety of DSBM models, and illustrates that DIGRAC indeed takes the least or second least computational time per epoch. The results are averaged over 10 runs for the first 200 epochs. DiGCL is also efficient in running time, but with worse performance than DIGRAC even as a supervised method, see the enlarged synthetic results in Sec. C.8 (Figure 13). The total number of epochs required until the validation loss does not decrease for 200 epochs (or the maximum number of 1000 epochs is reached) varies for different data sets.

**Table 3:** GNN average runtime (seconds per epoch) comparison. The results are averaged over 10 runs for the first 200 epochs. The fastest is highlighted in **bold red** while the second fastest is marked with underline blue.

Runtime (second per epoch on average)/GNN method	DiGCL	DGCN	DiGCN	MagNet	DIGRAC
DSBM(“complete”, T, $n = 1000, K = 5, p = 0.1, \rho = 1.5, \eta = 0.1$ )	<b>0.107</b>	0.606	0.469	0.369	<u>0.308</u>
DSBM(“path”, F, $n = 1000, K = 5, p = 0.02, \rho = 1, \eta = 0.15$ )	<b>0.061</b>	0.227	0.212	0.238	<u>0.201</u>
DSBM(“star”, F, $n = 1000, K = 5, p = 0.02, \rho = 1, \eta = 0.3$ )	<b>0.095</b>	0.305	0.294	0.324	<u>0.292</u>
DSBM(“star”, F, $n = 5000, K = 5, p = 0.02, \rho = 1, \eta = 0.4$ )	0.222	0.966	0.276	<u>0.116</u>	<b>0.101</b>
DSBM(“cycle”, F, $n = 5000, K = 5, p = 0.01, \rho = 1.5, \eta = 0$ )	0.177	0.330	0.099	<u>0.095</u>	<b>0.089</b>
DSBM(“cycle”, F, $n = 30000, K = 5, p = 0.001, \rho = 1, \eta = 0$ )	<b>0.070</b>	0.868	0.208	0.183	<u>0.156</u>

### 775 C.3 Data

#### 776 C.3.1 Data splits and preprocessing

777 The results comparing DIGRAC with other methods on synthetic data are averaged over 50 runs, five  
 778 synthetic networks under the same setting, each with 10 different data splits. For synthetic data, 10%  
 779 of all nodes are selected as test nodes for each cluster (the actual number is the ceiling of the total  
 780 number of nodes times 0.1, to avoid falling below 10% of test nodes), 10% are selected as validation  
 781 nodes (for model selection and early-stopping; again, we consider the ceiling for the actual number),  
 782 while the remaining roughly 80% are selected as training nodes (the actual number can never be  
 783 higher than 80% due to using the ceiling for both the test and validation splits). To further clarify  
 784 the training setup, we use 0% of the labels in training. As DIGRAC is a self-supervised method,  
 785 in principle, we could use all nodes for training. However, for a fair comparison with other GNN  
 786 methods, we use only 80% of the nodes for training. For supervised methods our split of 80% - 10% -  
 787 10% is a standard split. For the non-GNN methods, all nodes are used for training.

788 For both synthetic and real-world data sets, we extract the largest weakly connected component for  
 789 experiments, as our framework could be applied to different weakly connected components, if the  
 790 digraph is disconnected. Isolated nodes do not include any imbalance information. As customary in  
 791 community detection, they are often omitted in real networks. When “ground-truth” is given, test  
 792 results are averaged over 10 different data splits on one network. When no labels are available, results  
 793 are averaged over 10 different data splits.

794 Averaged results are reported with error bars representing one standard deviation in the figures, and  
 795 plus/minus one standard deviation in the tables.

#### 796 C.3.2 Synthetic data

797 Our synthetic data, DSBM, which we denote by  $\text{DSBM}(\mathcal{M}, \mathbb{1}(\text{ambient}), n, K, p, \rho, \eta)$ , is built  
 798 similarly to [12] but with possibly unequal cluster sizes: **•(1)** Assign cluster sizes  $n_0 \leq n_1 \leq$   
 799  $\dots \leq n_{K-1}$  with size ratio  $\rho \geq 1$ , as follows. If  $\rho = 1$  then the first  $K - 1$  clusters have the  
 800 same size  $\lfloor n/K \rfloor$  and the last cluster has size  $n - (K - 1)\lfloor n/K \rfloor$ . If  $\rho > 1$ , we set  $\rho_0 = \rho^{\frac{1}{K-1}}$ .  
 801 Solving  $\sum_{i=0}^{K-1} \rho_0^i n_0 = n$  and taking integer value gives  $n_0 = \lfloor n(1 - \rho_0)/(1 - \rho_0^K) \rfloor$ . Further, set  
 802  $n_i = \lfloor \rho_0 n_{i-1} \rfloor$ , for  $i = 1, \dots, K - 2$  if  $K \geq 3$ , and  $n_{K-1} = n - \sum_{i=0}^{K-2} n_i$ . Then the ratio of  
 803 the size of the largest to the smallest cluster is approximately  $\rho_0^{K-1} = \rho$ . **•(2)** Assign each node  
 804 randomly to one of  $K$  clusters, so that each cluster has the allocated size. **•(3)** For node  $v_i, v_j \in \mathcal{C}_k$ ,  
 805 independently sample an edge from node  $v_i$  to node  $v_j$  with probability  $p \cdot \tilde{\mathbf{F}}_{k,k}$ . **•(4)** For each pair

806 of different clusters  $C_k, C_l$  with  $k \neq l$ , for each node  $v_i \in C_k$ , and each node  $v_j \in C_l$ , independently  
 807 sample an edge from node  $v_i$  to node  $v_j$  with probability  $p \cdot \tilde{\mathbf{F}}_{k,l}$ .

### 808 C.3.3 Real-world data

809 For real-world data sets, we choose the number  $K$  of clusters in the meta-graph and the number  $\beta$  of  
 810 edges between clusters in the meta-graph as follows. As they are needed as input for DIGRAC, we  
 811 resort to Herm\_rw [12] as an initial view of the network clustering. When a suitable meta-graph is  
 812 suggested in a previous publication, then we use that choice. Otherwise, the number  $K$  of clusters  
 813 is determined using the clustering from Herm\_rw. First, we pick a range of  $K$ , and for each  $K$ , we  
 814 calculate the global imbalance scores and plot the predicted meta-graph flow matrix  $\mathbf{F}'$  based on the  
 815 clustering from Herm\_rw. Its entries are defined as

$$\mathbf{F}'(k, l) = \mathbb{1}(W(C_k, C_l) + W(C_l, C_k) > 0) \times \frac{W(C_k, C_l)}{W(C_k, C_l) + W(C_l, C_k)}. \quad (12)$$

816 These entries can be viewed as predicted probabilities of edge directions. Then, we choose  $K$  from  
 817 this range so that the predicted meta-graph flow matrix has the highest imbalance scores and strong  
 818 imbalance in the predicted meta-graph flow matrix.

819 The choice of  $\beta$ , which we assume should be equal to the number of edges in the meta-graph, is as  
 820 follows. We plot the ranked pairs of  $\text{CI}^{\text{plain}}$  values from Herm\_rw and select the  $\beta$  which is at least as  
 821 large as  $K - 2$ , to allow the meta-graph to be connected, and which corresponds to a large drop in  
 822 the plot.

823 Here we provide a brief description for each of the data sets; Table 4 gives the number,  $n$ , of nodes,  
 824 the number,  $|\mathcal{E}|$ , of directed edges, the number  $|\mathcal{E}^r|$ , of reciprocal edges (self-loops are counted once  
 825 and for  $u \neq v$ , a reciprocal edge  $u \rightarrow v, v \rightarrow u$  is counted twice) as well as their percentage among  
 826 all edges, for the real-world networks, illustrating the variability in network size and density (defined  
 827 as  $|\mathcal{E}|/[n(n-1)]$ ).

828 •*Telegram* [3] is a pairwise influence network between  $n = 245$  Telegram channels with  $|\mathcal{E}| = 8,912$   
 829 directed edges. It is found in [3] that this network reveals a core-periphery structure in the sense of  
 830 [5]. A directed core-periphery structure arises when there is a densely connected group of edges – a  
 831 core – and sparsely connected groups of peripheral nodes with edges leading into the core, as well as  
 832 sparsely connected groups of peripheral nodes with edges coming out of the core. Following [3] we  
 833 assume  $K = 4$  clusters, and the core-periphery structures gives  $\beta = 5$ .

834 •*Blog* [46] records  $|\mathcal{E}| = 19,024$  directed edges between  $n = 1,212$  political blogs from the 2004  
 835 US presidential election. In [46] it is found that there is an underlying structure with  $K = 2$  clusters  
 836 corresponding to the Republican and Democratic parties. Hence we choose  $K = 2$  and  $\beta = 1$ .

837 •*Migration* [4] reports the number of people that migrated between pairs of counties in the US during  
 838 1995-2000. It involves  $n = 3,075$  counties and  $|\mathcal{E}| = 721,432$  directed edges after obtaining the  
 839 largest weakly connected component. We choose  $K = 10$  and  $\beta = 9$ , following [12]. Since the  
 840 original digraph has extremely large entries, to cope with these outliers, we preprocess the input  
 841 network by

$$\mathbf{A}_{i,j} = \frac{\mathbf{A}_{i,j}}{\mathbf{A}_{i,j} + \mathbf{A}_{j,i}} \mathbb{1}(\mathbf{A}_{i,j} > 0), \forall i, j \in \{1, \dots, n\}, \quad (13)$$

842 which follows the preprocessing of [12]. The results for not doing this preprocessing is provided in  
 843 Table 12.

844 •*WikiTalk* [47] contains all users and discussion from the inception of Wikipedia until Jan. 2008. The  
 845  $n = 2,388,953$  nodes in the network represent Wikipedia users and a directed edge from node  $v_i$   
 846 to node  $v_j$  denotes that user  $i$  edited at least once a talk page of user  $j$ . There are  $|\mathcal{E}| = 5,018,445$   
 847 edges. We choose  $K = 10$  clusters among candidates  $\{2, 3, 5, 6, 8, 10\}$ , and  $\beta = 10$ .

848 •*Lead-Lag* [20] contains yearly lead-lag matrices from 269 stocks from 2001 to 2019. We choose  
 849  $K = 10$  clusters based on the GICS industry sectors [61], and choose  $\beta = 3$  to emphasize the top  
 850 three pairs of imbalance values. The lead-lag matrices are built from time series of daily price log  
 851 returns, as detailed in [20]. The lead-lag metric for entry  $(i, j)$  in the network encodes a measure of  
 852 the extent to which stock  $i$  leads stock  $j$ , and is obtained by applying a functional that computes the  
 853 signed normalized area under the curve (auc) of the standard cross-correlation function (ccf). The  
 854 resulting matrix is skew-symmetric, and entry  $(i, j)$  quantifies the extent to which stock  $i$  leads or  
 855 lags stocks  $j$ , thus leading to a directed network interpretation. Starting from the skew-symmetric  
 856 matrix, we further convert negative entries to zero, so that the resulting digraph can be directly fed

857 into other methods; note that this step does not throw away any information, and is pursued only to  
 858 render the representation of the digraph consistent with the format expected by all methods compared,  
 859 including DIGRAC. Note that the statistics given in Table 4 are averaged over the 19 years.

**Table 4:** Summary statistics for the real-world networks.

data set	$n$	$ \mathcal{E} $	density	weighted	$ \mathcal{E}^r $	$\frac{ \mathcal{E}^r }{ \mathcal{E} }(\%)$
<i>Telegram</i>	245	8,912	$1.28 \cdot 10^{-2}$	True	1,572	17.64
<i>Blog</i>	1,222	19,024	$1.49 \cdot 10^{-1}$	True	4,617	24.27
<i>Migration</i>	3,075	721,432	$7.63 \cdot 10^{-2}$	True	351,100	48.67
<i>WikiTalk</i>	2,388,953	5,018,445	$8.79 \cdot 10^{-7}$	False	723,526	14.42
<i>Lead-Lag</i>	269	29,159	$4.04 \cdot 10^{-1}$	True	0.00	0.00

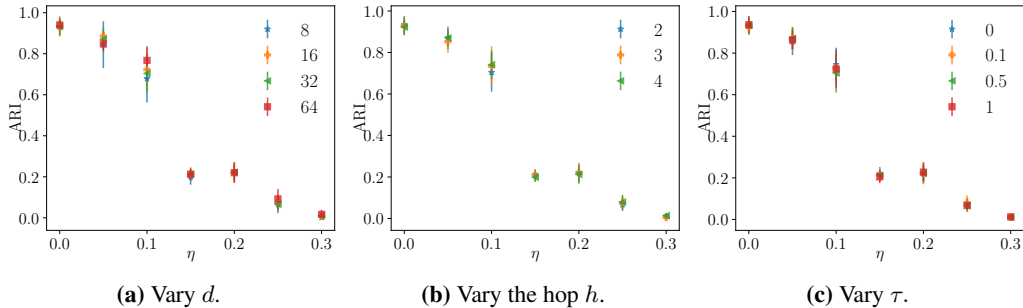
860 As input features, after obtaining eigenvectors from Hermitian matrices constructed as in [12], we  
 861 standardize each column vector so that it has mean zero and variance one. We use these features for  
 862 all GNN methods except MagNet, since MagNet has its own way of generating random features of  
 863 dimension one.

#### 864 C.4 Hyperparameter selection for DIMPA

865 We conduct hyperparameter selection via a greedy search, for DIGRAC implemented with DIMPA  
 866 as its aggregator. To explain the details, consider for example the following synthetic data setting:  
 867 DSBM with 1000 nodes, 5 clusters,  $\rho = 1$ , and  $p = 0.02$ , without ambient nodes under different  
 868 hyperparameter settings. By default, we use the loss function  $\mathcal{L}_{\text{vol\_sum}}^{\text{sort}}$ ,  $d = 32$  hidden units, hop  
 869  $h = 2$ , and no seed nodes. Instead of a grid search, we tune hyperparameters according to what  
 870 performs the best in the default setting of the respective GNN method. The procedure starts with a  
 871 random setting. For the next iteration, the hyperparameters are set to the current best setting (based  
 872 on the last iteration), independently. For example, if we start with  $a = 1, b = 2, c = 3$ , and we find  
 873 that under this default setting, the best  $a$  (when fixing  $b = 2, c = 3$ ) is 2 and the best  $b$  (when fixing  
 874  $a = 1, c = 3$ ) is 3, and the best  $c$  is 3 (when fixing  $a = 1, b = 2$ ), then for the next iteration, we set  
 875  $a = 2, b = 3, c = 3$ . If two settings give similar results, we choose the simpler setting, for example,  
 876 the smaller hop size. When we reach a local optimum, we stop searching. Indeed, just a few iterations  
 877 (less than five) were required for us to find the current setting, as DIGRAC tends to be robust to most  
 878 hyperparameters.

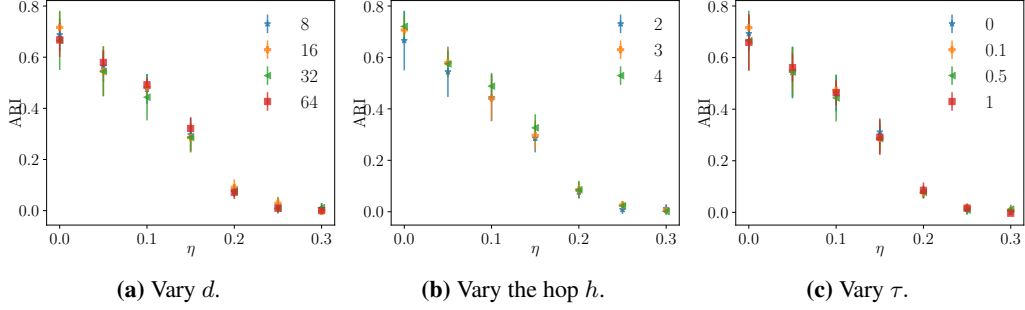
879 Fig. 10, 11 and 12 are plots corresponding to the same setting but for three different meta-graph  
 880 structures, namely the complete meta-graph structure, the cycle structure but with ambient nodes, and  
 881 the complete structure with ambient nodes, respectively.

882 In theory, more hidden units give better expressive power. To reduce complexity, we use 32 hidden  
 883 units throughout, which seems to have desirable performance. We observe that for low-noise  
 884 regimes, more hidden units actually hurt performance. We can draw a similar conclusion about the  
 885 hyperparameter selection. In terms of  $\tau$ , DIGRAC seems to be robust to different choices. Therefore,  
 886 we use  $\tau = 0.5$  throughout.

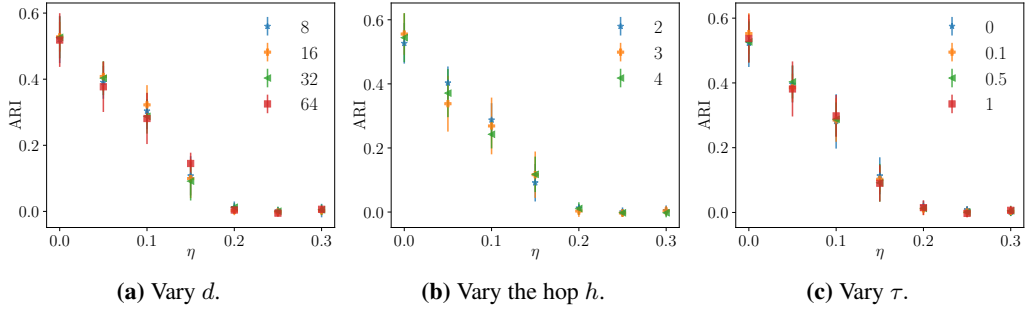


**Figure 10:** Hyperparameter analysis on different hyperparameter settings on the **complete** DSBM with 1000 nodes, 5 clusters,  $\rho = 1$ , and  $p = 0.02$  **without** ambient nodes.





**Figure 11:** Hyperparameter analysis on different hyperparameter settings on the **complete** DSBM with 1000 nodes, 5 clusters,  $\rho = 1$ , and  $p = 0.02$  **with** ambient nodes.



**Figure 12:** Hyperparameter analysis on different hyperparameter settings on the **cycle** DSBM with 1000 nodes, 5 clusters,  $\rho = 1$ , and  $p = 0.02$  **with** ambient nodes.

## 887 C.5 Use of seed nodes in a semi-supervised manner

### 888 C.5.1 Supervised loss

889 For seed nodes in  $\mathcal{V}^{\text{seed}}$ , similar to the loss function in [42], we use as a supervised loss function the  
890 sum of a cross-entropy loss and a triplet loss. The cross-entropy loss is given by

$$\mathcal{L}_{\text{CE}} = -\frac{1}{|\mathcal{V}^{\text{seed}}|} \sum_{v_i \in \mathcal{V}^{\text{seed}}} \sum_{k=1}^K \mathbb{1}(v_i \in \mathcal{C}_k) \log((\mathbf{p}_i)_k), \quad (14)$$

891 where  $\mathbb{1}$  is the indicator function,  $\mathcal{C}_k$  denotes the  $k^{\text{th}}$  cluster, and  $(\mathbf{p}_i)_k$  denotes the  $k^{\text{th}}$  entry of  
892 probability vector  $(\mathbf{p}_i)$ . With the function  $L : \mathbb{R}^2 \rightarrow \mathbb{R}$  given by  $L(x, y) = [x - y]_+$  (where the  
893 subscript  $+$  indicates taking the maximum of the expression value and 0), the triplet loss is defined as

$$\mathcal{L}_{\text{triplet}} = \frac{1}{|\mathcal{S}|} \sum_{(v_i, v_j, v_k) \in \mathcal{S}} L(\text{CS}(\mathbf{z}_i, \mathbf{z}_j), \text{CS}(\mathbf{z}_i, \mathbf{z}_k)), \quad (15)$$

894 where  $\mathcal{S} \subseteq \mathcal{V}^{\text{seed}} \times \mathcal{V}^{\text{seed}} \times \mathcal{V}^{\text{seed}}$  is a set of node triplets:  $v_i$  is an anchor seed node, and  $v_j$  is a seed  
895 node from the same cluster as the anchor, while  $v_k$  is from a different cluster; and  $\text{CS}(\mathbf{z}_i, \mathbf{z}_j)$  is the  
896 cosine similarity of the embeddings of nodes  $v_i$  and  $v_j$ . We choose cosine similarity so as to avoid  
897 sensitivity to the magnitude of the embeddings. The triplet loss is designed so that, given two seed  
898 nodes from the same cluster and one seed node from a different cluster, the respective embeddings of  
899 the pairs from different clusters should be farther away than the embedding of the pair within the  
900 same cluster.

901 We then consider the weighted sum  $\mathcal{L}_{\text{CE}} + \gamma_t \mathcal{L}_{\text{triplet}}$  as the supervised part of the loss function for  
902 DIGRAC, for some parameter  $\gamma_t > 0$ . The parameter  $\gamma_t$  arises as follows. The cosine similarity  
903 between two randomly picked vectors in  $d$  dimensions is bounded by  $\sqrt{\ln(d)/d}$  with high probability.  
904 In our experiments  $d = 32$ , and  $\sqrt{\ln(2d)/(2d)} \approx 0.25$ . In contrast, for fairly uniform clustering, the  
905 cross-entropy loss grows like  $\log n$ , which in our experiments ranges between 3 and 17. Thus some  
906 balancing of the contribution is required. Following [42], we choose  $\gamma_t = 0.1$  in our experiments.

### 907 C.5.2 Overall objective function

908 By combining Eq. (14), Eq. (15), and Eq. (3), our objective function for semi-supervised training  
909 with known seed nodes minimizes

$$\mathcal{L} = \mathcal{L}_{\text{vol\_sum}}^{\text{sort}} + \gamma_s(\mathcal{L}_{\text{CE}} + \gamma_t\mathcal{L}_{\text{triplet}}), \quad (16)$$

910 where  $\gamma_s, \gamma_t > 0$  are weights for the supervised part of the loss and triplet loss within the supervised  
911 part, respectively. We set  $\gamma_s = 50$  as we want our model to perform well on seed nodes. The weights  
912 could be tuned depending on how important each term is perceived to be.

### 913 C.6 Training

914 For all synthetic data, we train DIGRAC with a maximum of 1000 epochs, and stop training when no  
915 gain in validation performance is achieved for 200 epochs (early-stopping). For real-world data, no  
916 “ground-truth” labels are available; we use all nodes to train and stop training when the training loss  
917 does not decrease for 200 epochs, or when we reach the maximum number of epochs, 1000.

918 When using the “std” variant for training, for the initial 50 epochs, we apply the “sort” variant with  
919  $\beta = 3$  for a reasonable starting clustering probability matrix for training, as otherwise during the  
920 initial training epochs possibly no pairs could be picked out. During the epochs actually utilizing this  
921 “std” variant, if no pairs could be picked out, we temporarily switch to the “naive” variant to count all  
922 the pairs for that epoch.

923 For the two-layer MLP, we do not have a bias term for each layer, and we use Rectified Linear Unit  
924 (ReLU) followed by a dropout layer with 0.5 dropout probability between the two layers, following  
925 [42]. We use Adam [62] as the optimizer and  $\ell_2$  regularization with weight decay  $5 \cdot 10^{-4}$  to avoid  
926 overfitting. We use as learning rate 0.01 throughout.

### 927 C.7 Implementation details for the comparison methods

928 In our experiments, we compare DIGRAC against five spectral methods, InfoMap, and four GNN-  
929 based supervised methods on synthetic data, and spectral methods and InfoMap on real data. The  
930 reason we are not able to compare DIGRAC with the other GNNs (namely, DGCN, DiGCN, MagNet,  
931 and DiGCL) on these data sets is due to the fact that these data sets do not have labels, which are  
932 required by the other GNN methods. We use the same hyperparameter settings stated in these papers.  
933 Data splits for all models are the same; the comparison GNNs are trained with 80% nodes under label  
934 supervision.

935 For MagNet, we use  $q = 0.25$  for the phase matrix as in [29], because it is mentioned that  $q = 0.25$   
936 lays the most emphasis on directionality, which is our main focus in this paper. Code for MagNet is  
937 from <https://github.com/matthew-hirn/magnet>. For DiGCN, we use the code from [https://github.com/flyingtango/DiGCN/blob/main/code/digcn\\_ib.py](https://github.com/flyingtango/DiGCN/blob/main/code/digcn_ib.py) with option “adj\_type”  
938 equals “ib”. As a recommended option in [27], we use three layers for DiGCN. All other settings  
939 are the same as in the original paper [27]. Code for DiGCL is from <https://github.com/flyingtango/DiGCL>, where we adopt the settings for Cora\_ML for hyperparameters.  
941

### 942 C.8 Enlarged synthetic result figures

943 Figure 13 enlarges the results in the main text on synthetic data, with the same conclusions to be  
944 drawn.

### 945 C.9 NMI results example and reasons against using NMI

946 As NMI is an often used measure for assessing similarities between partitions, Fig. 14 provides NMI  
947 results on some synthetic models mentioned in the main text. The results are qualitatively similar to  
948 the ARI results in Fig. 3.

949 We do not use NMI in the main text to evaluate results as NMI is known to suffer from finite size  
950 effects [44, 63]. In particular NMI prefers a larger number of partitions. Moreover it has been  
951 observed that the NMI between two independent partitions can be much larger than zero. This feature  
952 makes NMI more difficult to interpret than for example ARI.

## D Additional results on real-world data

### D.1 Extended result tables

Tables 5, 6, 7 and 8 provide a detailed comparison of DIGRAC with spectral methods and InfoMap. Since no labeling information is available and all of the other competing GNN methods require labels, we do not compare DIGRAC with them on these real data sets.

In Tables 5, 6, 7 and 8, we report 12 combinations of global imbalance scores by data set. The naming convention of these imbalance scores is provided in Table 2. To assess how balanced our recovered clusters are in terms of sizes, we also report the size ratio, which is defined as the size of the largest predicted cluster to the smallest one, and the standard deviation of sizes, size std, in order to show how varied the sizes of predicted clusters are. For a relatively balanced clustering, we expect the latter two terms to be small.

**Table 5:** Performance comparison on *Telegram*. The best is marked in **bold red** and the second best is marked in underline blue.

Metric/Method	InfoMap	Bi_sym	DD_sym	DISG_LR	Herm	Herm_rw	DIGRAC
$\mathcal{O}^{\text{sort}}_{\text{vol\_sum}}$	0.04±0.00	<u>0.21±0.00</u>	<u>0.21±0.00</u>	<u>0.21±0.01</u>	0.20±0.01	0.14±0.00	<b>0.32±0.01</b>
$\mathcal{O}^{\text{sort}}_{\text{vol\_min}}$	0.47±0.00	<u>0.67±0.00</u>	0.61±0.00	0.66±0.02	0.66±0.02	0.19±0.00	<b>0.79±0.06</b>
$\mathcal{O}^{\text{sort}}_{\text{vol\_max}}$	0.03±0.00	<u>0.20±0.00</u>	<u>0.20±0.00</u>	<u>0.20±0.01</u>	0.19±0.01	0.12±0.00	<b>0.29±0.01</b>
$\mathcal{O}^{\text{sort}}_{\text{plain}}$	<b>1.00±0.00</b>	0.80±0.00	0.75±0.00	0.78±0.03	0.76±0.04	0.59±0.00	<u>0.96±0.01</u>
$\mathcal{O}^{\text{std}}_{\text{vol\_sum}}$	0.01±0.00	0.26±0.00	0.26±0.00	0.26±0.01	0.25±0.02	<b>0.35±0.00</b>	<u>0.28±0.01</u>
$\mathcal{O}^{\text{std}}_{\text{vol\_min}}$	0.16±0.00	<b>0.84±0.00</b>	0.76±0.00	<u>0.82±0.03</u>	<u>0.82±0.03</u>	0.49±0.00	0.73±0.03
$\mathcal{O}^{\text{std}}_{\text{vol\_max}}$	0.01±0.00	<u>0.25±0.00</u>	<u>0.25±0.00</u>	<u>0.25±0.01</u>	0.24±0.02	<b>0.29±0.00</b>	<u>0.25±0.01</u>
$\mathcal{O}^{\text{std}}_{\text{plain}}$	0.68±0.00	<b>1.00±0.00</b>	0.94±0.00	0.98±0.04	0.95±0.04	<u>0.99±0.00</u>	0.90±0.05
$\mathcal{O}^{\text{naive}}_{\text{vol\_sum}}$	0.01±0.00	<u>0.26±0.00</u>	<u>0.26±0.00</u>	<u>0.26±0.01</u>	0.25±0.02	0.23±0.00	<b>0.27±0.01</b>
$\mathcal{O}^{\text{naive}}_{\text{vol\_min}}$	0.11±0.00	<b>0.84±0.00</b>	0.76±0.00	<u>0.82±0.03</u>	<u>0.82±0.03</u>	0.32±0.00	0.72±0.04
$\mathcal{O}^{\text{naive}}_{\text{vol\_max}}$	0.00±0.00	<b>0.25±0.00</b>	<b>0.25±0.00</b>	<b>0.25±0.01</b>	0.24±0.02	0.20±0.00	0.24±0.01
$\mathcal{O}^{\text{naive}}_{\text{plain}}$	0.63±0.00	<b>1.00±0.00</b>	0.94±0.00	0.98±0.04	0.95±0.04	<u>0.99±0.00</u>	0.89±0.06
size ratio	<u>24.750</u>	242.000	242.000	242.000	242.00	53	<b>3.090</b>
size std	<u>35.57</u>	104.360	104.360	104.360	104.360	63.460	<b>26.39</b>

**Table 6:** Performance comparison on *Blog*. The best is marked in **bold red** and the second best is marked in underline blue.

Metric/Method	InfoMap	Bi_sym	DD_sym	DISG_LR	Herm	Herm_rw	DIGRAC
$\mathcal{O}^{\text{sort}}_{\text{vol\_sum}}$	0.07±0.00	0.07±0.00	0.00±0.00	0.05±0.00	<u>0.37±0.00</u>	0.00±0.00	<b>0.44±0.00</b>
$\mathcal{O}^{\text{sort}}_{\text{vol\_min}}$	0.02±0.00	0.33±0.00	0.05±0.00	0.31±0.00	<u>0.78±0.01</u>	<b>0.89±0.00</b>	0.76±0.00
$\mathcal{O}^{\text{sort}}_{\text{vol\_max}}$	0.05±0.00	0.05±0.00	0.00±0.00	0.04±0.00	<u>0.26±0.00</u>	0.00±0.00	<b>0.40±0.00</b>
$\mathcal{O}^{\text{sort}}_{\text{plain}}$	<b>1.00±0.00</b>	0.33±0.00	0.05±0.00	0.31±0.00	0.78±0.01	<u>0.89±0.00</u>	0.76±0.00
$\mathcal{O}^{\text{std}}_{\text{vol\_sum}}$	0.00±0.00	0.07±0.00	0.00±0.00	0.05±0.00	<u>0.37±0.00</u>	0.00±0.00	<b>0.44±0.00</b>
$\mathcal{O}^{\text{std}}_{\text{vol\_min}}$	0.00±0.00	0.33±0.00	0.05±0.00	0.31±0.00	<u>0.78±0.01</u>	<b>0.89±0.00</b>	0.76±0.00
$\mathcal{O}^{\text{std}}_{\text{vol\_max}}$	0.00±0.00	0.05±0.00	0.00±0.00	0.04±0.00	<u>0.26±0.00</u>	0.00±0.00	<b>0.40±0.00</b>
$\mathcal{O}^{\text{std}}_{\text{plain}}$	0.73±0.00	0.33±0.00	0.05±0.00	0.31±0.00	<u>0.78±0.01</u>	<b>0.89±0.00</b>	0.76±0.00
$\mathcal{O}^{\text{naive}}_{\text{vol\_sum}}$	0.00±0.00	0.07±0.00	0.00±0.00	0.05±0.00	<u>0.37±0.00</u>	0.00±0.00	<b>0.44±0.00</b>
$\mathcal{O}^{\text{naive}}_{\text{vol\_min}}$	0.00±0.00	0.33±0.00	0.05±0.00	0.31±0.00	<u>0.78±0.01</u>	<b>0.89±0.00</b>	0.76±0.00
$\mathcal{O}^{\text{naive}}_{\text{vol\_max}}$	0.00±0.00	0.05±0.00	0.00±0.00	0.04±0.00	<u>0.26±0.00</u>	0.00±0.00	<b>0.40±0.00</b>
$\mathcal{O}^{\text{naive}}_{\text{plain}}$	0.76±0.00	0.33±0.00	0.05±0.00	0.31±0.00	<u>0.78±0.01</u>	<b>0.89±0.00</b>	0.76±0.00
size ratio	<b>1.270</b>	8.700	2.450	6.100	11.93	44.26	<u>1.860</u>
size std	<b>64.50</b>	485	256.200	439	516.500	584	<u>183.20</u>

Tables 5, 6, 7, 8, 9, 10 and 11 reveal that DIGRAC provides competitive global imbalance scores in all of the 12 objectives introduced, and across all the real data sets, usually outperforming all the other methods. Among the tables, Table 11 provides results in terms of the distance to the best yearly performance, averaged across the 19 years; DIGRAC usually outperforms all the other methods across all the years. Note that Bi\_sym and DD\_sym are not able to generate results for *WikiTalk*, as

**Table 7:** Performance comparison on *Migration*. The best is marked in **bold red** and the second best is marked in underline blue. InfoMap results are omitted here as it predicts a single huge cluster and could not generate imbalance results.

Metric/Method	Bi_sym	DD_sym	DISG_LR	Herm	Herm_rw	DIGRAC
$\mathcal{O}_{vol\_sum}^{sort}$	0.03±0.00	0.01±0.00	0.02±0.00	<u>0.04±0.00</u>	0.02±0.00	<b>0.05±0.00</b>
$\mathcal{O}_{vol\_min}^{sort}$	<b>0.19±0.00</b>	0.08±0.00	0.08±0.00	0.15±0.02	0.05±0.00	<u>0.18±0.03</u>
$\mathcal{O}_{vol\_max}^{sort}$	<u>0.03±0.00</u>	0.01±0.00	0.01±0.00	<u>0.03±0.00</u>	0.02±0.00	<b>0.04±0.00</b>
$\mathcal{O}_{plain}^{sort}$	0.24±0.00	0.20±0.00	0.17±0.00	<u>0.40±0.01</u>	<b>0.49±0.06</b>	0.29±0.04
$\mathcal{O}_{vol\_sum}^{std}$	0.01±0.00	0.01±0.00	0.01±0.00	<u>0.02±0.00</u>	<u>0.02±0.00</u>	<b>0.04±0.01</b>
$\mathcal{O}_{vol\_min}^{std}$	<u>0.10±0.00</u>	0.05±0.00	0.05±0.00	0.08±0.01	0.04±0.00	<b>0.16±0.03</b>
$\mathcal{O}_{vol\_max}^{std}$	<u>0.01±0.00</u>	<u>0.01±0.00</u>	<u>0.01±0.00</u>	<u>0.01±0.00</u>	<u>0.01±0.00</u>	<b>0.03±0.01</b>
$\mathcal{O}_{plain}^{std}$	0.13±0.00	0.12±0.00	0.11±0.00	<u>0.20±0.01</u>	<u>0.20±0.01</u>	<b>0.26±0.01</b>
$\mathcal{O}_{naive\_vol\_sum}$	0.01±0.00	0.01±0.00	0.01±0.00	<u>0.02±0.00</u>	0.01±0.00	<b>0.04±0.01</b>
$\mathcal{O}_{naive\_vol\_min}$	<u>0.09±0.00</u>	0.04±0.00	0.04±0.00	0.08±0.01	0.01±0.00	<b>0.16±0.03</b>
$\mathcal{O}_{naive\_vol\_max}$	<u>0.01±0.00</u>	0.00±0.00	<u>0.01±0.00</u>	<u>0.01±0.00</u>	0.00±0.00	<b>0.03±0.01</b>
$\mathcal{O}_{naive\_plain}$	0.12±0.00	0.10±0.00	0.08±0.00	<u>0.19±0.00</u>	<u>0.19±0.03</u>	<b>0.26±0.01</b>
size ratio	7.780	6.070	<b>4.360</b>	36.05	1035.90	<u>4.420</u>
size std	135.210	<u>132.76</u>	<b>103.43</b>	335.790	353.060	264.500

**Table 8:** Performance comparison on *WikiTalk*. The best is marked in **bold red** and the second best is marked in underline blue. InfoMap results are omitted here as its large number of predicted clusters leads to memory error in imbalance calculation.

Metric/Method	DISG_LR	Herm	Herm_rw	DIGRAC
$\mathcal{O}_{vol\_sum}^{sort}$	<u>0.18±0.03</u>	0.15±0.02	0.00±0.00	<b>0.24±0.05</b>
$\mathcal{O}_{vol\_min}^{sort}$	0.10±0.03	0.22±0.05	<u>0.26±0.00</u>	<b>0.28±0.13</b>
$\mathcal{O}_{vol\_max}^{sort}$	<u>0.16±0.03</u>	0.09±0.01	0.00±0.00	<b>0.19±0.04</b>
$\mathcal{O}_{plain}^{sort}$	0.87±0.08	<u>0.99±0.01</u>	0.98±0.00	<b>1.00±0.00</b>
$\mathcal{O}_{vol\_sum}^{std}$	<b>0.17±0.04</b>	0.06±0.01	0.01±0.00	<u>0.14±0.02</u>
$\mathcal{O}_{vol\_min}^{std}$	0.09±0.02	0.09±0.02	<b>0.27±0.00</b>	<u>0.18±0.08</u>
$\mathcal{O}_{vol\_max}^{std}$	<b>0.15±0.04</b>	0.04±0.00	0.00±0.00	<u>0.11±0.02</u>
$\mathcal{O}_{plain}^{std}$	0.72±0.03	0.70±0.05	<b>0.98±0.00</b>	<u>0.84±0.06</u>
$\mathcal{O}_{naive\_vol\_sum}$	<u>0.10±0.02</u>	0.04±0.00	0.00±0.00	<b>0.12±0.01</b>
$\mathcal{O}_{naive\_vol\_min}$	0.06±0.03	0.07±0.02	<b>0.26±0.00</b>	<u>0.15±0.07</u>
$\mathcal{O}_{naive\_vol\_max}$	<b>0.09±0.02</b>	0.03±0.00	0.00±0.00	<b>0.09±0.01</b>
$\mathcal{O}_{naive\_plain}$	0.64±0.04	0.61±0.04	<b>0.98±0.00</b>	<u>0.76±0.06</u>
size ratio	1190162.25	2217434.50	<b>250.48</b>	<u>71765.14</u>
size std	713813.72	660060.33	<u>657941.88</u>	<b>643220.37</b>

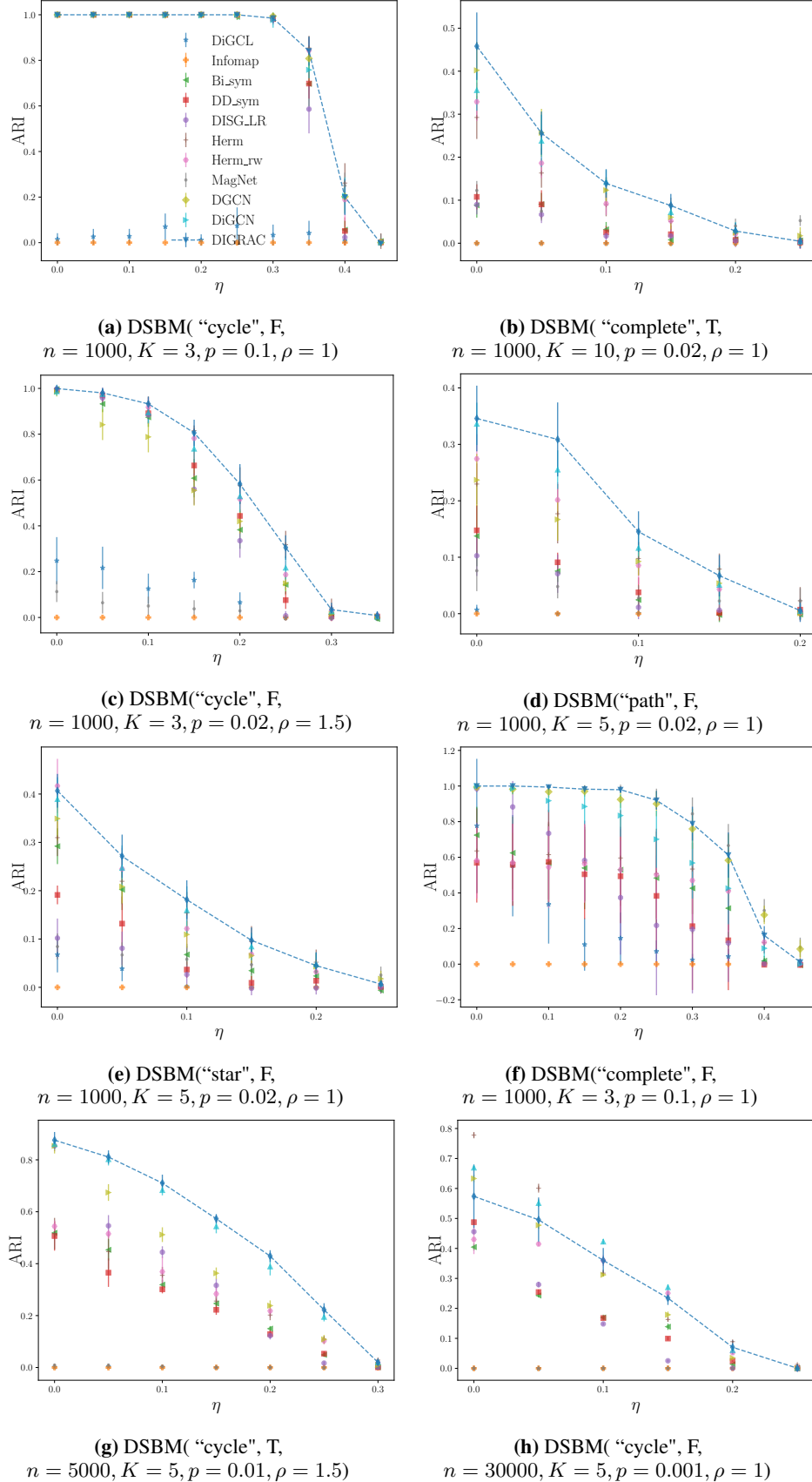
969 large  $n \times n$  matrix multiplication with its transpose causes memory issue, when  $n = 2,388,953$ .  
 970 Small values of the size ratio and size standard deviation suggest that the normalization in the loss  
 971 function penalizes tiny clusters, and that DIGRAC tends to predict balanced cluster sizes.

**Table 9:** Performance comparison on *Lead-Lag* for year 2015. The best is marked in **bold red** and the second best is marked in **underline blue**. InfoMap results are omitted here as it usually predicts a single huge cluster and could not generate imbalance results.

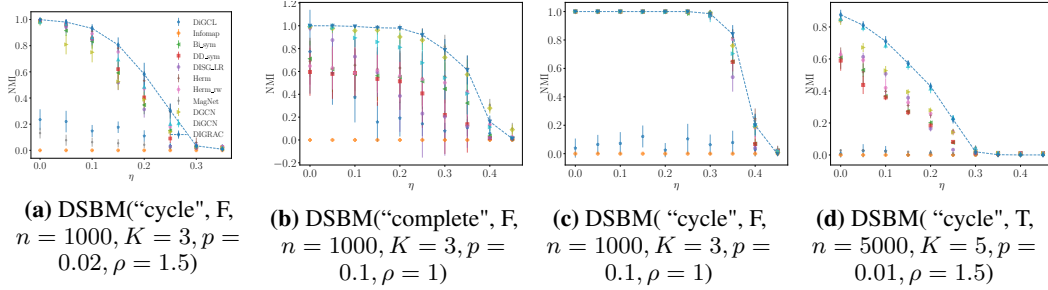
Metric/Method	Bi_sym	DD_sym	DISG_LR	Herm	Herm_rw	DIGRAC
$\mathcal{O}^{\text{sort}}_{\text{vol\_sum}}$	<u>0.07±0.00</u>	<u>0.07±0.00</u>	0.06±0.00	<u>0.07±0.00</u>	0.06±0.01	<b>0.15±0.00</b>
$\mathcal{O}^{\text{sort}}_{\text{vol\_min}}$	<b>0.53±0.06</b>	<u>0.50±0.02</u>	0.45±0.07	<u>0.50±0.03</u>	0.46±0.06	<u>0.50±0.02</u>
$\mathcal{O}^{\text{sort}}_{\text{vol\_max}}$	<u>0.07±0.00</u>	0.06±0.00	0.06±0.00	0.06±0.00	0.06±0.00	<b>0.15±0.01</b>
$\mathcal{O}^{\text{sort}}_{\text{plain}}$	<u>0.65±0.03</u>	<b>0.67±0.03</b>	0.59±0.03	<u>0.65±0.03</u>	<u>0.65±0.02</u>	0.55±0.07
$\mathcal{O}^{\text{std}}_{\text{vol\_sum}}$	<u>0.04±0.00</u>	<u>0.04±0.00</u>	<u>0.04±0.00</u>	<u>0.04±0.00</u>	<u>0.04±0.00</u>	<b>0.11±0.02</b>
$\mathcal{O}^{\text{std}}_{\text{vol\_min}}$	<u>0.27±0.03</u>	<u>0.27±0.02</u>	0.24±0.02	<u>0.27±0.02</u>	0.26±0.04	<b>0.35±0.04</b>
$\mathcal{O}^{\text{std}}_{\text{vol\_max}}$	<u>0.03±0.00</u>	<u>0.03±0.00</u>	<u>0.03±0.00</u>	<u>0.03±0.00</u>	<u>0.03±0.00</u>	<b>0.10±0.02</b>
$\mathcal{O}^{\text{std}}_{\text{plain}}$	<u>0.39±0.02</u>	<u>0.39±0.01</u>	0.37±0.02	<u>0.39±0.02</u>	<b>0.40±0.02</b>	0.38±0.04
$\mathcal{O}^{\text{naive}}_{\text{vol\_sum}}$	<u>0.03±0.00</u>	<u>0.03±0.00</u>	<u>0.03±0.00</u>	<u>0.03±0.00</u>	<u>0.03±0.00</u>	<b>0.08±0.03</b>
$\mathcal{O}^{\text{naive}}_{\text{vol\_min}}$	<u>0.20±0.02</u>	<u>0.20±0.02</u>	0.17±0.03	<u>0.20±0.02</u>	<u>0.20±0.03</u>	<b>0.25±0.08</b>
$\mathcal{O}^{\text{naive}}_{\text{vol\_max}}$	0.02±0.00	<u>0.03±0.00</u>	0.02±0.00	<u>0.03±0.00</u>	<u>0.03±0.00</u>	<b>0.08±0.03</b>
$\mathcal{O}^{\text{naive}}_{\text{plain}}$	0.29±0.01	0.29±0.01	0.26±0.02	<u>0.30±0.01</u>	<u>0.30±0.01</u>	<b>0.31±0.05</b>
size ratio	3.070	3.110	3.060	<b>2.89</b>	<u>2.95</u>	15.640
size std	8.390	<u>7.94</u>	8.680	<b>7.28</b>	8.050	18.680

**Table 10:** Performance comparison on *Lead-Lag*. Results in each year is averaged over ten runs. Mean and standard deviation (after ±) are calculated over the 19 years. The best is marked in **bold red** and the second best is marked in **underline blue**. InfoMap results are omitted here as it usually predicts a single huge cluster and could not generate imbalance results.

Metric/Method	Bi_sym	DD_sym	DISG_LR	Herm	Herm_rw	DIGRAC
$\mathcal{O}^{\text{sort}}_{\text{vol\_sum}}$	<u>0.07±0.01</u>	<u>0.07±0.01</u>	<u>0.07±0.01</u>	<u>0.07±0.02</u>	<u>0.07±0.02</u>	<b>0.15±0.03</b>
$\mathcal{O}^{\text{sort}}_{\text{vol\_min}}$	<b>0.51±0.10</b>	0.48±0.09	0.47±0.10	<b>0.51±0.11</b>	0.50±0.10	0.47±0.09
$\mathcal{O}^{\text{sort}}_{\text{vol\_max}}$	<u>0.07±0.01</u>	0.06±0.01	0.06±0.01	<u>0.07±0.01</u>	<u>0.07±0.01</u>	<b>0.14±0.03</b>
$\mathcal{O}^{\text{sort}}_{\text{plain}}$	<b>0.66±0.09</b>	0.64±0.08	0.63±0.08	<b>0.66±0.09</b>	0.65±0.09	0.53±0.09
$\mathcal{O}^{\text{std}}_{\text{vol\_sum}}$	<u>0.04±0.01</u>	<u>0.04±0.01</u>	<u>0.04±0.01</u>	<u>0.04±0.01</u>	<u>0.04±0.01</u>	<b>0.12±0.03</b>
$\mathcal{O}^{\text{std}}_{\text{vol\_min}}$	<u>0.27±0.04</u>	<u>0.27±0.04</u>	0.25±0.04	<u>0.27±0.03</u>	<u>0.27±0.03</u>	<b>0.38±0.07</b>
$\mathcal{O}^{\text{std}}_{\text{vol\_max}}$	<u>0.04±0.00</u>	0.03±0.00	0.03±0.00	0.03±0.00	0.03±0.00	<b>0.11±0.02</b>
$\mathcal{O}^{\text{std}}_{\text{plain}}$	<u>0.40±0.05</u>	0.39±0.05	0.38±0.05	<u>0.40±0.05</u>	<u>0.40±0.05</u>	<b>0.44±0.07</b>
$\mathcal{O}^{\text{naive}}_{\text{vol\_sum}}$	<u>0.03±0.01</u>	<u>0.03±0.01</u>	<u>0.03±0.01</u>	<u>0.03±0.01</u>	<u>0.03±0.01</u>	<b>0.08±0.04</b>
$\mathcal{O}^{\text{naive}}_{\text{vol\_min}}$	<u>0.20±0.05</u>	0.19±0.05	0.18±0.05	0.19±0.04	0.19±0.04	<b>0.26±0.10</b>
$\mathcal{O}^{\text{naive}}_{\text{vol\_max}}$	<u>0.03±0.01</u>	0.02±0.01	0.02±0.01	0.02±0.00	0.02±0.00	<b>0.08±0.03</b>
$\mathcal{O}^{\text{naive}}_{\text{plain}}$	<u>0.30±0.06</u>	0.28±0.06	0.27±0.06	0.29±0.05	0.29±0.05	<b>0.32±0.11</b>
size ratio	<u>3.67</u>	<b>3.34</b>	3.900	4.110	3.880	8.070
size std	<u>9.31</u>	<b>9.14</b>	10.090	10.490	10.360	17.060



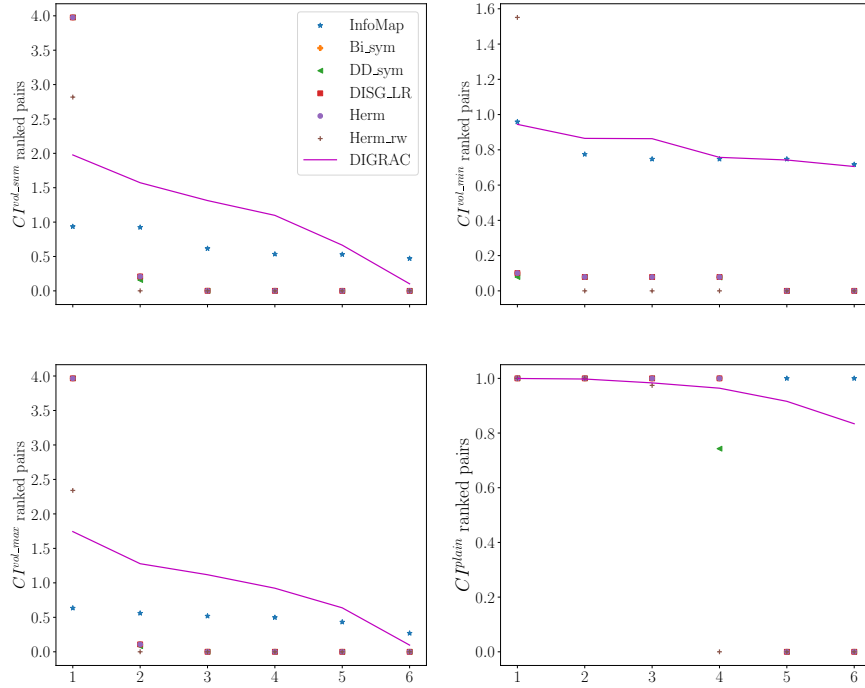
**Figure 13:** Test ARI comparison on synthetic data. Dashed lines highlight DIGRAC's performance. Error bars are given by one standard error.



**Figure 14:** Test NMI comparison on some synthetic data. Dashed lines highlight DIGRAC's performance. Error bars are given by one standard error.

**Table 11:** Performance comparison on *Lead-Lag*, where we evaluate the performance distance to the best one in each year. Results in each year is averaged over ten runs. Mean and standard deviation (after  $\pm$ ) are calculated over the 19 years. The best is marked in **bold red** and the second best is marked in underline blue. InfoMap results are omitted here as it usually predicts a single huge cluster and could not generate imbalance results.

Metric/Method	Bi_sym	DD_sym	DISG_LR	Herm	Herm_rw	DIGRAC
$\mathcal{O}^{\text{sort}}_{\text{vol\_sum}}$	<u>0.07±0.02</u>	0.08±0.02	0.08±0.02	<u>0.07±0.02</u>	<u>0.07±0.02</u>	<b>0.00±0.00</b>
$\mathcal{O}^{\text{sort}}_{\text{vol\_min}}$	<b>0.01±0.01</b>	0.05±0.03	0.06±0.03	<u>0.02±0.02</u>	<u>0.02±0.02</u>	0.06±0.04
$\mathcal{O}^{\text{sort}}_{\text{vol\_max}}$	<u>0.07±0.02</u>	<u>0.07±0.02</u>	<u>0.07±0.02</u>	<u>0.07±0.02</u>	<u>0.07±0.02</u>	<b>0.00±0.00</b>
$\mathcal{O}^{\text{sort}}_{\text{plain}}$	<b>0.01±0.02</b>	0.03±0.03	0.05±0.03	<b>0.01±0.02</b>	0.02±0.02	0.14±0.03
$\mathcal{O}^{\text{std}}_{\text{vol\_sum}}$	<u>0.08±0.02</u>	<u>0.08±0.02</u>	<u>0.08±0.02</u>	<u>0.08±0.02</u>	<u>0.08±0.02</u>	<b>0.00±0.00</b>
$\mathcal{O}^{\text{std}}_{\text{vol\_min}}$	<u>0.10±0.05</u>	0.11±0.04	0.13±0.05	0.11±0.05	0.11±0.05	<b>0.00±0.00</b>
$\mathcal{O}^{\text{std}}_{\text{vol\_max}}$	<u>0.07±0.02</u>	0.08±0.02	0.08±0.02	0.08±0.02	0.08±0.02	<b>0.00±0.00</b>
$\mathcal{O}^{\text{std}}_{\text{plain}}$	<u>0.04±0.03</u>	0.05±0.04	0.06±0.04	<u>0.04±0.04</u>	<u>0.04±0.03</u>	<b>0.00±0.00</b>
$\mathcal{O}^{\text{naive}}_{\text{vol\_sum}}$	<u>0.05±0.03</u>	0.06±0.03	0.06±0.03	<u>0.05±0.03</u>	<u>0.05±0.03</u>	<b>0.00±0.00</b>
$\mathcal{O}^{\text{naive}}_{\text{vol\_min}}$	<u>0.06±0.07</u>	0.07±0.06	0.08±0.07	0.07±0.08	0.07±0.08	<b>0.00±0.00</b>
$\mathcal{O}^{\text{naive}}_{\text{vol\_max}}$	<u>0.05±0.03</u>	<u>0.05±0.03</u>	<u>0.05±0.03</u>	<u>0.05±0.03</u>	<u>0.05±0.03</u>	<b>0.00±0.00</b>
$\mathcal{O}^{\text{naive}}_{\text{plain}}$	<u>0.03±0.06</u>	0.05±0.05	0.06±0.06	0.04±0.06	0.04±0.06	<b>0.01±0.02</b>
size ratio	<u>1.04</u>	<b>0.71</b>	1.270	1.480	1.250	5.440
size std	<u>0.58</u>	<b>0.41</b>	1.360	1.770	1.630	8.340

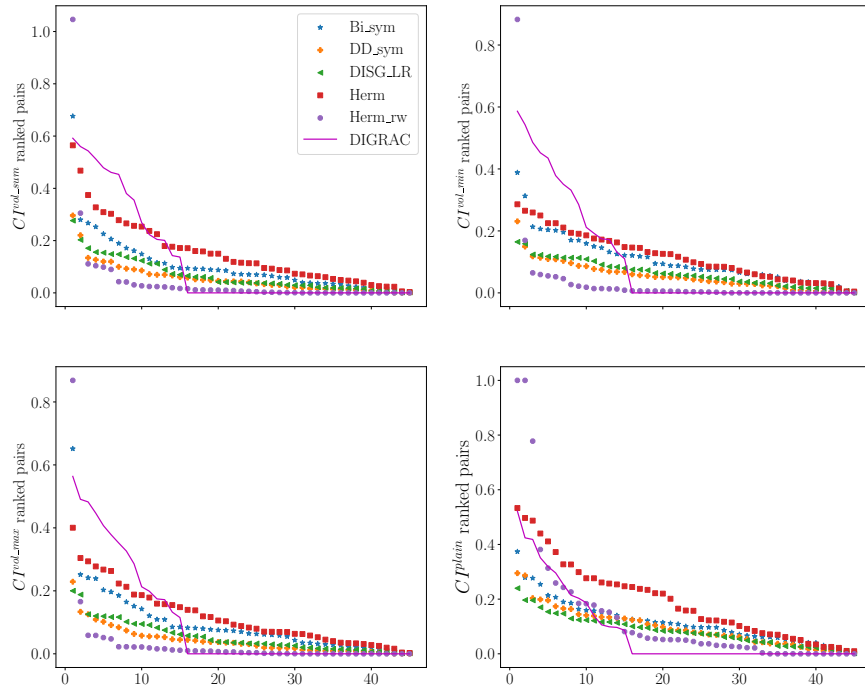


**Figure 15:** Ranked pairs of pairwise imbalance recovered by comparing methods for different choices of normalization on the *Telegram* data set. Lines are used to highlight DIGRAC’s performance.

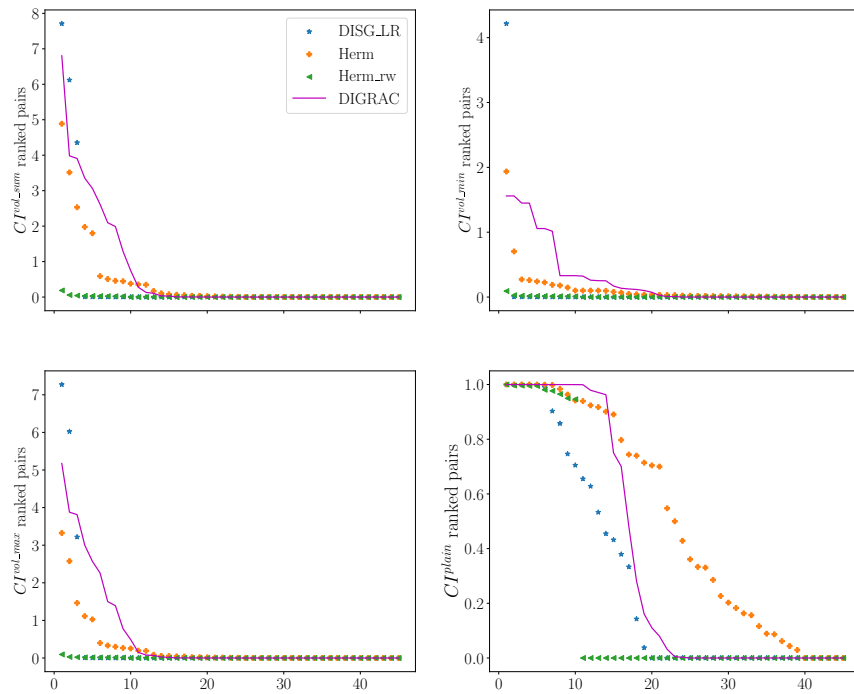
972 **D.2 Ranked pairwise imbalance scores**

973 We also plot the ranked pairwise imbalance scores for all data sets except *Blog*, which has only one  
 974 possible pairwise imbalance score. For *Lead-Lag*, we only plot the year 2015 as an example; the plots  
 975 for the other years are similar. Figures 15, 16, 17 and 18 illustrate that DIGRAC is able to provide  
 976 comparable or higher pairwise imbalance scores for the leading pairs, especially on  $CI^{vol\_min}$  pairs.  
 977 We also observe that except for  $CI^{plain}$ , DIGRAC has a less rapid drop in pairwise imbalance scores  
 978 after the first leading pair compared to Herm and Herm\_rw, which can have a few pairs with higher  
 979 imbalance scores than DIGRAC.

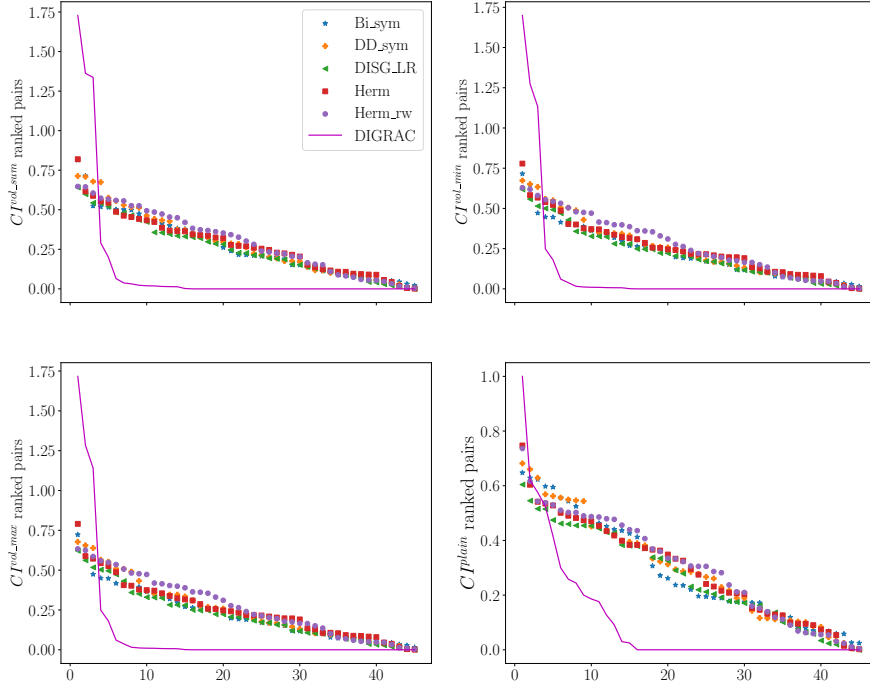




**Figure 16:** Ranked pairs of pairwise imbalance recovered by comparing methods for different choices of normalization on the *Migration* data set. Lines are used to highlight DIGRAC’s performance. InfoMap results are omitted as it predicts one single huge cluster and could not produce imbalance results.



**Figure 17:** Ranked pairs of pairwise imbalance recovered by comparing methods for different choices of normalization on *WikiTalk* data set. Lines are used to highlight DIGRAC’s performance. InfoMap results are omitted here because it triggers memory error due to the large number of predicted clusters.



**Figure 18:** Ranked pairs of pairwise imbalance recovered by comparing methods for different choices of normalization on *Lead-Lag* data set. Lines are used to highlight DIGRAC’s performance. InfoMap results are omitted here because it only predicts a single cluster.

### 980 D.3 Predicted meta-graph flow matrix plots

981 For each data set, we plot the predicted meta-graph flow matrix  $\mathbf{F}'$  defined in Eq. (12).

982 From Fig. 19, we conclude that DIGRAC is able to recover a directed flow imbalance between  
 983 clusters in all of the selected data sets. Fig. 19a shows a clear cut imbalance between two clusters,  
 984 possibly corresponding to the Republican and Democratic parties. Fig. 19b plots imbalance flows in  
 985 the real data set *Telegram*, where cluster 3 is a core-transient cluster, cluster 0 is a core-sink cluster,  
 986 cluster 2 is a periphery-upstream cluster, while cluster 1 is a periphery-downstream cluster [3, 5]. For  
 987 *WikiTalk*, illustrated in Fig. 19d, the lower-triangular part entries are typically source nodes for edges,  
 988 while the upper-triangular part are target nodes. For *Lead-Lag*, taking the year 2015 as an example,  
 989 DIGRAC is also able to recover high imbalance in the data.

990 We also note that DIGRAC would not necessarily predict the same number of clusters as assumed, so  
 991 that we do not need to specify the exact number of clusters before training DIGRAC; specifying the  
 992 maximum number of possible clusters suffices.

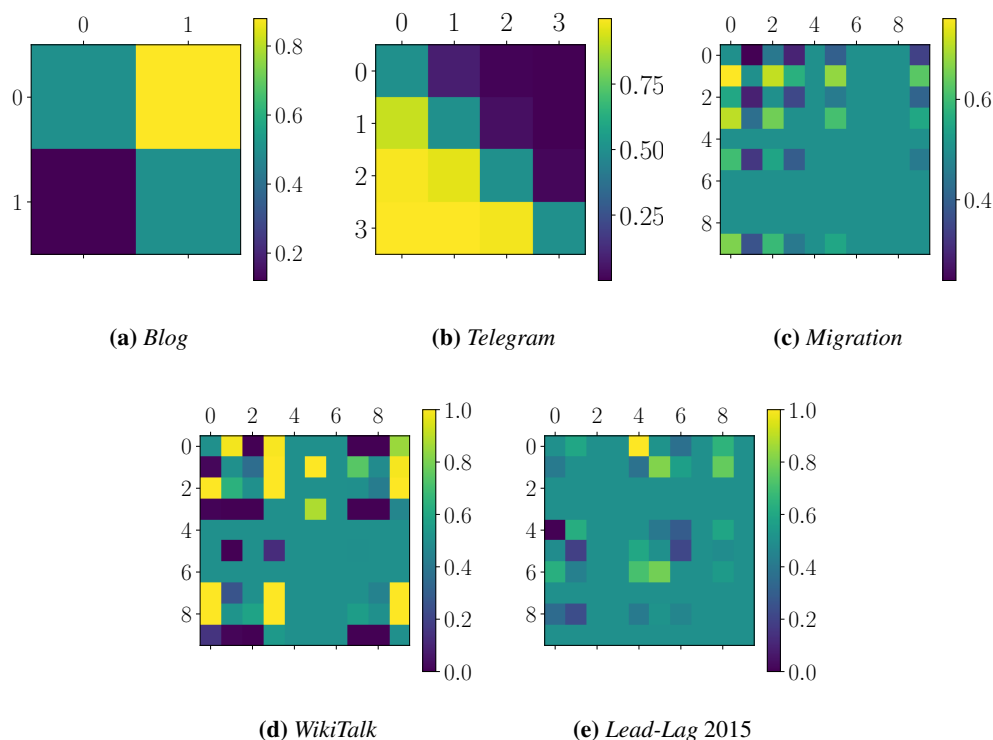


Figure 19: Predicted meta-graph flow matrix from DIGRAC of five real-world data sets.

#### 993 D.4 Migration plots

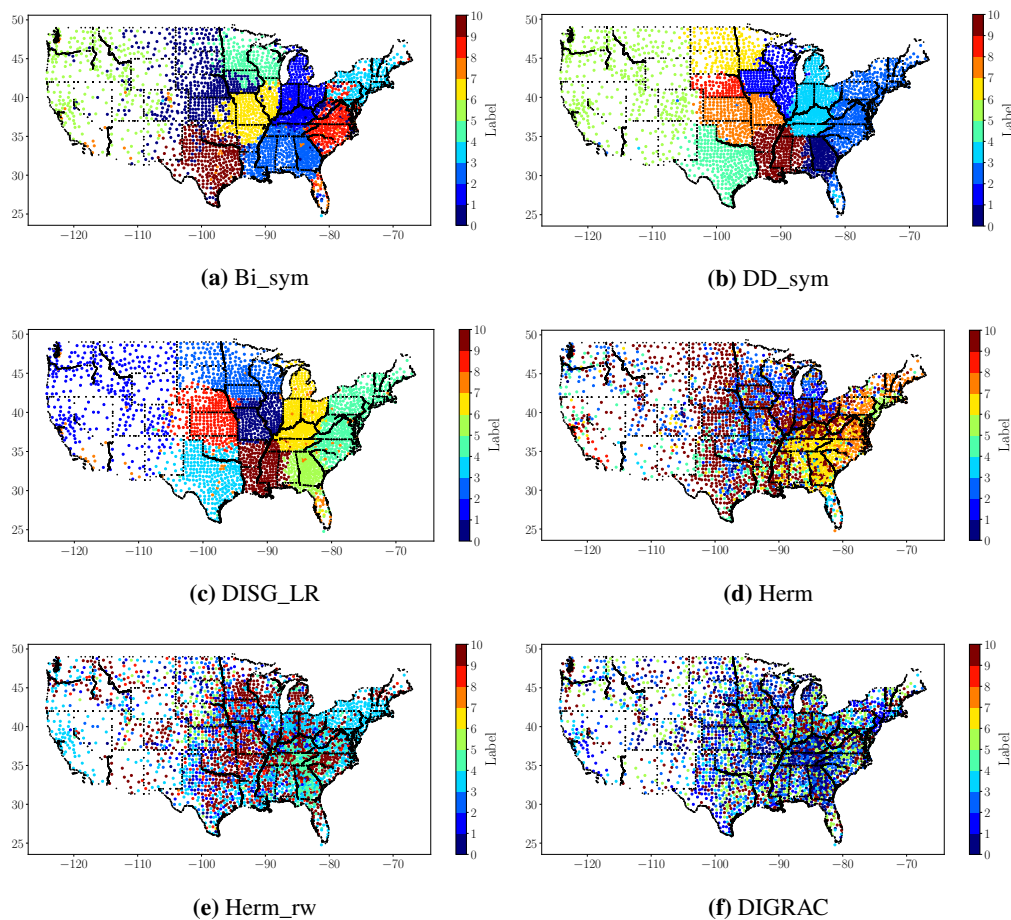
994 We compare DIGRAC to five spectral methods for recovering clusters for the US migration data set,  
995 and plot the recovered clusters on a map, in Fig. 20.

996 The visualization in Fig. (a-c) shows that clusters align particularly well with the political and  
997 administrative boundaries of the US states, as previously observed in [64]. This outcome is not  
998 deemed too insightful, as it trivially reveals the fact there is significant intra-state and inter-state  
999 migration, and does not uncover any of the information on latent migration patterns between far-  
1000 away states, and more generally, between regions which are not necessarily geographically cohesive.  
1001 DIGRAC outcomes, however, reveal nontrivial migration patterns, for example migration from New  
1002 York to Florida, and from California to Arizona, which is consistent with the patterns discovered by  
1003 [4]. Fig. 21 details on the top pair migration patterns uncovered by DIGRAC.

#### 1004 D.5 Coping with outliers

1005 As mentioned in Section C.3, the preprocessing step to use ratio of migration instead of absolute  
1006 migration numbers is a way to cope with outliers (here, extremely large entries in the original digraph)  
1007 in *Migration*. To validate the effectiveness of this approach to cope with outliers, Table 12 provides  
1008 imbalance results for *Migration* when we do not transform the nonzero entries into ratios. Comparing  
1009 with Table 7, we witness an overall decrease in the performance. In this case InfoMap no longer  
1010 predicts a single huge cluster. However, its predicted number of clusters is about 44, which is too  
1011 large. This also implies that InfoMap is very sensitive to the magnitude of digraph entries, while  
1012 DIGRAC is not. Indeed, InfoMap gives 43 (too many) clusters for *Blog*, 19 (too many) for *Telegram*,  
1013 1 (too small) for *Migration*, and 17498 (far too many) for *WikiTalk*.

1014 We compare DIGRAC to five spectral methods as well as InfoMap for recovering clusters for the US  
1015 migration data set without the preprocessing step discussed earlier, and plot the recovered clusters on  
1016 a map in Fig. 22. Note that all methods, except DIGRAC, recover either clusters which are trivially  
1017 small in size or contain one very large dominant cluster (as in (a), (b), (e) and to some extent, also  
1018 (f)). The DISG\_LR clustering and InfoMap clustering provide clear geographic boundaries, but were



**Figure 20:** US migration predicted clusters, along with the geographic locations of the counties as well as state boundaries (in black). InfoMap results are omitted here because it only produces one huge cluster. The input data is normalized, following Eq. (13).

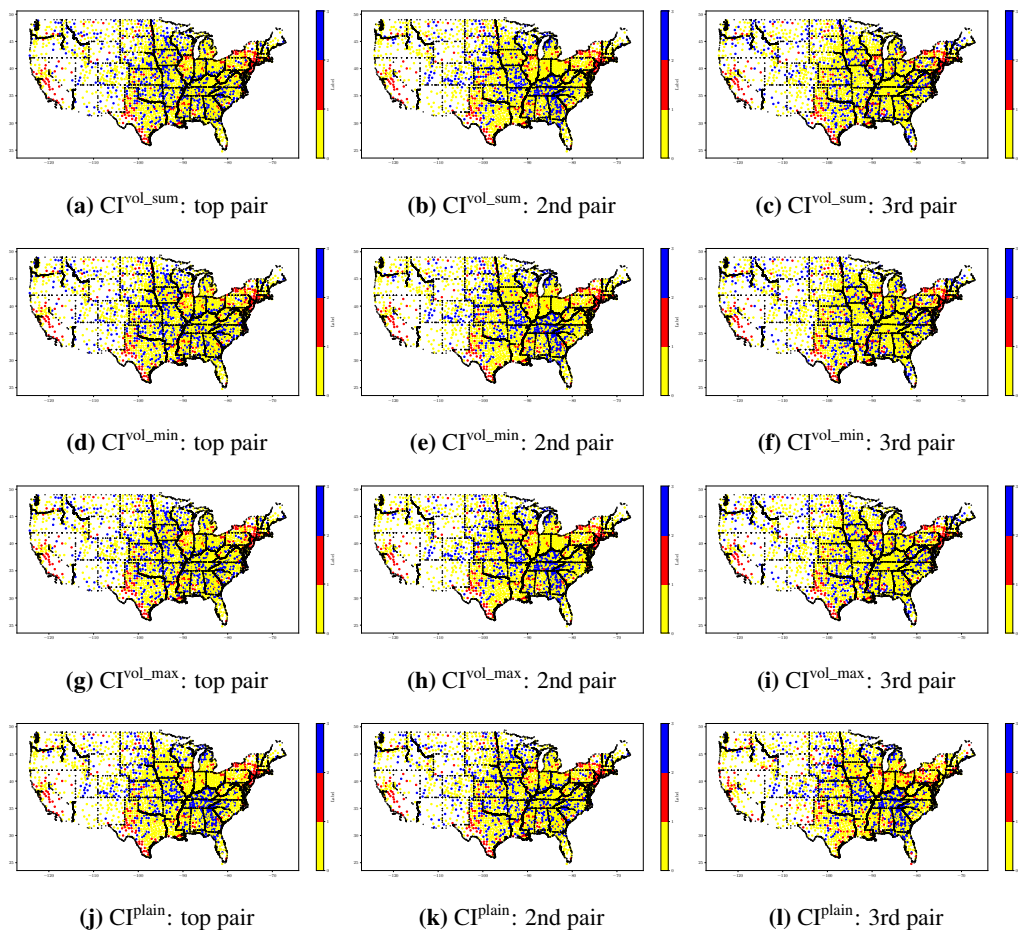
1019 not able to recover the imbalance among clusters. Other spectral methods generally have a dominant  
 1020 cluster containing most of the nodes, whereas DIGRAC has more balanced cluster sizes.

1021 When employing methods that symmetrize the adjacency matrix (as in (a) and (b)), the migration  
 1022 flows between counties in different states will be lost in the process. Furthermore, the visualization in  
 1023 Fig. (c) shows that clusters align particularly well with the political and administrative boundaries  
 1024 of the US states, as previously observed in [64]. The same is for Fig. (d). This outcome is not  
 1025 deemed very insightful, as it trivially reveals the fact that there is significant intra-state and inter-state  
 1026 migration, and does not uncover any of the information on latent migration patterns between far-away  
 1027 states, and more generally, between regions which are not necessarily geographically cohesive.

1028 Fig. 21 further plots the top three pairs of clusters based on four different imbalance scores given by  
 1029 DIGRAC. As shown in the figure, DIGRAC uncovers the migration trend from coastal to interior,  
 1030 across states. This trend of the directed flow agrees with that discussed in [4], with many people  
 1031 migrating from New York and California to the interior states.

## 1032 E Discussion of related methods that are not compared against in the main 1033 text

1034 To further emphasize the importance of directionality, our synthetic data sets have no difference in  
 1035 density between clusters; their sole signal is in the directionality of the edges. If all edge directions



**Figure 21:** US migration predicted cluster pairs with top imbalance, along with the geographic locations of the counties as well as state boundaries (in black). Red (label 1) is the sending cluster while blue (label 2) is the receiving cluster. Yellow (label 0) denotes all the other locations being considered. Subcaptions show the imbalance score and the rank based on that score.

1036 were to be removed, then no algorithm should be available to detect the clusters. To further support our  
 1037 claim why some methods mentioned in Section 2 in the main text are not appropriate for comparison,  
 1038 we have applied the default setting versions of the Louvain method [35], the Leiden algorithm [36]  
 1039 and OSLOM [34], to our synthetic data sets, and find that they do not detect the structure in the data,  
 1040 with ARI and NMI values very close to zero, and very low imbalance values. In particular, Louvain  
 1041 and Leiden tend to give a larger number of clusters than the ground truth which is designed to have  
 1042 small cluster sizes. OSLOM outputs clusters with extreme sizes, either a huge cluster containing  
 1043 (almost always) all the nodes, or every node forming a cluster by itself. To further demonstrate that  
 1044 comparing DIGRAC with density-based methods is unfair, We report the test ARI results for Infomap,  
 1045 Louvain and Leiden in Figure 23. We can see that Infomap, Louvain and Leiden normally produces  
 1046 zero-zero ARI values, which are much worse than the results from DIGRAC given in Figure 3.

1047 On the real-world data sets, these methods often give numbers of clusters that do not match our  
 1048 expectations. (*Blog* has two underlying parties, *Telegram* has a four-cluster core-periphery struc-  
 1049 ture). Louvain clusters nodes from *Blog* into 8-13 clusters (too many), *Telegram* into 4-5 clusters  
 1050 (acceptable), *Migration* into 5-7 clusters (acceptable), *WikiTalk* into 150-219 clusters (too many), and  
 1051 *Lead-Lag* into 10-55 clusters (acceptable or a bit too many). Leiden gives 12 (too many) clusters  
 1052 for *Blog*, 4-5 for *Telegram*, 5-6 for *Migration*, 170-248 (too many) for *WikiTalk*, and 10-55 clusters  
 1053 (acceptable or a bit too many) for *Lead-Lag*. OSLOM gives 6 clusters for *Blog* (too many), 16 for  
 1054 *Telegram* (too many), and 46 for *Migration* (too many). It could not generate results for *WikiTalk*

**Table 12:** Performance comparison on *Migration* (without preprocessing). The best is marked in **bold red** and the second best is marked in underline blue.

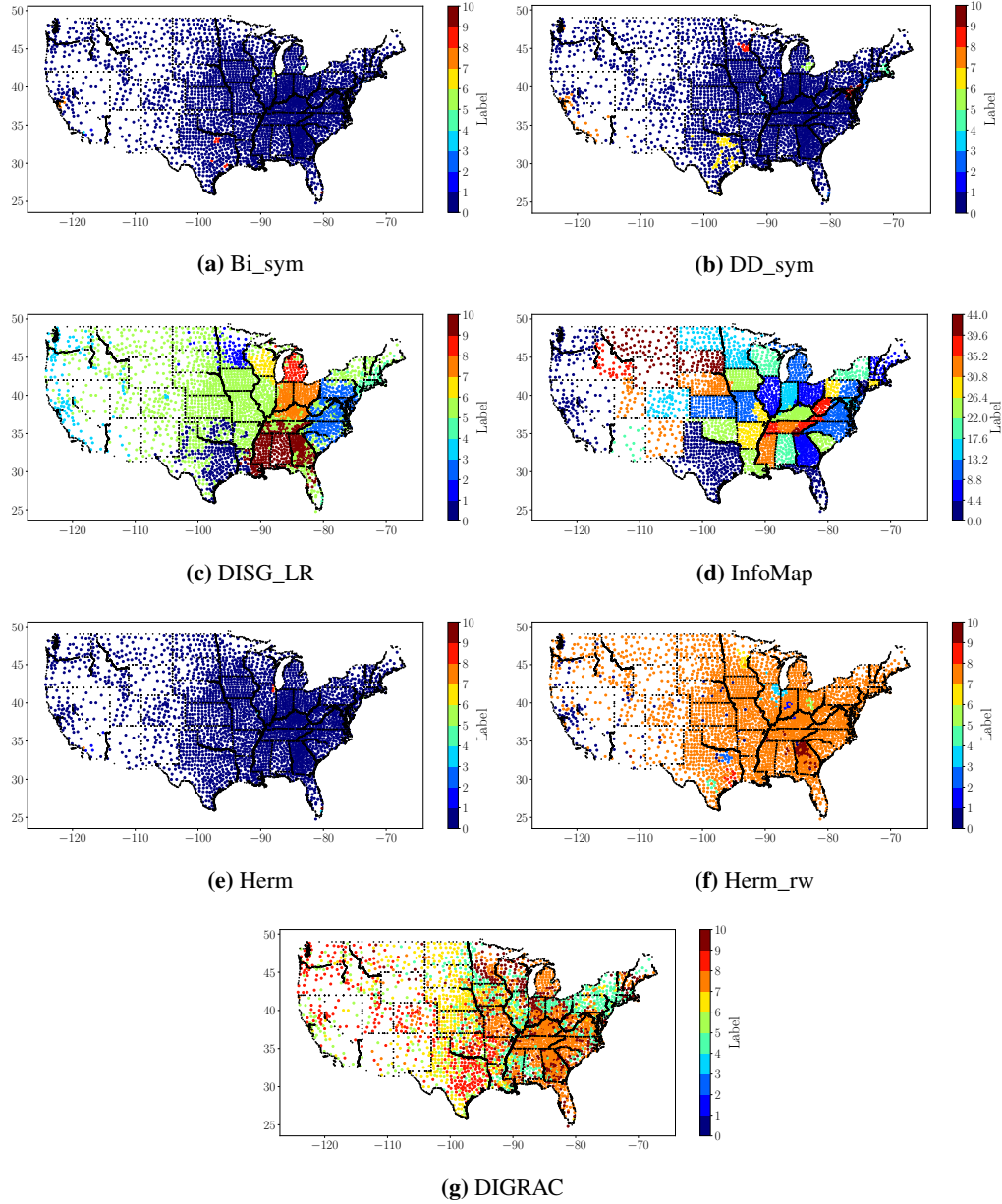
Metric/Method	InfoMap	Bi_sym	DD_sym	DISG_LR	Herm	Herm_rw	DIGRAC
$\mathcal{O}^{\text{sort}}_{\text{vol\_sum}}$	0.02±0.00	0.03±0.00	0.01±0.00	0.01±0.00	<b>0.07±0.00</b>	0.01±0.00	<u>0.04±0.00</u>
$\mathcal{O}^{\text{sort}}_{\text{vol\_min}}$	<b>0.24±0.00</b>	0.20±0.01	0.12±0.02	0.14±0.00	<u>0.21±0.01</u>	0.05±0.02	0.18±0.02
$\mathcal{O}^{\text{sort}}_{\text{vol\_max}}$	0.02±0.00	0.03±0.00	0.01±0.00	0.01±0.00	<b>0.06±0.00</b>	0.00±0.00	<u>0.04±0.00</u>
$\mathcal{O}^{\text{sort}}_{\text{plain}}$	<u>0.61±0.00</u>	0.46±0.00	0.29±0.02	0.26±0.00	<b>0.62±0.02</b>	0.40±0.00	0.32±0.11
$\mathcal{O}^{\text{std}}_{\text{vol\_sum}}$	0.00±0.00	0.01±0.00	0.00±0.00	0.00±0.00	<u>0.02±0.00</u>	0.00±0.00	<b>0.03±0.01</b>
$\mathcal{O}^{\text{std}}_{\text{vol\_min}}$	0.03±0.00	<u>0.09±0.00</u>	0.04±0.01	0.05±0.00	0.08±0.01	0.02±0.01	<b>0.11±0.03</b>
$\mathcal{O}^{\text{std}}_{\text{vol\_max}}$	0.00±0.00	<u>0.01±0.00</u>	0.00±0.00	0.00±0.00	<u>0.01±0.00</u>	0.00±0.00	<b>0.02±0.01</b>
$\mathcal{O}^{\text{std}}_{\text{plain}}$	0.19±0.00	0.23±0.00	0.14±0.01	0.12±0.00	<b>0.32±0.01</b>	<u>0.25±0.01</u>	0.21±0.03
$\mathcal{O}^{\text{naive}}_{\text{vol\_sum}}$	0.00±0.00	0.01±0.00	0.00±0.00	0.00±0.00	<u>0.02±0.00</u>	0.00±0.00	<b>0.03±0.01</b>
$\mathcal{O}^{\text{naive}}_{\text{vol\_min}}$	0.02±0.00	<u>0.08±0.00</u>	0.04±0.01	0.05±0.00	<u>0.08±0.01</u>	0.02±0.01	<b>0.11±0.04</b>
$\mathcal{O}^{\text{naive}}_{\text{vol\_max}}$	0.00±0.00	<u>0.01±0.00</u>	0.00±0.00	0.00±0.00	<u>0.01±0.00</u>	0.00±0.00	<b>0.02±0.01</b>
$\mathcal{O}^{\text{naive}}_{\text{plain}}$	0.16±0.00	<u>0.22±0.00</u>	0.13±0.01	0.11±0.00	<b>0.31±0.01</b>	<u>0.22±0.00</u>	0.21±0.03
size ratio	8.500	<u>3043.80</u>	722.620	25.780	<b>3059.20</b>	415.880	203.230
size std	<b>58.96</b>	912.100	861.280	409.900	917.230	844.750	<u>342.38</u>

1055 after running for 12 hours, and hence we omit its discussion here. On *Lead-Lag*, OSLOM places  
1056 every node in a single cluster for most of the years, and clusters the rest of the years into either a huge  
1057 single cluster or two clusters.

1058 None of the methods outperform DIGRAC on our chosen performance measures from Table 1 ,  
1059 except on the *Lead-Lag* data set (See Tables 14, 15, 16 and 17 for the other results). With regards  
1060 to the 12 imbalance measures from Appendix Table 6, leaving out OSLOM as before, Louvain and  
1061 Leiden perform poorly on all of the real data sets, except on *Lead-Lag*. Indeed, for *Lead-Lag*, the  
1062 number of clusters we use for DIGRAC is ten according to the GICS sector memberships. However,  
1063 if we use the sector memberships as labels, the imbalance values are poor, which implies that ten may  
1064 not be a desirable choice of the number of clusters. Further, DIGRAC usually clusters the nodes into  
1065 smaller number of clusters, while Louvain and Leiden usually cluster the nodes into a larger number  
1066 of clusters (usually around 30, and sometimes above 50 clusters).

**Table 13:** Performance comparison on *Lead-Lag*, including Louvain and Leiden. Results in each year is averaged over ten runs. Mean and standard deviation (after ±) are calculated over the 19 years. The best is marked in **bold red** and the second best is marked in underline blue. InfoMap results are omitted here as it usually predicts a single huge cluster and could not generate imbalance results. Louvain and Leiden yield essentially identical results and often attain the highest objectives, while DIGRAC almost always places either first or second across all methods considered.

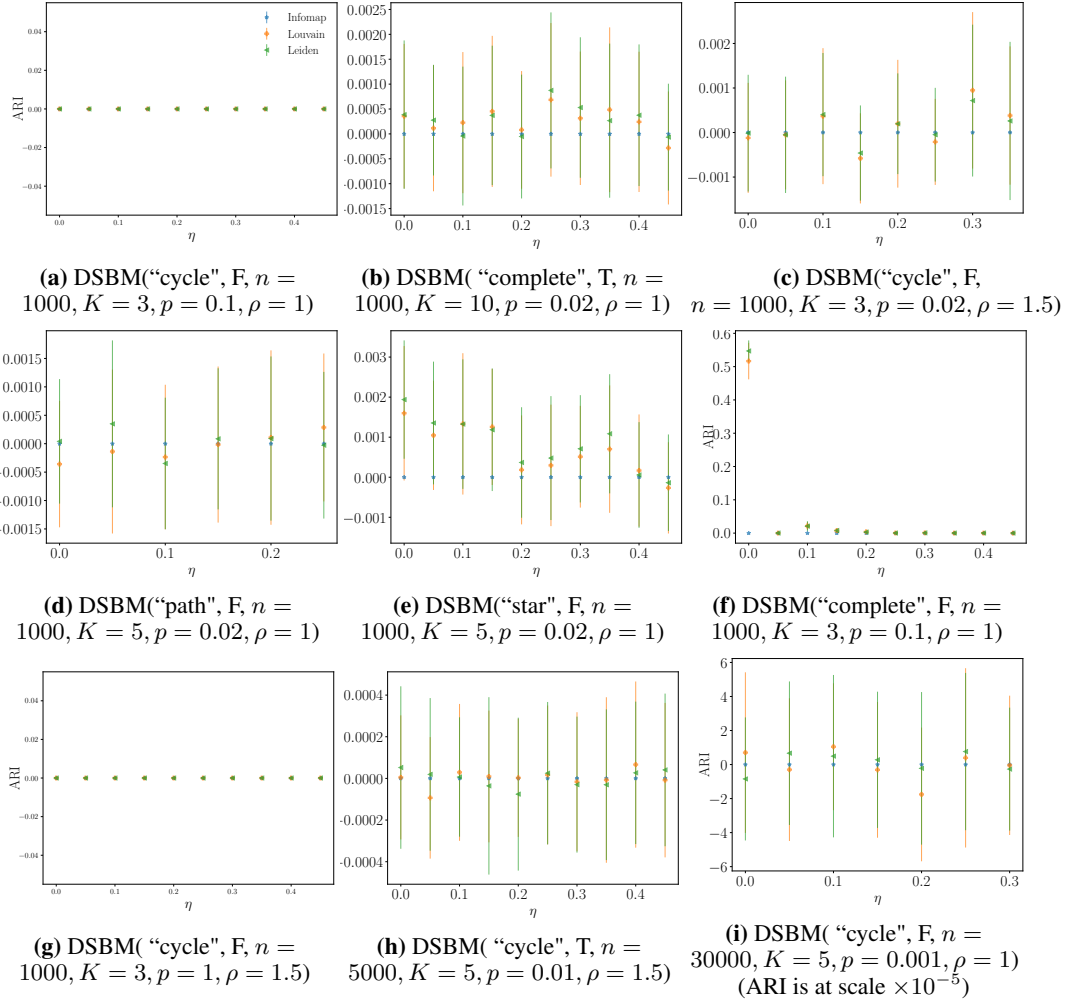
Metric/Method	Louvain/Leiden	Bi_sym	DD_sym	DISG_LR	Herm	Herm_rw	DIGRAC
$\mathcal{O}^{\text{sort}}_{\text{vol\_sum}}$	<u>0.08±0.02</u>	0.07±0.01	0.07±0.01	0.07±0.01	0.07±0.02	0.07±0.02	<b>0.15±0.03</b>
$\mathcal{O}^{\text{sort}}_{\text{vol\_min}}$	0.15±0.04	<b>0.51±0.10</b>	0.48±0.09	0.47±0.10	<b>0.51±0.11</b>	0.50±0.10	0.47±0.09
$\mathcal{O}^{\text{sort}}_{\text{vol\_max}}$	<u>0.08±0.02</u>	0.07±0.01	0.06±0.01	0.06±0.01	0.07±0.01	0.07±0.01	<b>0.14±0.03</b>
$\mathcal{O}^{\text{sort}}_{\text{plain}}$	0.15±0.04	<b>0.66±0.09</b>	0.64±0.08	0.63±0.08	<b>0.66±0.09</b>	0.65±0.09	0.53±0.09
$\mathcal{O}^{\text{std}}_{\text{vol\_sum}}$	<b>0.23±0.06</b>	0.04±0.01	0.04±0.01	0.04±0.01	0.04±0.01	0.04±0.01	<u>0.12±0.03</u>
$\mathcal{O}^{\text{std}}_{\text{vol\_min}}$	<b>0.46±0.11</b>	0.27±0.04	0.27±0.04	0.25±0.04	0.27±0.03	0.27±0.03	<u>0.38±0.07</u>
$\mathcal{O}^{\text{std}}_{\text{vol\_max}}$	<b>0.23±0.05</b>	0.04±0.00	0.03±0.00	0.03±0.00	0.03±0.00	0.03±0.00	<u>0.11±0.02</u>
$\mathcal{O}^{\text{std}}_{\text{plain}}$	<b>0.46±0.11</b>	0.40±0.05	0.39±0.05	0.38±0.05	0.40±0.05	0.40±0.05	<u>0.44±0.07</u>
$\mathcal{O}^{\text{naive}}_{\text{vol\_sum}}$	<b>0.23±0.06</b>	0.03±0.01	0.03±0.01	0.03±0.01	0.03±0.01	0.03±0.01	<u>0.08±0.04</u>
$\mathcal{O}^{\text{naive}}_{\text{vol\_min}}$	<b>0.46±0.11</b>	0.20±0.05	0.19±0.05	0.18±0.05	0.19±0.04	0.19±0.04	<u>0.26±0.10</u>
$\mathcal{O}^{\text{naive}}_{\text{vol\_max}}$	<b>0.23±0.05</b>	0.03±0.01	0.02±0.01	0.02±0.01	0.02±0.00	0.02±0.00	<u>0.08±0.03</u>
$\mathcal{O}^{\text{naive}}_{\text{plain}}$	<b>0.46±0.11</b>	0.30±0.06	0.28±0.06	0.27±0.06	0.29±0.05	0.29±0.05	<u>0.32±0.11</u>
size ratio	124.530	<u>3.67</u>	<b>3.34</b>	3.900	4.110	3.880	8.070
size std	47.960	<u>9.31</u>	<b>9.14</b>	10.090	10.490	10.360	17.060



**Figure 22:** US migration predicted clusters, along with the geographic locations of the counties as well as state boundaries (in black). The input digraph has extremely large entries; unlike in Fig. 20, we do not employ here the normalization given by Eq. (13). Altogether, this demonstrates the robustness of DIGRAC to outliers in the data, which is not a characteristic of other state-of-the-art methods such as Herm and Herm\_rw.

1067 Finally, we provide more examples/explanations on why these density-based methods or even other  
 1068 methods that are based on random-walk should fail. We would mainly like to point out a family of  
 1069 illustrative examples demonstrating the subtle nuance concerning edge density.

1070 Consider a meta graph with  $K = 3$  nodes (clusters) A,B,C with directed edges AB, BC, CA, hence  
 1071 a directed cycle (our "cycle" DSBM models). Each pair of nodes  $(v_i, v_j)$  in the graph of size  $n$  is  
 1072 connected by an edge independently with probability  $p$  (which can even be equal to 1, in the case of  
 1073 a complete graph), hence the graph has the same density throughout. Now suppose we consider a  
 1074 pair of nodes  $(v_i, v_j)$  such that  $v_i$  belongs to cluster A, and  $v_j$  to cluster B. Since this edge is part of



**Figure 23:** Test ARI comparison on synthetic data for Infomap, Louvain and Leiden. Error bars are given by one standard error.

1075 the metagraph, with probability  $1-\eta$ , it is directed from  $v_i$  to  $v_j$ , and with probability  $\eta$ ,  $v_j$  sends  
 1076 an edge to  $v_i$  (here,  $\eta$  is the noise level parameter). Similar arguments can be made when  $v_i$  (resp  
 1077  $v_j$ ) belongs to cluster B (resp C); and when  $v_i$  (resp  $v_j$ ) belongs to cluster C (resp A). See Figure 24  
 1078 for an illustration. We also see that when the network is complete (see Figure 23 (g) and Table 9),  
 1079 InfoMap [14] fails empirically as it produces a single huge cluster. As a method based on random  
 1080 walks, this failure might occur as the chain could hardly be trapped inside a cluster as in the usual  
 1081 setting.

1082 In such synthetic DSBM models with a "cycle" meta-graph structure, it can be shown that all nodes  
 1083 have the same in-degree and out-degree in expectation. Therefore, any density-based methods or  
 1084 modularity-based methods should fail. As the simplest possible example, one could just consider  
 1085  $K = 3$  clusters as above, without any noise (thus  $\eta = 0$ ). InfoMap [14] tries to minimize the  
 1086 description length, but as no description length difference occurs in the ground-truth clustering  
 1087 structure for such "cycle" DSBMs, if we consider a brute-force optimization of the map equation.  
 1088 Indeed, for any method that is based on a random walk, the probability of the random walker going  
 1089 from one cluster to another is the same as staying within the cluster. Therefore, we could hardly  
 1090 optimize anything if we base our clustering structure on a random walker's visit frequencies/path  
 1091 lengths. Similarly, the Markov clustering algorithm [65] is based on the intuition that higher-length  
 1092 paths would be relatively more likely to stay within clusters – an assumption that is not warranted  
 1093 when there is no density difference. [15] and [16] are two interesting Markov aggregation algorithms  
 1094 based on information theory and automatic control ideas that might be able to cover the above



**Table 14:** Performance comparison on *Telegram*, including Louvain and Leiden. The best is marked in **bold red** and the second best is marked in underline blue.

Metric/Method	Louvain/Leiden	InfoMap	Bi_sym	DD_sym	DISG_LR	Herm	Herm_rw	DIGRAC
$\mathcal{O}^{\text{sort}}_{\text{vol\_sum}}$	0.08±0.01	0.04±0.00	<u>0.21±0.00</u>	<u>0.21±0.00</u>	<u>0.21±0.01</u>	0.20±0.01	0.14±0.00	<b>0.32±0.01</b>
$\mathcal{O}^{\text{sort}}_{\text{vol\_min}}$	0.39±0.07	0.47±0.00	<u>0.67±0.00</u>	<u>0.61±0.00</u>	<u>0.66±0.02</u>	0.66±0.02	0.19±0.00	<b>0.79±0.06</b>
$\mathcal{O}^{\text{sort}}_{\text{vol\_max}}$	0.06±0.01	0.03±0.00	<u>0.20±0.00</u>	<u>0.20±0.00</u>	<u>0.20±0.01</u>	0.19±0.01	0.12±0.00	<b>0.29±0.01</b>
$\mathcal{O}^{\text{plain}}_{\text{vol\_sum}}$	0.71±0.05	<b>1.00±0.00</b>	0.80±0.00	0.75±0.00	0.78±0.03	0.76±0.04	0.59±0.00	<u>0.96±0.01</u>
$\mathcal{O}^{\text{std}}_{\text{vol\_sum}}$	0.07±0.01	0.01±0.00	0.26±0.00	0.26±0.00	0.26±0.01	0.25±0.02	<b>0.35±0.00</b>	<u>0.28±0.01</u>
$\mathcal{O}^{\text{std}}_{\text{vol\_min}}$	0.33±0.08	0.16±0.00	<b>0.84±0.00</b>	0.76±0.00	<u>0.82±0.03</u>	<u>0.82±0.03</u>	0.49±0.00	0.73±0.03
$\mathcal{O}^{\text{std}}_{\text{vol\_max}}$	0.05±0.01	0.01±0.00	<u>0.25±0.00</u>	<u>0.25±0.00</u>	<u>0.25±0.01</u>	0.24±0.02	<b>0.29±0.00</b>	<u>0.25±0.01</u>
$\mathcal{O}^{\text{plain}}_{\text{vol\_sum}}$	0.59±0.05	0.68±0.00	<b>1.00±0.00</b>	0.94±0.00	0.98±0.04	0.95±0.04	<u>0.99±0.00</u>	0.90±0.05
$\mathcal{O}^{\text{naive}}_{\text{vol\_sum}}$	0.06±0.02	0.01±0.00	<u>0.26±0.00</u>	<u>0.26±0.00</u>	<u>0.26±0.01</u>	0.25±0.02	0.23±0.00	<b>0.27±0.01</b>
$\mathcal{O}^{\text{naive}}_{\text{vol\_min}}$	0.28±0.11	0.11±0.00	<b>0.84±0.00</b>	0.76±0.00	<u>0.82±0.03</u>	<u>0.82±0.03</u>	0.32±0.00	0.72±0.04
$\mathcal{O}^{\text{naive}}_{\text{vol\_max}}$	0.04±0.01	0.00±0.00	<b>0.25±0.00</b>	<b>0.25±0.00</b>	<b>0.25±0.01</b>	0.24±0.02	0.20±0.00	0.24±0.01
$\mathcal{O}^{\text{plain}}_{\text{vol\_sum}}$	0.56±0.01	0.63±0.00	<b>1.00±0.00</b>	0.94±0.00	0.98±0.04	0.95±0.04	<u>0.99±0.00</u>	0.89±0.06

**Table 15:** Performance comparison on *Blog*, including Louvain and Leiden. The best is marked in **bold red** and the second best is marked in underline blue.

Metric/Method	Louvain/Leiden	InfoMap	Bi_sym	DD_sym	DISG_LR	Herm	Herm_rw	DIGRAC
$\mathcal{O}^{\text{sort}}_{\text{vol\_sum}}$	0.00±0.00	0.07±0.00	0.07±0.00	0.00±0.00	0.05±0.00	<u>0.37±0.00</u>	0.00±0.00	<b>0.44±0.00</b>
$\mathcal{O}^{\text{sort}}_{\text{vol\_min}}$	0.01±0.01	0.02±0.00	0.33±0.00	0.05±0.00	0.31±0.00	<u>0.78±0.01</u>	<b>0.89±0.00</b>	0.76±0.00
$\mathcal{O}^{\text{sort}}_{\text{vol\_max}}$	0.01±0.01	0.05±0.00	0.05±0.00	0.00±0.00	0.04±0.00	<u>0.26±0.00</u>	0.00±0.00	<b>0.40±0.00</b>
$\mathcal{O}^{\text{plain}}_{\text{vol\_sum}}$	<b>1.00±0.00</b>	<b>1.00±0.00</b>	0.33±0.00	0.05±0.00	0.31±0.00	0.78±0.01	0.89±0.00	0.76±0.00
$\mathcal{O}^{\text{std}}_{\text{vol\_sum}}$	0.00±0.00	0.00±0.00	0.07±0.00	0.00±0.00	0.05±0.00	<u>0.37±0.00</u>	0.00±0.00	<b>0.44±0.00</b>
$\mathcal{O}^{\text{std}}_{\text{vol\_min}}$	0.00±0.00	0.00±0.00	0.33±0.00	0.05±0.00	0.31±0.00	<u>0.78±0.01</u>	<b>0.89±0.00</b>	0.76±0.00
$\mathcal{O}^{\text{std}}_{\text{vol\_max}}$	0.00±0.00	0.00±0.00	0.05±0.00	0.00±0.00	0.04±0.00	<u>0.26±0.00</u>	0.00±0.00	<b>0.40±0.00</b>
$\mathcal{O}^{\text{plain}}_{\text{vol\_sum}}$	0.56±0.13	0.73±0.00	0.33±0.00	0.05±0.00	0.31±0.00	<u>0.78±0.01</u>	<b>0.89±0.00</b>	0.76±0.00
$\mathcal{O}^{\text{naive}}_{\text{vol\_sum}}$	0.00±0.00	0.00±0.00	0.07±0.00	0.00±0.00	0.05±0.00	<u>0.37±0.00</u>	0.00±0.00	<b>0.44±0.00</b>
$\mathcal{O}^{\text{naive}}_{\text{vol\_min}}$	0.00±0.00	0.00±0.00	0.33±0.00	0.05±0.00	0.31±0.00	<u>0.78±0.01</u>	<b>0.89±0.00</b>	0.76±0.00
$\mathcal{O}^{\text{naive}}_{\text{vol\_max}}$	0.00±0.00	0.00±0.00	0.05±0.00	0.00±0.00	0.04±0.00	<u>0.26±0.00</u>	0.00±0.00	<b>0.40±0.00</b>
$\mathcal{O}^{\text{plain}}_{\text{vol\_sum}}$	0.76±0.00	0.76±0.00	0.33±0.00	0.05±0.00	0.31±0.00	<u>0.78±0.01</u>	<b>0.89±0.00</b>	0.76±0.00

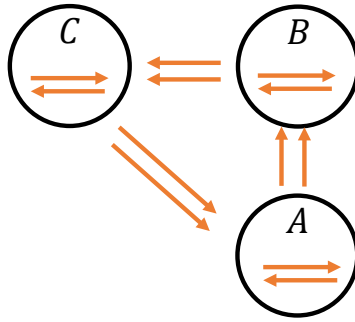
1095 example and may inspire some further comparison, but we omit comparison to them for now as we  
1096 already have more than ten comparison methods and that InfoMap shares similar ideas to these two  
1097 papers. As another example, as shown in [19], using belief propagation, in our model community  
1098 structure should not be detectable (the right-hand side of (20) in [19] is zero for our "cycle" DSBMs).  
1099 Therefore, at least methods that rely on belief propagation will fail on our benchmark models.

**Table 16:** Performance comparison on *Migration*, including Louvain and Leiden. The best is marked in **bold red** and the second best is marked in underline blue. InfoMap results are omitted here as it predicts a single huge cluster and could not generate imbalance results.

Metric/Method	Louvain/Leiden	Bi_sym	DD_sym	DISG_LR	Herm	Herm_rw	DIGRAC
$\mathcal{O}^{\text{sort}}_{\text{vol\_sum}}$	0.01±0.00	0.03±0.00	0.01±0.00	0.02±0.00	<u>0.04±0.00</u>	0.02±0.00	<b>0.05±0.00</b>
$\mathcal{O}^{\text{sort}}_{\text{vol\_min}}$	0.05±0.01	<b>0.19±0.00</b>	0.08±0.00	0.08±0.00	<u>0.15±0.02</u>	0.05±0.00	<u>0.18±0.03</u>
$\mathcal{O}^{\text{sort}}_{\text{vol\_max}}$	0.01±0.00	<u>0.03±0.00</u>	0.01±0.00	0.01±0.00	<u>0.03±0.00</u>	0.02±0.00	<b>0.04±0.00</b>
$\mathcal{O}^{\text{sort}}_{\text{plain}}$	0.09±0.02	0.24±0.00	0.20±0.00	0.17±0.00	<u>0.40±0.01</u>	<b>0.49±0.06</b>	0.29±0.04
$\mathcal{O}^{\text{std}}_{\text{vol\_sum}}$	0.00±0.00	0.01±0.00	0.01±0.00	0.01±0.00	<u>0.02±0.00</u>	<u>0.02±0.00</u>	<b>0.04±0.01</b>
$\mathcal{O}^{\text{std}}_{\text{vol\_min}}$	0.04±0.01	<u>0.10±0.00</u>	0.05±0.00	0.05±0.00	0.08±0.01	0.04±0.00	<b>0.16±0.03</b>
$\mathcal{O}^{\text{std}}_{\text{vol\_max}}$	0.00±0.00	<u>0.01±0.00</u>	<u>0.01±0.00</u>	<u>0.01±0.00</u>	<u>0.01±0.00</u>	<u>0.01±0.00</u>	<b>0.03±0.01</b>
$\mathcal{O}^{\text{std}}_{\text{plain}}$	0.07±0.01	0.13±0.00	0.12±0.00	0.11±0.00	<u>0.20±0.01</u>	<u>0.20±0.01</u>	<b>0.26±0.01</b>
$\mathcal{O}^{\text{naive}}_{\text{vol\_sum}}$	0.00±0.00	0.01±0.00	0.01±0.00	0.01±0.00	<u>0.02±0.00</u>	0.01±0.00	<b>0.04±0.01</b>
$\mathcal{O}^{\text{naive}}_{\text{vol\_min}}$	0.04±0.01	<u>0.09±0.00</u>	0.04±0.00	0.04±0.00	0.08±0.01	0.01±0.00	<b>0.16±0.03</b>
$\mathcal{O}^{\text{naive}}_{\text{vol\_max}}$	0.00±0.00	<u>0.01±0.00</u>	0.00±0.00	<u>0.01±0.00</u>	<u>0.01±0.00</u>	0.00±0.00	<b>0.03±0.01</b>
$\mathcal{O}^{\text{naive}}_{\text{plain}}$	0.07±0.00	0.12±0.00	0.10±0.00	0.08±0.00	<u>0.19±0.00</u>	<u>0.19±0.03</u>	<b>0.26±0.01</b>

**Table 17:** Performance comparison on *WikiTalk*, including Louvain and Leiden. The best is marked in **bold red** and the second best is marked in underline blue. InfoMap results are omitted here as its large number of predicted clusters leads to memory error in imbalance calculation.

Metric/Method	Louvain/Leiden	DISG_LR	Herm	Herm_rw	DIGRAC
$\mathcal{O}^{\text{sort}}_{\text{vol\_sum}}$	0.01±0.00	<u>0.18±0.03</u>	0.15±0.02	0.00±0.00	<b>0.24±0.05</b>
$\mathcal{O}^{\text{sort}}_{\text{vol\_min}}$	0.15±0.00	<u>0.10±0.03</u>	0.22±0.05	<u>0.26±0.00</u>	<b>0.28±0.13</b>
$\mathcal{O}^{\text{sort}}_{\text{vol\_max}}$	0.01±0.00	<u>0.16±0.03</u>	0.09±0.01	0.00±0.00	<b>0.19±0.04</b>
$\mathcal{O}^{\text{sort}}_{\text{plain}}$	<b>1.00±0.00</b>	0.87±0.08	0.99±0.01	0.98±0.00	<b>1.00±0.00</b>
$\mathcal{O}^{\text{std}}_{\text{vol\_sum}}$	0.00±0.00	<b>0.17±0.04</b>	0.06±0.01	0.01±0.00	<u>0.14±0.02</u>
$\mathcal{O}^{\text{std}}_{\text{vol\_min}}$	0.01±0.00	0.09±0.02	0.09±0.02	<b>0.27±0.00</b>	<u>0.18±0.08</u>
$\mathcal{O}^{\text{std}}_{\text{vol\_max}}$	0.00±0.00	<b>0.15±0.04</b>	0.04±0.00	0.00±0.00	<u>0.11±0.02</u>
$\mathcal{O}^{\text{std}}_{\text{plain}}$	0.42±0.00	0.72±0.03	0.70±0.05	<b>0.98±0.00</b>	<u>0.84±0.06</u>
$\mathcal{O}^{\text{naive}}_{\text{vol\_sum}}$	0.00±0.00	<u>0.10±0.02</u>	0.04±0.00	0.00±0.00	<b>0.12±0.01</b>
$\mathcal{O}^{\text{naive}}_{\text{vol\_min}}$	0.01±0.00	0.06±0.03	0.07±0.02	<b>0.26±0.00</b>	<u>0.15±0.07</u>
$\mathcal{O}^{\text{naive}}_{\text{vol\_max}}$	0.00±0.00	<b>0.09±0.02</b>	0.03±0.00	0.00±0.00	<b>0.09±0.01</b>
$\mathcal{O}^{\text{naive}}_{\text{plain}}$	0.43±0.00	0.64±0.04	0.61±0.04	<b>0.98±0.00</b>	<u>0.76±0.06</u>



**Figure 24:** An example of a "cycle" meta-graph.