

# HENCLER: NODE CLUSTERING IN HETEROPHILOUS GRAPHS VIA LEARNED ASYMMETRIC SIMILARITY

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Clustering nodes in heterophilous graphs is challenging as traditional methods assume that effective clustering is characterized by high intra-cluster and low inter-cluster connectivity. To address this, we introduce HeNCler—a novel approach for **H**eterophilous **N**ode **C**lustering. HeNCler *learns* a similarity graph by optimizing a clustering-specific objective based on weighted kernel singular value decomposition. Our approach enables spectral clustering on an *asymmetric* similarity graph, providing flexibility for both directed and undirected graphs. By solving the primal problem directly, our method overcomes the computational difficulties of traditional adjacency partitioning-based approaches. Experimental results show that HeNCler significantly improves node clustering performance in heterophilous graph settings, highlighting the advantage of its asymmetric graph-learning framework.

## 1 INTRODUCTION

Graph neural networks (GNNs) have substantially advanced machine learning applications to graph-structured data by effectively propagating node attributes end-to-end. Typically, GNNs rely on the assumption of homophily, where nodes with similar labels are more likely to be connected (Zheng et al., 2024; Wu et al., 2021). The homophily assumption holds true in contexts such as social networks and citation graphs, where models like GCN (Kipf & Welling, 2017), GIN (Xu et al., 2019), and GraphSAGE (Hamilton et al., 2017) excel at tasks like node classification and graph prediction.

However, in heterophilous datasets, such as web page and transaction networks, edges often link nodes with differing labels. Models like GAT (Veličković et al., 2018) and various graph transformers (Ying et al., 2022; Dwivedi & Bresson, 2021) have demonstrated improved performance on these datasets. Their attention mechanisms learning edge importance provide a straightforward way to reduce the reliance on homophily for supervised tasks.

Our work specifically addresses unsupervised attributed node clustering tasks. Such tasks necessitate entirely unsupervised or self-supervised learning approaches. For instance, auto-encoder type models (Park et al., 2019; Pan et al., 2020) are primarily focused on node representation learning rather than clustering, making them less suited for directly improving cluster-ability. Various self-supervised, contrastive learning techniques (Hassani & Ahmadi, 2020; You et al., 2020) enhance node representation learning in homophilous settings only and lack a specific clustering objective. At the same time, several self-supervised methods have been developed to handle heterophilous graphs (Chen et al., 2022; Xiao et al., 2022; Yuan et al., 2023). For example, MUSE (Yuan et al., 2023) extracts semantic and contextual views for contrastive learning. However, these methods are designed for the general node representation learning task and lack a clustering objective.

In contrast,  $S^3GC$  (Devvrit et al., 2022) employs a self-supervised approach specifically designed for clustering. It however assumes homophily by leveraging random walk co-occurrences to infer proximity-based similarities. MinCutPool (Bianchi et al., 2020) and DMoN (Tsitsulin et al., 2023) introduce unsupervised losses linked to graph structure, with strong theoretical ties to spectral clustering and graph modularity, respectively. These methods are suited for undirected graphs only, and moreover rely on partitioning the adjacency matrix where effective clustering correlates with high intra-cluster and low inter-cluster similarity—a premise often invalid in heterophilous graphs.

This paper introduces HeNCler, a novel approach for node clustering in heterophilous graphs, illustrated in Figure 1. Existing works overlook the asymmetric relationships in heterophilous

Table 1: Qualitative comparison of HeNCler with several baselines. In the table,  $|\mathcal{V}|$ ,  $|\mathcal{B}|$ , and  $|\mathcal{E}|$  denote the total number of nodes, the mini-batch size, and the number of edges respectively.

	BASELINES				OURS
	MINCUTP.	DMoN	S <sup>3</sup> GC	MUSE	HeNCler
CAN HANDLE HETEROPHILY	✗	✗	✗	✓	✓
DIRECTED GRAPHS	✗	✗	✓	✓	✓
HAS CLUSTERING OBJECTIVE	✓	✓	✓	✗	✓
SPACE COMPLEXITY	$\mathcal{O}( \mathcal{V} ^2)$	$\mathcal{O}( \mathcal{V}  +  \mathcal{E} )$	$\mathcal{O}( \mathcal{B} )$	$\mathcal{O}( \mathcal{V}  +  \mathcal{E} )$	$\mathcal{O}( \mathcal{B} )$
TIME COMPLEXITY	$\mathcal{O}( \mathcal{V}  +  \mathcal{E} )$	$\mathcal{O}( \mathcal{V}  +  \mathcal{E} )$	$\mathcal{O}( \mathcal{V} )$	$\mathcal{O}( \mathcal{V}  +  \mathcal{E} )$	$\mathcal{O}( \mathcal{V} )$

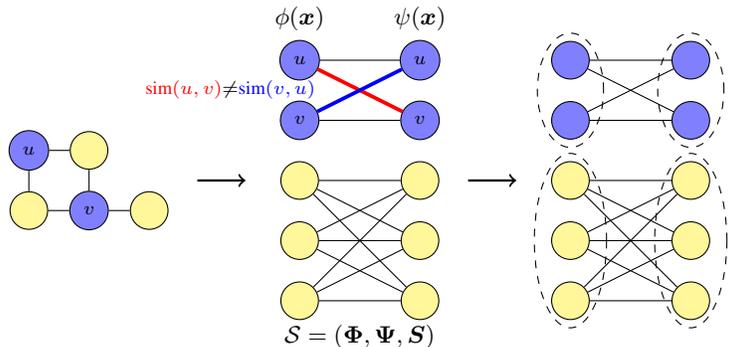


Figure 1: **HeNCler Overview.** Starting from a heterophilous graph, where nodes with the same label are not close to each other (left), HeNCler learns two sets of node representations,  $\{\phi(\mathbf{x}_v)\}_{v \in \mathcal{V}}$  and  $\{\psi(\mathbf{x}_v)\}_{v \in \mathcal{V}}$ , forming a bipartite graph  $\mathcal{S}$  (middle), where the similarity between nodes is defined as  $S_{uv} = \text{sim}(u, v) = \phi(\mathbf{x}_u)^\top \psi(\mathbf{x}_v)$ . Due to the clustering objective, nodes that should belong to the same cluster are positioned closer together in the learned graph. These clusters are then identified using spectral biclustering through wKSVD (right).

graphs, as shown in Table 1. HeNCler addresses this by using weighted kernel singular value decomposition (wKSVD) to induce a learned asymmetric similarity graph for both directed and undirected graphs. The dual problem of wKSVD aligns with asymmetric kernel spectral clustering, enabling the interpretation of similarities without homophily. By solving the primal problem directly, HeNCler overcomes computational difficulties and shows superior performance in node clustering tasks within heterophilous graphs.

**Contributions:** Our contributions in this work can be summarized as follows:

- We introduce HeNCler, a kernel spectral biclustering framework designed to *learn* an induced *asymmetric* similarity graph suited for node clustering of heterophilous graphs, applicable to both directed and undirected graphs.
- We develop a primal-dual framework for a generic weighted kernel singular value decomposition (wKSVD) model.
- We show that the dual wKSVD formulation allows for biclustering of bipartite/asymmetric graphs, while we employ a computationally feasible implementation in the primal wKSVD formulation.
- We further generalize our approach with trainable feature mappings, using node and edge decoders, such that the similarity matrix to cluster is learned.
- We train HeNCler in the primal setting and demonstrate its superior performance on the node clustering task for heterophilous attributed graphs. Our implementation is available in supplementary materials.

## 2 PRELIMINARIES AND RELATED WORK

We use lowercase symbols (e.g.,  $x$ ) for scalars, lowercase bold (e.g.,  $\mathbf{x}$ ) for vectors and uppercase bold (e.g.,  $\mathbf{X}$ ) for matrices. A single entry of a matrix is represented by  $X_{ij}$ .  $\phi(\cdot)$  denotes a mapping and  $\phi_v = \phi(\mathbf{x}_v)$  represents the mapping of node  $v$  in the induced feature space. We represent a graph  $\mathcal{G}$  by its vertices (i.e., nodes)  $\mathcal{V}$  and edges  $\mathcal{E}$ ,  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ , or by its node feature matrix and adjacency matrix  $\mathcal{G} = (\mathbf{X}, \mathbf{A})$ . For a bipartite graph, we have  $\mathcal{G} = (\mathcal{I}, \mathcal{J}, \mathcal{E})$  or  $\mathcal{G} = (\mathbf{X}_{\mathcal{I}}, \mathbf{X}_{\mathcal{J}}, \mathbf{S})$  where  $S_{ij}$  is the edge weight between nodes  $i \in \mathcal{I}$  and  $j \in \mathcal{J}$ . Note that  $\mathbf{S}$  is generally asymmetric and rectangular, and that the adjacency matrix of the bipartite graph is given by  $\mathbf{A} = \begin{bmatrix} \mathbf{0} & \mathbf{S} \\ \mathbf{S}^{\top} & \mathbf{0} \end{bmatrix}$ .

**Kernel singular value decomposition** KSVD (Suykens, 2016) sets up a primal-dual framework, based on Lagrange duality, that formulates a variational principle in the primal formulation that corresponds to the matrix singular value decomposition (SVD) in the dual for linear feature maps. By employing non-linear feature mappings or asymmetric kernel functions, this framework allows for non-linear extensions of the SVD problem. The KSVD framework can be applied on data structures such as row and column features, directed graphs, and/or can exploit asymmetric similarity information such as conditional probabilities (He et al., 2023). Interestingly, KSVD often outperforms the similar though symmetric kernel principal component analysis model on tasks where the asymmetry is not immediately apparent (Tao et al., 2024). A different connection is shown in Primal-Attention (Chen et al., 2023), where the authors demonstrate the relation between canonical self-attention, which is asymmetric, and KSVD. They show how to gain computational efficiency by considering a primal equivalent of the attention mechanism.

**Spectral clustering** generalizations have been proposed in many settings. Spectral graph biclustering (Dhillon, 2001) formulates the spectral clustering problem of a bipartite graph  $\mathcal{G} = (\mathcal{I}, \mathcal{J}, \mathbf{S})$  and shows the equivalence with the SVD of the normalized matrix  $\mathbf{S}_n = \mathbf{D}_1^{-1/2} \mathbf{S} \mathbf{D}_2^{-1/2}$ , where  $D_{1,ii} = \sum_j S_{ij}$  and  $D_{2,jj} = \sum_i S_{ij}$ . Cluster assignments for nodes  $\mathcal{I}$  and nodes  $\mathcal{J}$  can be inferred from the left and right singular vectors respectively. Further, kernel spectral clustering (KSC) (Alzate & Suykens, 2010) proposes a weighted kernel principal component analysis in which the dual formulation corresponds to the random walks interpretation of the spectral clustering problem. KSC and the aforementioned spectral biclustering formulation lack asymmetry and a primal formulation respectively, which are limitations that our model will address.

**Restricted kernel machines** (RKM) (Suykens, 2017) possess primal and dual model formulations, based on the concept of conjugate feature duality. It is an energy-based framework for (deep) kernel machines, that shows relations with least-squares support vector machines (Suykens et al., 2002) and restricted Boltzmann machines (Salakhutdinov, 2015). The RKM framework encompasses many model classes, including classification, regression, kernel principal component analysis and KSVD, and allows for deep kernel learning (Tonin et al., 2021) and deep kernel learning on graphs (Achten et al., 2024). One possibility to represent the feature maps in RKMs is by means of deep neural networks, e.g., for unsupervised representation learning (Pandey et al., 2021; 2022). RKM models can work in either primal or dual setting, and with decomposition or gradient based algorithms (Achten et al., 2023).

**Homophilous node clustering** methods like MinCutPool (Bianchi et al., 2020) and DMoN (Tsitulin et al., 2023) introduce unsupervised loss functions within a graph neural network framework. MinCutPool employs a relaxed version of the minimal cut loss applied to the adjacency matrix, while DMoN optimizes the modularity score of clustering labels with respect to the adjacency structure. Both of these methods rely on partitioning the adjacency matrix and inherently assume homophily. Additionally, due to their theoretical underpinnings, these losses are only applicable to undirected graphs. Beyond these adjacency partitioning-based approaches, self-supervised or contrastive methods have also been proposed (You et al., 2020; Hassani & Ahmadi, 2020; Devvrit et al., 2022). These methods typically use graph proximity as their supervision signal, which similarly assumes homophily. For example, S<sup>3</sup>GC (Devvrit et al., 2022) employs a self-supervised loss based on random walk co-occurrences.

**Heterophilous node clustering** methods typically rely on self-supervised or contrastive techniques. Gong et al. (Gong et al., 2023) propose Sparse Graph Anomaly Detection (SparseGAD), a method that sparsifies graph structures to effectively reduce noise from irrelevant edges and enhance the detection of closely related nodes. This technique reveals underlying node dependencies, accommodating both homophilous and heterophilous relationships. Similarly, HGRL (Chen et al., 2022) employs

self-supervised learning on heterophilous graphs by utilizing graph augmentation techniques to capture global and higher-order structural information. MUSE (Yuan et al., 2023), on the other hand, constructs semantic and contextual views to capture both node-level and neighborhood information for contrastive learning, subsequently integrating these multi-view representations through a fusion controller.

While adjacency partitioning-based methods have demonstrated both theoretical and empirical success for homophilous graphs, they have not been effectively extended to heterophilous graph learning. On the other hand, self-supervised clustering approaches, though promising, often lack a clear clustering interpretation. In the following section, we introduce HeNCler, which bridges these gaps.

### 3 METHOD

**Model motivation** Our approach employs an RKM auto-encoder framework, which has been shown to be effective in unsupervised representation learning by jointly optimizing feature mappings and projection matrices within a kernel-based setting (Pandey et al., 2022). To capture long-range relational dependencies in heterophilous graphs, we utilize a KSVD loss, where a double feature mapping yields a learned asymmetric similarity matrix. To further enhance the cluster-ability of this matrix, we extend the loss function to a weighted KSVD (wKSVD) loss, which not only boosts clustering performance but also offers a spectral graph biclustering interpretation. We next introduce a general wKSVD framework, after which we introduce our HeNCler model that operates in the primal setting while jointly learning the feature mappings end-to-end.

#### 3.1 KERNEL SPECTRAL BICLUSTERING WITH ASYMMETRIC SIMILARITIES

Consider a dataset with two, possibly different, input sources  $\{\mathbf{x}_i\}_{i=1}^n$  and  $\{\mathbf{z}_j\}_{j=1}^m$ , on which we want to define an unsupervised learning task. To this end, we introduce a weighted kernel singular value decomposition model (wKSVD), starting from the following primal optimization problem, which is a weighted variant of the KSVD formulation:

$$\begin{aligned} \min_{\mathbf{U}, \mathbf{V}, \mathbf{e}, \mathbf{r}} J \triangleq & \text{Tr}(\mathbf{U}^\top \mathbf{V}) - \frac{1}{2} \sum_{i=1}^n w_{1,i} \mathbf{e}_i^\top \boldsymbol{\Sigma}^{-1} \mathbf{e}_i - \frac{1}{2} \sum_{j=1}^m w_{2,j} \mathbf{r}_j^\top \boldsymbol{\Sigma}^{-1} \mathbf{r}_j \\ \text{s.t. } & \{\mathbf{e}_i = \mathbf{U}^\top \phi(\mathbf{x}_i), \forall i = 1, \dots, n; \quad \mathbf{r}_j = \mathbf{V}^\top \psi(\mathbf{z}_j), \forall j = 1, \dots, m\}, \end{aligned} \quad (1)$$

with projection matrices  $\mathbf{U}, \mathbf{V} \in \mathbb{R}^{d_f \times s}$ ; strictly positive weighting scalars  $w_{1,i}, w_{2,j}$ ; latent variables  $\mathbf{e}_i, \mathbf{r}_j \in \mathbb{R}^s$ ; diagonal and positive definite hyperparameter matrix  $\boldsymbol{\Sigma} \in \mathbb{R}^{s \times s}$ ; and centered feature maps  $\phi(\cdot) : \mathbb{R}^{d_x} \mapsto \mathbb{R}^{d_f}$  and  $\psi(\cdot) : \mathbb{R}^{d_z} \mapsto \mathbb{R}^{d_f}$ ; details on centering of the feature maps are provided in Appendix A. The following derivation shows the equivalence with the spectral biclustering problem.

**Proposition 1.** *The solution to the primal problem (1) can be obtained by solving the singular value decomposition of*

$$\mathbf{W}_1^{1/2} \mathbf{S} \mathbf{W}_2^{1/2} = \mathbf{H}_e \boldsymbol{\Sigma} \mathbf{H}_r^\top, \quad (2)$$

where  $\mathbf{W}_1$  and  $\mathbf{W}_2$  are diagonal matrices such that  $W_{1,ii} = w_{1,i}$  and  $W_{2,jj} = w_{2,j}$ ,  $\mathbf{S} = \boldsymbol{\Phi} \boldsymbol{\Psi}^\top$  is an asymmetric similarity matrix where  $S_{ij} = \phi(\mathbf{x}_i)^\top \psi(\mathbf{z}_j)$ ,  $\boldsymbol{\Phi} = [\phi(\mathbf{x}_1) \dots \phi(\mathbf{x}_n)]^\top$ ,  $\boldsymbol{\Psi} = [\psi(\mathbf{z}_1) \dots \psi(\mathbf{z}_m)]^\top$ , and where  $\mathbf{H}_e = [\mathbf{h}_{e_1} \dots \mathbf{h}_{e_n}]^\top$ , and  $\mathbf{H}_r = [\mathbf{h}_{r_1} \dots \mathbf{h}_{r_m}]^\top$  are the left and right singular vectors respectively; and by applying  $\mathbf{r}_j = \boldsymbol{\Sigma} \mathbf{h}_{r_j} / \sqrt{w_{2,j}}$  and  $\mathbf{e}_i = \boldsymbol{\Sigma} \mathbf{h}_{e_i} / \sqrt{w_{1,i}}$ .

*Proof.* We now introduce dual variables  $\mathbf{h}_{e_i}$  and  $\mathbf{h}_{r_j}$  using a case of Fenchel-Young inequality (Rockafellar, 1974):

$$\frac{1}{2} w_{1,i} \mathbf{e}_i^\top \boldsymbol{\Sigma}^{-1} \mathbf{e}_i + \frac{1}{2} \mathbf{h}_{e_i}^\top \boldsymbol{\Sigma} \mathbf{h}_{e_i} \geq \sqrt{w_{1,i}} \mathbf{e}_i^\top \mathbf{h}_{e_i}, \quad \frac{1}{2} w_{2,j} \mathbf{r}_j^\top \boldsymbol{\Sigma}^{-1} \mathbf{r}_j + \frac{1}{2} \mathbf{h}_{r_j}^\top \boldsymbol{\Sigma} \mathbf{h}_{r_j} \geq \sqrt{w_{2,j}} \mathbf{r}_j^\top \mathbf{h}_{r_j}, \quad (3)$$

$\forall \mathbf{e}_i, \mathbf{r}_j, \mathbf{h}_{e_i}, \mathbf{h}_{r_j} \in \mathbb{R}^s, \forall w_{1,i}, w_{2,j} \in \mathbb{R}_{>0}, \forall \boldsymbol{\Sigma} \in \mathbb{R}_{>0}^{s \times s}$ . The above inequalities can be verified by writing it in quadratic form:  $\frac{1}{2} \begin{bmatrix} \mathbf{e}_i^\top & \mathbf{h}_{e_i}^\top \end{bmatrix} \begin{bmatrix} w_{1,i} \boldsymbol{\Sigma}^{-1} & -\sqrt{w_{1,i}} \mathbf{I}_s \\ -\sqrt{w_{1,i}} \mathbf{I}_s & \boldsymbol{\Sigma} \end{bmatrix} \begin{bmatrix} \mathbf{e}_i \\ \mathbf{h}_{e_i} \end{bmatrix} \geq 0, \quad \forall i$ , with  $\mathbf{I}_s$  the  $s$ -dimensional identity matrix, which follows immediately from the Schur complement form:

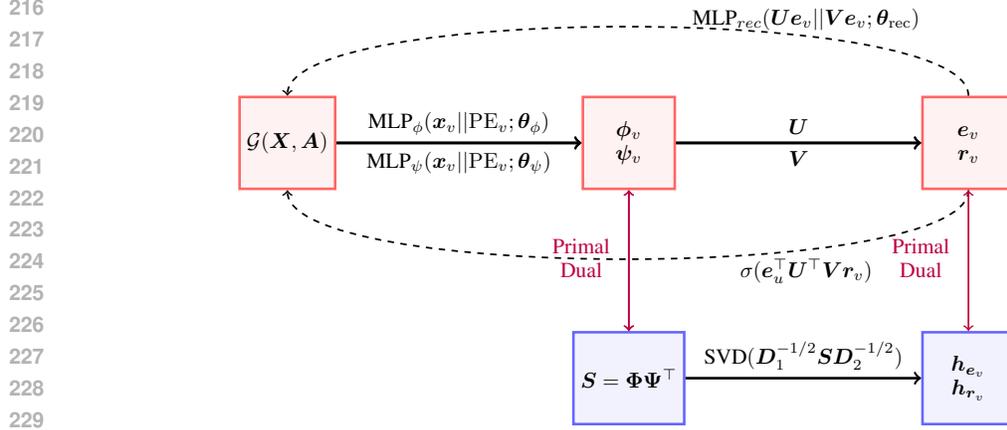


Figure 2: **The HeNCler model.** HeNCler operates in the primal setting (top of the figure in red) and uses a double multilayer perceptron (MLP) to map node representations to a feature space. The obtained representations  $\phi_v$  and  $\psi_v$  are then projected to latent representations  $e_v$  and  $r_v$  respectively. The wKSVD loss ensures that these latent representations correspond to the dual equivalent (bottom of the figure in blue) i.e., a biclustering of the asymmetric similarity graph defined by  $S$ . The node and edge reconstructions (dashed arrows) aid in the feature map learning.

for a matrix  $Q = \begin{bmatrix} Q_1 & Q_2 \\ Q_2^\top & Q_3 \end{bmatrix}$ , one has  $Q \succeq 0$  if and only if  $Q_1 \succ 0$  and the Schur complement  $Q_3 - Q_2^\top Q_1^{-1} Q_2 \succeq 0$  (Boyd & Vandenberghe, 2004).

By substituting the constraints of (1) and inequalities (3) into the objective function of (1), we obtain an objective in primal and dual variables as an upper bound on the primal objective  $\bar{J} \geq J$ :

$$\min_{U, V, h_e, h_r} \bar{J} \triangleq \text{Tr}(U^\top V) - \sum_{i=1}^n \sqrt{w_{1,i}} \phi(x_i)^\top U h_{e_i} + \frac{1}{2} \sum_{i=1}^n h_{e_i}^\top \Sigma h_{e_i} - \sum_{j=1}^m \sqrt{w_{2,j}} \psi(z_j)^\top V h_{r_j} + \frac{1}{2} \sum_{j=1}^m h_{r_j}^\top \Sigma h_{r_j}. \quad (4)$$

Next, we formulate the stationarity conditions of problem (4):

$$\begin{aligned} \frac{\partial \bar{J}}{\partial V} = 0 &\Rightarrow U = \sum_{j=1}^m \sqrt{w_{2,j}} \psi(z_j) h_{r_j}^\top, & \frac{\partial \bar{J}}{\partial h_{e_i}} = 0 &\Rightarrow \Sigma h_{e_i} = \sqrt{w_{1,i}} U^\top \phi(x_i), \\ \frac{\partial \bar{J}}{\partial U} = 0 &\Rightarrow V = \sum_{i=1}^n \sqrt{w_{1,i}} \phi(x_i) h_{e_i}^\top, & \frac{\partial \bar{J}}{\partial h_{r_j}} = 0 &\Rightarrow \Sigma h_{r_j} = \sqrt{w_{2,j}} V^\top \psi(z_j), \end{aligned} \quad (5)$$

from which we then eliminate the primal variables  $U$  and  $V$ . This yields the eigenvalue problem:

$$\begin{bmatrix} \mathbf{0} & W_1^{1/2} S W_2^{1/2} \\ W_2^{1/2} S^\top W_1^{1/2} & \mathbf{0} \end{bmatrix} \begin{bmatrix} H_e \\ H_r \end{bmatrix} = \begin{bmatrix} H_e \\ H_r \end{bmatrix} \Sigma, \quad (6)$$

where  $\mathbf{0}$  is an all-zeros matrix. Note that, by Lanczos' Theorem (Lanczos, 1958), the above eigenvalue problem is equivalent with (2), and that the stationarity conditions (5) provide the relationships between primal and dual variables, which concludes the proof.  $\square$

We have thus shown the connection between the primal (1) and dual formulation (2). Similarly to the KSVD framework, the wKSVD framework can be used for learning with asymmetric kernel functions and/or rectangular data sources. The spectral biclustering problem can now easily be obtained by choosing the weights  $w_{1,i}$  and  $w_{2,j}$  appropriately.

**Corollary 2.** Given Proposition 1, and by choosing  $\mathbf{W}_1$  and  $\mathbf{W}_2$  to equal  $\mathbf{D}_1^{-1/2}$  and  $\mathbf{D}_2^{-1/2}$ , where  $D_{1,ii} = \sum_j S_{ij}$  and  $D_{2,jj} = \sum_i S_{ij}$ , we obtain the random walk interpretation  $\mathbf{D}_1^{-1/2} \mathbf{S} \mathbf{D}_2^{-1/2} = \mathbf{H}_c \mathbf{\Sigma} \mathbf{H}_r^\top$  of the spectral graph bipartitioning problem for the bipartite graph  $\mathcal{S} = (\Phi, \Psi, \mathcal{S})$ .

Moreover, the wKSVD framework is more general as, on the one hand, one can use a given similarity matrix (e.g. adjacency matrix of a graph) or (asymmetric) kernel function in the dual, or, on the other hand, one can choose to use explicitly defined (deep) feature maps in both primal or dual.

### 3.2 THE HeNCLER MODEL

HeNCLer employs the wKSVD framework in a graph setting, where the dataset is a node set  $\mathcal{V}$  and where the asymmetry arises from employing to different mappings that operate on the nodes given the entire graph  $\mathcal{G} = (\mathcal{X}, \mathcal{A})$ . Our method is visualized in Figure 2, where red indicates the primal setting of the framework and blue the dual.

In the preceding subsection, we showed that problem (1) has an equivalent dual problem corresponding to the graph bipartitioning problem, when  $w_{1,i}$  and  $w_{2,j}$  are chosen to equal the square root of the inverse of the out-degree and in-degree of a similarity graph  $\mathcal{S}$  respectively. This similarity graph  $\mathcal{S}$  depends on the feature mappings  $\phi(\cdot)$  and  $\psi(\cdot)$ , which for our method does not only depend on the node of interest, but also on the rest of the input graph and the learnable parameters. The mappings for node  $v$  thus become  $\phi(\mathbf{x}_v, \mathcal{G}; \theta_\phi)$  and  $\psi(\mathbf{x}_v, \mathcal{G}; \theta_\psi)$  and we will ease these notations to  $\phi(\mathbf{x}_v)$  and  $\psi(\mathbf{x}_v)$ . The ability of our method to learn these feature mappings is an important aspect of our contribution, as a key motivation behind our model is that we need to learn new similarities for clustering heterophilous graphs. The loss function is comprised of three terms: the wKSVD-loss, a node-reconstruction loss, and an edge-reconstruction loss:

$$\mathcal{L}_{\text{wKSVD}}(\mathbf{U}, \mathbf{V}, \mathbf{\Sigma}, \theta_\phi, \theta_\psi) + \mathcal{L}_{\text{NodeRec}}(\mathbf{U}, \mathbf{V}, \theta_\phi, \theta_\psi, \theta_{\text{rec}}) + \mathcal{L}_{\text{EdgeRec}}(\mathbf{U}, \mathbf{V}, \theta_\phi, \theta_\psi),$$

where the trainable parameters of the model are in the the multilayer perceptron (MLP) feature maps ( $\theta_\phi$  and  $\theta_\psi$ ), the MLP node decoder ( $\theta_{\text{rec}}$ ), in the  $\mathbf{U}$  and  $\mathbf{V}$  projection matrices, and in the singular values  $\mathbf{\Sigma}$ . All these parameters are trained end-to-end and we next explain the losses in more detail.

**wKSVD-loss** Instead of solving the SVD in the dual formulation, HeNCLer leverages the primal formulation (1) of the wKSVD framework for greater computational efficiency. While equation (1) assumes that the feature maps  $\phi(\cdot)$  and  $\psi(\cdot)$  are fixed, HeNCLer utilizes parametric functions  $\phi(\cdot; \theta_\phi)$  and  $\psi(\cdot; \theta_\psi)$ , enabling it to learn new similarities between nodes. By incorporating regularization terms for these functions and defining the weighting scalars as  $w_{1,v} = D_{1,vv}^{-1} = 1/\sum_u \phi(\mathbf{x}_v)^\top \psi(\mathbf{x}_u)$  and  $w_{2,v} = D_{2,vv}^{-1} = 1/\sum_u \phi(\mathbf{x}_u)^\top \psi(\mathbf{x}_v)$ , we derive the wKSVD-loss:

$$\begin{aligned} \mathcal{L}_{\text{wKSVD}} \triangleq & - \sum_{v=1}^{|\mathcal{V}|} D_{1,vv}^{-1} \phi(\mathbf{x}_v)^\top \mathbf{U} \mathbf{\Sigma}^{-1} \mathbf{U}^\top \phi(\mathbf{x}_v) - \sum_{v=1}^{|\mathcal{V}|} D_{2,vv}^{-1} \psi(\mathbf{x}_v)^\top \mathbf{V} \mathbf{\Sigma}^{-1} \mathbf{V}^\top \psi(\mathbf{x}_v) \\ & + \text{Tr}(\mathbf{U}^\top \mathbf{V}) + \sum_{v=1}^{|\mathcal{V}|} \sqrt{D_{1,vv}^{-1} D_{2,vv}^{-1}} \phi(\mathbf{x}_v)^\top \psi(\mathbf{x}_v). \quad (7) \end{aligned}$$

The primal formulation of HeNCLer (7) can be understood as follows: The first two terms aim to maximize the weighted variance of the learned node representations  $\mathbf{e}$  and  $\mathbf{r}$ . The third and fourth terms act as regularizers, encouraging asymmetry by penalizing the similarity between  $\mathbf{U}$  and  $\mathbf{V}$ , and between  $\phi(\mathbf{x}_v)$  and  $\psi(\mathbf{x}_v)$ , respectively.

For the two feature maps  $\phi(\cdot)$  and  $\psi(\cdot)$ , we employ two MLPs:  $\phi(\mathbf{x}_v, \mathcal{G}; \theta_\phi) \equiv \text{MLP}_\phi(\mathbf{x}_v \| \text{PE}_v; \theta_\phi)$  and  $\psi(\mathbf{x}_v, \mathcal{G}; \theta_\psi) \equiv \text{MLP}_\psi(\mathbf{x}_v \| \text{PE}_v; \theta_\psi)$ . We construct a random walks positional encoding (PE) (Dwivedi et al., 2022) to embed the network’s structure and concatenate this encoding with the node attributes. The MLPs have two linear layers with a LeakyReLU activation function in between, followed by a batch normalization layer. The singular values in  $\mathbf{\Sigma}$  are jointly learned, constrained to lie between 0 and 1, with the additional condition that  $\text{Tr}(\mathbf{\Sigma}^{-\frac{1}{2}}) = 1$ .

**Reconstruction losses** Since the feature maps  $\phi(\cdot)$  and  $\psi(\cdot)$  need to be learned, an additional loss function beyond the above regularization term is required to effectively optimize the parameters of the MLPs. As the node clustering setting is completely unsupervised, we add a decoder network and

a reconstruction loss. This technique has been proven to be effective for unsupervised learning in the RKM-framework (Pandey et al., 2022), as well as for unsupervised node representation learning (Sun et al., 2021). For heterophilous graphs, we argue that it is particularly important to also reconstruct node features and not only the graph structure.

For the node reconstruction, we first project the  $e$  and  $r$  variables back to feature space, concatenate these and then map to input space with another MLP. This MLP has also two layers and a leaky ReLU activation function. The hidden layer size is set to the average of the latent dimension and input dimension. With the mean-squared-error as the associated loss, this gives:

$$\mathcal{L}_{\text{NodeRec}} = \frac{1}{|\mathcal{V}|} \sum_{v \in \mathcal{V}} \|\text{MLP}_{\text{rec}}(\mathbf{U}e_v \| \mathbf{V}r_v; \theta_{\text{rec}}) - \mathbf{x}_v\|^2. \quad (8)$$

To reconstruct edges, we use a simple dot-product decoder  $\sigma(e_u^\top \mathbf{U}^\top \mathbf{V}r_v)$  where  $\sigma$  is the sigmoid function. By using the  $e$  representation for source nodes and  $r$  for target nodes, this reconstruction is asymmetric and can reconstruct directed graphs. We use a binary cross-entropy loss:

$$\mathcal{L}_{\text{EdgeRec}} = \frac{1}{|\mathcal{U}|} \sum_{(u,v) \in \mathcal{U}} \text{BCE}(\sigma(e_u^\top \mathbf{U}^\top \mathbf{V}r_v), \mathcal{E}_{uv}), \quad (9)$$

where  $\mathcal{U}$  is a node-tuple set, resampled every epoch, containing  $2|\mathcal{V}|$  positive edges from  $\mathcal{E}$  and  $2|\mathcal{V}|$  negative edges from  $\mathcal{E}^C$ , and  $\mathcal{E}_{uv} \in \{0, 1\}$  indicates whether an edge  $(u, v)$  exist:  $(u, v) \in \mathcal{E}$ .

**Optimizer, constraints, and cluster assignment** We use Adam (Kingma & Ba, 2015) for the training of all parameters. The batch normalization in the MLP’s keeps the wKSVD-loss bounded and the constraints on the singular values is enforced with a softmax function. Cluster assignments are obtained by KMeans clustering on the concatenation of learned  $e$  and  $r$  node representations.

HeNCler jointly learns the wKSVD projection matrices,  $\mathbf{U}$  and  $\mathbf{V}$ , along with the feature map parameters,  $\theta_\phi$  and  $\theta_\psi$ . The wKSVD loss improves the cluster-ability of the learned similarity graph, ensuring that  $e$  and  $r$  function as spectral biclustering embeddings. The two distinct feature maps enable asymmetric learning, effectively capturing potential asymmetric relationships in the data, while the reconstruction losses ensure robust and meaningful representation learning.

Table 2: Dataset statistics of the employed heterophilous graphs.

Dataset	short	# Nodes	# Edges	# Classes	Directed	$\mathcal{H}(\mathcal{G})$
Texas	tex	183	325	5	✓	0.000
Cornell	corn	183	298	5	✓	0.150
Wisconsin	wis	251	515	5	✓	0.084
Chameleon	cha	2,277	31,371	5	✗	0.042
Squirrel	squi	5,201	198,353	5	✗	0.031
Roman-empire	rom	22,662	32,927	18	✗	0.021
Minesweeper	mine	10,000	39,402	2	✗	0.009
Tolokers	tol	11,758	519,000	2	✗	0.180

## 4 EXPERIMENTS

**Datasets** We assess the performance of HeNCler on heterophilous attributed graphs that are available in literature. we use Texas, Cornell, and Wisconsin (Pei et al., 2020)<sup>1</sup>, which are directed webpage networks where edges encode hyperlinks between pages. Next, we use Chameleon and Squirrel (Rozemberczki et al., 2021), which are undirected Wikipedia webpage networks where edges encode mutual links. We further assess our model on the undirected graphs: Roman-empire, Minesweeper, and Tolokers (Platonov et al., 2023), which are a graph representation of a Wikipedia article, a grid graph based on the minesweeper game, and a crowd-sourcing network respectively. We include

<sup>1</sup><http://www.cs.cmu.edu/afs/cs.cmu.edu/project/theo-11/www/wwkb>

Table 3: Experimental results on heterophilous graphs. We report NMI and F1 scores for 10 runs (mean  $\pm$  standard deviation), where higher values indicate better performance. The best results for each metric are highlighted in bold.

Dataset		Baselines				Ours	
		KMeans	MinCutPool	DMoN	S <sup>3</sup> GC	MUSE	HeNCler
tex	NMI	4.97 $\pm$ 1.00	11.60 $\pm$ 2.19	9.06 $\pm$ 2.11	11.56 $\pm$ 1.46	39.23 $\pm$ 4.91	<b>43.65</b> $\pm$ 2.52
	F1	59.27 $\pm$ 0.83	55.26 $\pm$ 0.56	47.76 $\pm$ 4.79	43.69 $\pm$ 2.74	65.96 $\pm$ 3.52	<b>71.39</b> $\pm$ 2.16
corn	NMI	5.42 $\pm$ 2.04	17.04 $\pm$ 1.61	12.49 $\pm$ 2.51	14.48 $\pm$ 1.79	38.99 $\pm$ 2.73	<b>41.52</b> $\pm$ 4.35
	F1	52.97 $\pm$ 0.24	51.21 $\pm$ 5.06	43.83 $\pm$ 6.23	33.13 $\pm$ 0.83	60.58 $\pm$ 3.61	<b>63.40</b> $\pm$ 3.67
wis	NMI	6.84 $\pm$ 4.39	13.38 $\pm$ 2.36	12.56 $\pm$ 1.23	13.07 $\pm$ 0.61	39.71 $\pm$ 2.22	<b>47.13</b> $\pm$ 1.76
	F1	56.16 $\pm$ 0.58	55.63 $\pm$ 2.96	45.72 $\pm$ 7.85	31.71 $\pm$ 2.25	58.94 $\pm$ 3.09	<b>68.30</b> $\pm$ 2.17
cha	NMI	0.44 $\pm$ 0.11	11.88 $\pm$ 1.99	12.87 $\pm$ 1.86	15.83 $\pm$ 0.26	23.06 $\pm$ 0.28	<b>23.89</b> $\pm$ 0.84
	F1	<b>53.23</b> $\pm$ 0.07	50.40 $\pm$ 5.65	45.05 $\pm$ 4.30	36.51 $\pm$ 0.24	52.10 $\pm$ 0.48	44.14 $\pm$ 1.83
squi	NMI	1.40 $\pm$ 2.12	6.35 $\pm$ 0.32	3.08 $\pm$ 0.38	3.83 $\pm$ 0.11	8.30 $\pm$ 0.23	<b>9.67</b> $\pm$ 0.13
	F1	54.05 $\pm$ 2.72	<b>55.26</b> $\pm$ 0.57	49.21 $\pm$ 2.74	35.08 $\pm$ 0.18	50.07 $\pm$ 5.99	36.51 $\pm$ 2.39
rom	NMI	35.20 $\pm$ 1.79	9.97 $\pm$ 2.02	13.14 $\pm$ 0.53	14.48 $\pm$ 0.21	<b>40.50</b> $\pm$ 0.73	36.99 $\pm$ 0.61
	F1	37.17 $\pm$ 2.12	<b>42.19</b> $\pm$ 0.26	22.69 $\pm$ 3.91	17.76 $\pm$ 0.53	38.34 $\pm$ 0.35	35.43 $\pm$ 1.07
mine	NMI	0.02 $\pm$ 0.02	6.16 $\pm$ 2.17	<b>6.87</b> $\pm$ 2.91	6.53 $\pm$ 0.17	0.06 $\pm$ 0.01	0.06 $\pm$ 0.00
	F1	73.63 $\pm$ 3.58	71.76 $\pm$ 8.86	70.42 $\pm$ 9.47	48.78 $\pm$ 0.63	75.77 $\pm$ 2.24	<b>76.48</b> $\pm$ 1.56
tol	NMI	3.04 $\pm$ 2.83	6.68 $\pm$ 0.98	6.69 $\pm$ 0.20	5.99 $\pm$ 0.05	6.67 $\pm$ 0.55	<b>6.73</b> $\pm$ 0.59
	F1	65.56 $\pm$ 10.49	72.10 $\pm$ 10.38	67.87 $\pm$ 4.74	59.17 $\pm$ 0.27	73.56 $\pm$ 1.94	<b>73.66</b> $\pm$ 2.10

experimental results for additional homophilous datasets in Appendix B. The dataset statistics can be consulted in Table 2, where the class insensitive edge homophily ratio  $\mathcal{H}(\mathcal{G})$  (Lim et al., 2021) is a homophily measure.

**Model selection and metrics** Model selection in this unsupervised setting is non-trivial, and the best metric depends on the task at hand. Therefore, this is not the scope of this paper and we assess our model agnostically to the model selection, and fairly w.r.t. to the baselines. We fix the hyperparameter configuration of the models across all datasets. We train for a fixed number of epochs and keep track of the evaluation metrics to report the best observed result. We repeat the training process 10 times and report average best results with standard deviations. We report the normalized mutual information (NMI) and pairwise F1-scores, based on the class labels.

**Baselines and hyperparameters** We compare our model against several methods, including a simple KMeans based on node attributes, adjacency partitioning-based approaches such as MinCutPool (Bianchi et al., 2020) and DMoN (Tsitsulin et al., 2023), as well as S<sup>3</sup>GC (Devvrit et al., 2022) and MUSE (Yuan et al., 2023), which represent the current state-of-the-art in homophilous and heterophilous node clustering, respectively. For HeNCler, we fix the hyperparameters to: MLP hidden dimensions 256, output dimensions 128, latent dimension  $s = 2 \times \#classes$ , learning rate 0.01, and epochs 300. For the baselines, we used their code implementations and the default hyperparameter settings as proposed by the authors. The number of clusters to infer is set to the number of classes cfr. Table 2 for all methods. The experiments are run on a Nvidia V100 GPU.

**Experimental results** Table 3 presents the experimental results for heterophilous graphs. HeNCler consistently demonstrates superior performance, significantly outperforming KMeans, MinCutPool, DMoN, S<sup>3</sup>GC, and MUSE, especially on the directed graphs. For undirected graphs, HeNCler also shows strong results, achieving the best performance in 5 out of 10 cases, compared to KMeans (1/10), MinCutPool (2/10), DMoN (1/10), S<sup>3</sup>GC (0/10), and MUSE (1/10). These results highlight HeNCler’s versatility and effectiveness in handling heterophilous graph structures.

**Ablation studies** We conduct several ablation studies, presented in Table 4. The ‘Undirected’ variant refers to a simplified, symmetric version of the model that uses a single MLP for both the  $\phi(\cdot)$  and  $\psi(\cdot)$  mappings, i.e.,  $\phi(\cdot) \equiv \psi(\cdot)$ . In this version, the model loses its asymmetry. The

Table 4: Ablation study results. We report mean NMI and F1 scores for 10 runs (higher is better) for different model configurations. Best results are highlighted in bold.

Metric	tex		corn		cha		rom		tol	
	NMI	F1	NMI	F1	NMI	F1	NMI	F1	NMI	F1
Undirected	27.58	65.20	18.12	53.69	19.91	44.08	33.17	33.57	6.33	73.89
Reconstr only	29.54	66.64	27.76	54.70	22.02	43.42	<b>40.05</b>	35.16	6.18	68.60
wKSVD only	31.64	62.83	20.63	47.12	22.60	42.98	35.99	35.30	4.42	68.45
HeNCler	<b>43.65</b>	<b>71.39</b>	<b>41.52</b>	<b>63.40</b>	<b>23.89</b>	<b>44.14</b>	36.99	<b>35.43</b>	<b>6.73</b>	<b>73.66</b>

'wKSVD only' and 'Reconstr only' variations reflect models that incorporate only the wKSVD loss ( $\mathcal{L}_{\text{wKSVD}}$ ) and the reconstruction losses ( $\mathcal{L}_{\text{NodeRec}} + \mathcal{L}_{\text{EdgeRec}}$ ), respectively. Interestingly, as shown in Table 4, even for undirected graphs, introducing asymmetry in HeNCler enhances clustering performance. Furthermore, all loss components are shown to contribute positively to HeNCler's overall performance. For a comprehensive analysis, including results across all datasets and standard deviations, we refer the reader to Table 7 in Appendix B.

## 5 DISCUSSION

A key motivation behind HeNCler is to learn a new graph representation where nodes belonging to the same cluster are positioned closer together, driven by the clustering objective. This results in spectral biclustering embeddings that exhibit improved cluster-ability. Note that HeNCler uses KMeans to obtain cluster assignments. Therefore, the comparisons between HeNCler and KMeans, as shown in Tables 3 and 6, demonstrate that our model enhances the cluster-ability of the node representations relative to the original input features.

The asymmetry in HeNCler eliminates the undirected constraints of traditional adjacency partitioning-based models, enabling superior performance on directed graphs, as shown in Table 3. Furthermore, our ablation study in Table 4 shows that, while most of the performance on undirected graphs stem from the graph learning component, HeNCler is additionally able to capture and learn meaningful asymmetric information. This capacity to extract valuable asymmetric insights from symmetric data is a common occurrence in KSVD frameworks (He et al., 2023; Tao et al., 2024). Importantly, thanks to the added performance boost from asymmetry, on top of the benefits from similarity learning, HeNCler outperforms state-of-the-art models, even when applied to undirected graphs.

We visualize the learned similarity matrix  $S = \Phi\Psi^T$  for two datasets in Figure 3. These matrices are generally asymmetric, with the asymmetry particularly pronounced in the directed graph of the Wisconsin dataset. In contrast, the Roman-Empire dataset, which is represented by an undirected graph, exhibits less asymmetry in the learned similarity matrix. This demonstrates the adaptability of HeNCler to handle both directed and undirected graphs. Further, given the observable block structures, the learned similarities are meaningful w.r.t. to the ground truth node labels. Note however that our model operates in the primal setting and directly projects the learned mappings  $\phi$  and  $\psi$  to their final embeddings  $e$  and  $r$  using  $U$  and  $V$  respectively, avoiding quadratic space complexity and cubic time complexity of the SVD. This is the motivation of employing a kernel based method, and exploiting the primal-dual framework that comes with it. In fact, the matrices in Figure 3 are only constructed for the sake of this visualization.

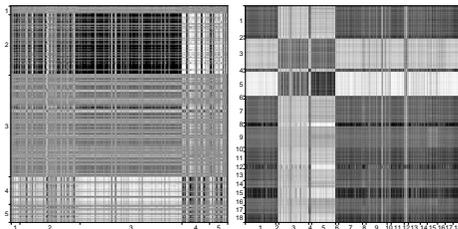


Figure 3: The learned matrix  $S = \Phi\Psi^T$  for the Wisconsin (left) and Roman-empire (right) dataset. Rows and columns are grouped according to ground-truth node labels.

**Computational complexity** The space and time complexity of the current implementation of HeNCler are both linear w.r.t. the number of nodes  $\mathcal{O}(|\mathcal{V}|)$ . Whereas MinCutPool and DMoN need

all the node attributes in memory to calculate the loss w.r.t. the full adjacency matrix, HeNCler is easily adaptable to work with minibatches which reduces space complexity to the minibatch size  $\mathcal{O}(|\mathcal{B}|)$ . Although HeNCler relies on edge reconstruction, the edge sampling avoids quadratic complexity w.r.t. number of nodes, and is specifically designed to scale with the number of nodes, rather than the number of edges. Assuming the graphs are sparse, we add an overview of space and time complexity w.r.t. the number of nodes and edges for all methods in Table 1. A detailed table with measured computation times is provided in Appendix C.

## 6 CONCLUSION AND FUTURE WORK

We tackle three limitations of current node clustering algorithms, that prevent these methods from effectively clustering nodes in heterophilous graphs: they assume homophily in their loss, they are only defined for undirected graphs and/or they lack a specific focus on clustering.

To this end, we introduce a weighted kernel SVD framework and harness its primal-dual equivalences. HeNCler relies on the dual interpretation for its theoretical motivation, while it benefits from the computational advantages of its implementation in the primal. In an end-to-end fashion, it learns new similarities, which are asymmetric where necessary, and node embeddings resulting from the spectral biclustering interpretation of these learned similarities. As empirical evidence shows, our approach effectively eliminates the aforementioned limitations, significantly outperforming current state-of-the-art alternatives.

HeNCler is the first heterophilous node clustering model that does not rely on contrastive learning techniques. Future research could explore the integration of contrastive learning into HeNCler, potentially combining the strengths of both approaches. Another next step can be to investigate how to do the cluster assignments in a graph pooling setting (i.e., differentiable graph coarsening), to enable end-to-end learning for downstream graph prediction tasks.

## REFERENCES

- Sonny Achten, Arun Pandey, Hannes De Meulemeester, Bart De Moor, and Johan A. K. Suykens. Duality in Multi-View Restricted Kernel Machines. *ICML Workshop on Duality for Modern Machine Learning*, 2023. arXiv:2305.17251 [cs].
- Sonny Achten, Francesco Tonin, Panagiotis Patrinos, and Johan A.K. Suykens. Unsupervised Neighborhood Propagation Kernel Layers for Semi-supervised Node Classification. *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(10):10766–10774, Mar. 2024.
- Carlos Alzate and Johan A. K. Suykens. Multiway Spectral Clustering with Out-of-Sample Extensions through Weighted Kernel PCA. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 32(2):335–347, 2010.
- Filippo Maria Bianchi, Daniele Grattarola, and Cesare Alippi. Spectral Clustering with Graph Neural Networks for Graph Pooling. In *International Conference on Machine Learning*, 2020.
- Stephen Boyd and Lieven Vandenberghe. *Convex Optimization*. Cambridge University Press, New York, NY, USA, 2004. ISBN 0521833787.
- Jingfan Chen, Guanghui Zhu, Yifan Qi, Chunfeng Yuan, and Yihua Huang. Towards self-supervised learning on graphs with heterophily. In *ACM International Conference on Information & Knowledge Management*, pp. 201–211, 2022.
- Yingyi Chen, Qinghua Tao, Francesco Tonin, and Johan Suykens. Primal-Attention: Self-attention through Asymmetric Kernel SVD in Primal Representation. In *Advances in Neural Information Processing Systems*, 2023.
- Fnu Devvrit, Aditya Sinha, Inderjit Dhillon, and Prateek Jain. S3GC: Scalable Self-Supervised Graph Clustering. In *Advances in Neural Information Processing Systems*, 2022.
- Inderjit S. Dhillon. Co-clustering documents and words using bipartite spectral graph partitioning. In *ACM SIGKDD international conference on Knowledge discovery and data mining*, 2001.

- 540 Vijay Prakash Dwivedi and Xavier Bresson. A generalization of transformer networks to  
541 graphs. *AAAI Workshop on Deep Learning on Graphs: Methods and Applications*, 2021.  
542 arXiv:2012.09699 [cs].
- 543 Vijay Prakash Dwivedi, Anh Tuan Luu, Thomas Laurent, Yoshua Bengio, and Xavier Bresson.  
544 Graph neural networks with learnable structural and positional representations. In *International  
545 Conference on Learning Representations*, 2022.  
546
- 547 Zheng Gong, Guifeng Wang, Ying Sun, Qi Liu, Yuting Ning, Hui Xiong, and Jingyu Peng. Beyond  
548 Homophily: Robust Graph Anomaly Detection via Neural Sparsification. In *International Joint  
549 Conference on Artificial Intelligence*, 2023.
- 550 William L. Hamilton, Rex Ying, and Jure Leskovec. Inductive Representation Learning on Large  
551 Graphs. In *Advances in Neural Information Processing Systems*, 2017.  
552
- 553 Kaveh Hassani and Amir Hosein Khas Ahmadi. Contrastive multi-view representation learning on  
554 graphs. In *International Conference on Machine Learning*, 2020.
- 555 Mingzhen He, Fan He, Lei Shi, Xiaolin Huang, and Johan A. K. Suykens. Learning With Asymmetric  
556 Kernels: Least Squares and Feature Interpretation. *IEEE Transactions on Pattern Analysis and  
557 Machine Intelligence*, 45(8):10044–10054, 2023.  
558
- 559 Diederik P. Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization. In *International  
560 Conference on Learning Representations*, 2015.
- 561 Thomas N. Kipf and Max Welling. Semi-Supervised Classification with Graph Convolutional  
562 Networks. In *International Conference on Learning Representations*, 2017.  
563
- 564 Cornelius Lanczos. Linear systems in self-adjoint form. *The American Mathematical Monthly*, 9(65):  
565 665–679, 1958.
- 566 Derek Lim, Felix Hohne, Xiuyu Li, Sijia Linda Huang, Vaishnavi Gupta, Omkar Bhalerao, and  
567 Ser-Nam Lim. Large Scale Learning on Non-Homophilous Graphs: New Benchmarks and Strong  
568 Simple Methods. In *Advances in Neural Information Processing Systems*, 2021.
- 569 Shirui Pan, Ruiqi Hu, Sai-Fu Fung, Guodong Long, Jing Jiang, and Chengqi Zhang. Learning  
570 graph embedding with adversarial training methods. *IEEE Transactions on Cybernetics*, 50(6):  
571 2475–2487, June 2020. ISSN 2168-2275.  
572
- 573 Arun Pandey, Joachim Schreurs, and Johan A.K. Suykens. Generative restricted kernel machines: A  
574 framework for multi-view generation and disentangled feature learning. *Neural Networks*, 135:  
575 177–191, 2021.
- 576 Arun Pandey, Michaël Fanuel, Joachim Schreurs, and Johan A. K. Suykens. Disentangled representa-  
577 tion learning and generation with manifold optimization. *Neural Computation*, 34(10):2009–2036,  
578 09 2022. ISSN 0899-7667.
- 579 Jiwoong Park, Minsik Lee, Hyung Jin Chang, Kyuewang Lee, and Jin Young Choi. Symmetric  
580 graph convolutional autoencoder for unsupervised graph representation learning. In *IEEE/CVF  
581 International Conference on Computer Vision (ICCV)*, October 2019.  
582
- 583 Hongbin Pei, Bingzhe Wei, Kevin Chen-Chuan Chang, Yu Lei, and Bo Yang. Geom-GCN: Geometric  
584 Graph Convolutional Networks. In *International Conference on Learning Representations*, 2020.  
585
- 586 Oleg Platonov, Denis Kuznedelev, Michael Diskin, Artem Babenko, and Liudmila Prokhorenkova.  
587 A critical look at the evaluation of GNNs under heterophily: Are we really making progress? In  
588 *International Conference on Learning Representations*, 2023.
- 589 R. Tyrrell Rockafellar. *Conjugate Duality and Optimization*. SIAM, 1974.
- 590 Benedek Rozemberczki, Carl Allen, and Rik Sarkar. Multi-Scale attributed node embedding. *Journal  
591 of Complex Networks*, 9(1):1–22, 2021.  
592
- 593 Ruslan Salakhutdinov. Learning Deep Generative Models. *Annual Review of Statistics and Its  
Application*, 2(1):361–385, 2015.

- 594 Prithviraj Sen, Galileo Namata, Mustafa Bilgic, Lise Getoor, Brian Galligher, and Tina Eliassi-Rad.  
595 Collective Classification in Network Data. *AI magazine*, 29(3):93–93, 2008.  
596
- 597 Dengdi Sun, Dashuang Li, Zhuanlian Ding, Xingyi Zhang, and Jin Tang. Dual-decoder graph  
598 autoencoder for unsupervised graph representation learning. *Knowledge-Based Systems*, 234:  
599 107564, December 2021. ISSN 09507051.  
600
- 601 Johan A. K. Suykens. SVD revisited: A new variational principle, compatible feature maps and  
602 nonlinear extensions. *Applied and Computational Harmonic Analysis*, 40(3):600–609, May 2016.  
603 ISSN 1063-5203.
- 604 Johan A. K. Suykens. Deep Restricted Kernel Machines Using Conjugate Feature Duality. *Neural  
605 Computation*, 29(8):2123–2163, 2017.  
606
- 607 Johan A. K. Suykens, Tony Van Gestel, Jos De Brabanter, Bart De Moor, and Joos Vandewalle. *Least  
608 Squares Support Vector Machines*. World Scientific, Singapore, 2002.  
609
- 610 Qinghua Tao, Francesco Tonin, Alex Lambert, Yingyi Chen, Panagiotis Patrinos, and Johan A. K.  
611 Suykens. Learning in Feature Spaces via Coupled Covariances: Asymmetric Kernel SVD and  
612 Nyström method. International Conference on Machine Learning, 2024.
- 613 Francesco Tonin, Panagiotis Patrinos, and Johan A. K. Suykens. Unsupervised learning of disen-  
614 tangled representations in deep restricted kernel machines with orthogonality constraints. *Neural  
615 Networks*, 142:661–679, 2021.  
616
- 617 Anton Tsitsulin, John Palowitch, Bryan Perozzi, and Emmanuel Müller. Graph Clustering with Graph  
618 Neural Networks. *Journal of Machine Learning Research*, 24(127):1–21, 2023.  
619
- 620 Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua  
621 Bengio. Graph Attention Networks. In *International Conference on Learning Representations*,  
622 2018.
- 623 Zonghan Wu, Shirui Pan, Fengwen Chen, Guodong Long, Chengqi Zhang, and Philip S. Yu. A  
624 Comprehensive Survey on Graph Neural Networks. *IEEE Transactions on Neural Networks and  
625 Learning Systems*, 32(1):4–24, 2021.  
626
- 627 Teng Xiao, Zhengyu Chen, Zhimeng Guo, Zeyang Zhuang, and Suhang Wang. Decoupled self-  
628 supervised learning for graphs. In *Advances in Neural Information Processing Systems*, 2022.  
629
- 630 Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How Powerful are Graph Neural  
631 Networks? In *International Conference on Learning Representations*, 2019.  
632
- 633 Zhilin Yang, William Cohen, and Ruslan Salakhudinov. Revisiting Semi-Supervised Learning with  
634 Graph Embeddings. In *International Conference on Machine Learning*, 2016.
- 635 Chengxuan Ying, Tianle Cai, Shengjie Luo, Shuxin Zheng, Guolin Ke, Di He, Yanming Shen, and  
636 Tie-Yan Liu. Do Transformers Really Perform Badly for Graph Representation? In *Advances in  
637 Neural Information Processing Systems*, 2022.  
638
- 639 Yuning You, Tianlong Chen, Yongduo Sui, Ting Chen, Zhangyang Wang, and Yang Shen. Graph  
640 contrastive learning with augmentations. In *Advances in Neural Information Processing Systems*,  
641 2020.  
642
- 643 Mengyi Yuan, Minjie Chen, and Xiang Li. MUSE: Multi-View Contrastive Learning for Heterophilic  
644 Graphs. In *ACM International Conference on Information and Knowledge Management*, pp.  
645 3094–3103, 2023.
- 646 Xin Zheng, Yi Wang, Yixin Liu, Ming Li, Miao Zhang, Di Jin, Philip S. Yu, and Shirui Pan. Graph  
647 neural networks for graphs with heterophily: A survey, 2024. arXiv:2202.07082 [cs].

## A NOTE ON FEATURE MAP CENTERING

In the wKSVD framework, we assume that the feature maps are centered. More precisely, given two arbitrary mappings  $\phi(\cdot)$  and  $\psi(\cdot)$ , the centered mappings are obtained by subtracting the weighted mean:

$$\begin{aligned}\phi_c(\mathbf{x}_i) &= \phi(\mathbf{x}_i) - \frac{\sum_{k=1}^n w_{1,k} \phi(\mathbf{x}_k)}{\sum_{k=1}^n w_{1,k}}, \\ \psi_c(\mathbf{z}_j) &= \psi(\mathbf{z}_j) - \frac{\sum_{l=1}^m w_{2,l} \psi(\mathbf{z}_l)}{\sum_{l=1}^m w_{2,l}}.\end{aligned}$$

Although we use the primal formulation in this paper, we next show how to obtain this centering in the dual for the sake of completeness. When using a kernel function or a given similarity matrix, one has no access to the explicit mappings and has to do an equivalently centering in the dual using:

$$S_c = M_1 S M_2^\top,$$

where  $M_1$  and  $M_2$  are the centering matrices:

$$\begin{aligned}M_1 &= I_n - \frac{1}{\mathbf{1}_n^\top W_1 \mathbf{1}_n} \mathbf{1}_n \mathbf{1}_n^\top W_1 \\ M_2 &= I_m - \frac{1}{\mathbf{1}_m^\top W_2 \mathbf{1}_m} \mathbf{1}_m \mathbf{1}_m^\top W_2,\end{aligned}$$

with  $I_n$  and  $\mathbf{1}_n$  a  $n \times n$  identity matrix and a  $n$ -dimensional all-ones vector respectively. We omit the subscript  $c$  in the paper and assume the feature maps are always centered. Note that this can easily be achieved in the implementations by using the above equations.

## B ADDITIONAL EXPERIMENTS

**Homophilous experiments** Although our work primarily focuses on heterophilous graphs, we further evaluate our model on homophilous citation networks Cora, Citeseer, and PubMed (Sen et al., 2008; Yang et al., 2016). The dataset statistics can be consulted in Table 5. We employ the same experimental setup as for the heterophilous datasets and report the experimental results in Table 6. While S<sup>3</sup>GC achieves the best overall performance due to its alignment with the homophily assumption, HeNCler outperforms adjacency partitioning methods like MinCutPool and DMoN. Additionally, HeNCler demonstrates competitive performance with MUSE, the state-of-the-art in heterophilous node clustering, further validating its robustness across different graph types.

Table 5: Dataset statistics of the employed homophilous graphs.

Dataset	short	# Nodes	# Edges	# Classes	Directed	$\mathcal{H}(\mathcal{G})$
Cora	cora	2,708	5,278	7	✗	0.765
CiteSeer	cite	3,327	4,614	6	✗	0.627
Pubmed	pub	19,717	44,325	3	✗	0.664

**Comprehensive ablation study** We provide the full ablation study results in Table 7, including all datasets and standard deviations. We compare HeNCler with three simplified versions. 'Undirected' reflects an undirected variant of the model with a single MLP decoder. 'wKSVD only' and 'Reconstr only' is the model where only the wKSVD loss  $\mathcal{L}_{wKSVD}$  and the reconstruction losses  $\mathcal{L}_{\text{NodeRec}} + \mathcal{L}_{\text{EdgeRec}}$  are used respectively. We observe that HeNCler performs better than its undirected version, even for undirected graphs, and that all loss terms contribute to HeNCler's performance.

## C COMPUTATION TIMES

We trained MinCutPool, DMoN, and HeNCler for 300 iterations; and S<sup>3</sup>GC for 30 iterations on a Nvidia V100 GPU, and report the computation times in Table 8. Figure 4 visualises these result w.r.t. the number of nodes in the graph, showing the linear time complexity of HeNCler and that it is insensitive to the number of edges. We conclude that HeNCler demonstrates fast computation times.

Table 6: Experimental results on homophilous graphs. We report NMI and F1 scores for 10 runs (mean  $\pm$  standard deviation), where higher values indicate better performance. The best results for each metric are highlighted in bold.

Dataset		Baselines				Ours	
		KMeans	MinCutPool	DMoN	S <sup>3</sup> GC	MUSE	HeNCler
cora	NMI	35.0 $\pm$ 3.21	49.0 $\pm$ 2.24	51.7 $\pm$ 1.63	<b>53.62</b> $\pm$ 0.55	36.45 $\pm$ 2.71	38.81 $\pm$ 2.26
	F1	36.0 $\pm$ 2.12	47.1 $\pm$ 1.78	51.8 $\pm$ 2.02	<b>60.12</b> $\pm$ 0.46	50.78 $\pm$ 2.79	47.93 $\pm$ 2.60
cite	NMI	19.9 $\pm$ 2.90	29.5 $\pm$ 3.21	30.3 $\pm$ 1.09	<b>43.56</b> $\pm$ 0.65	39.03 $\pm$ 1.99	34.83 $\pm$ 2.21
	F1	39.4 $\pm$ 3.07	47.1 $\pm$ 1.21	57.4 $\pm$ 3.42	<b>64.12</b> $\pm$ 0.28	52.89 $\pm$ 1.68	48.70 $\pm$ 2.79
pub	NMI	31.4 $\pm$ 2.18	21.4 $\pm$ 1.46	25.7 $\pm$ 2.46	31.01 $\pm$ 2.35	<b>36.09</b> $\pm$ 3.26	27.26 $\pm$ 1.72
	F1	59.2 $\pm$ 2.32	44.5 $\pm$ 2.47	34.3 $\pm$ 2.05	<b>69.12</b> $\pm$ 1.39	61.26 $\pm$ 1.50	51.17 $\pm$ 1.75

Table 7: Full Ablation study results. We report NMI and F1 scores for 10 runs (mean  $\pm$  standard deviation in %) where higher is better. Best results are highlighted in bold.

dataset		ablations			full model
		Undirected	Reconstr only	wKSVD only	HeNCler
tex	NMI	27.58 $\pm$ 4.75	29.54 $\pm$ 2.27	31.64 $\pm$ 2.14	<b>43.65</b> $\pm$ 2.52
	F1	65.20 $\pm$ 2.06	66.64 $\pm$ 1.83	62.83 $\pm$ 3.91	<b>71.39</b> $\pm$ 2.16
corn	NMI	18.12 $\pm$ 2.57	27.76 $\pm$ 3.29	20.63 $\pm$ 5.92	<b>41.52</b> $\pm$ 4.35
	F1	53.69 $\pm$ 0.98	54.70 $\pm$ 1.88	47.12 $\pm$ 2.61	<b>63.40</b> $\pm$ 3.67
wis	NMI	25.08 $\pm$ 3.54	34.65 $\pm$ 1.86	39.86 $\pm$ 4.63	<b>47.13</b> $\pm$ 1.76
	F1	57.13 $\pm$ 1.34	62.28 $\pm$ 1.58	63.60 $\pm$ 2.46	<b>68.30</b> $\pm$ 2.17
cha	NMI	19.91 $\pm$ 0.48	22.02 $\pm$ 0.25	22.60 $\pm$ 0.57	<b>23.89</b> $\pm$ 0.84
	F1	44.08 $\pm$ 1.79	43.42 $\pm$ 1.63	42.98 $\pm$ 0.37	<b>44.14</b> $\pm$ 1.83
squi	NMI	9.59 $\pm$ 0.21	9.59 $\pm$ 0.27	9.56 $\pm$ 0.19	<b>9.67</b> $\pm$ 0.13
	F1	<b>55.43</b> $\pm$ 0.03	53.74 $\pm$ 3.77	36.42 $\pm$ 1.85	36.51 $\pm$ 2.39
rom	NMI	33.17 $\pm$ 1.25	<b>40.05</b> $\pm$ 0.82	35.99 $\pm$ 0.95	36.99 $\pm$ 0.61
	F1	33.57 $\pm$ 2.15	35.16 $\pm$ 1.34	35.30 $\pm$ 0.97	<b>35.43</b> $\pm$ 1.07
mine	NMI	<b>0.08</b> $\pm$ 0.02	0.07 $\pm$ 0.02	0.04 $\pm$ 0.01	0.06 $\pm$ 0.00
	F1	76.15 $\pm$ 2.25	76.05 $\pm$ 2.16	73.77 $\pm$ 3.40	<b>76.48</b> $\pm$ 1.56
tol	NMI	6.33 $\pm$ 0.94	6.18 $\pm$ 0.67	4.42 $\pm$ 0.54	<b>6.73</b> $\pm$ 0.59
	F1	73.89 $\pm$ 4.00	68.60 $\pm$ 5.95	68.45 $\pm$ 7.57	<b>73.66</b> $\pm$ 2.10

Table 8: Computation times in seconds.

DATASET	BASELINES			OURS
	MINCUTP.	DMoN	S <sup>3</sup> GC	HENCLEP
CHAMELEON	8	20	89	24
SQUIRREL	14	86	105	49
ROMAN-EMPIRE	67	71	312	57
AMAZON-RATING	109	93	195	63
MINESWEEPER	55	27	98	21
TOLOKERS	71	198	100	35
QUESTIONS	340	215	217	125

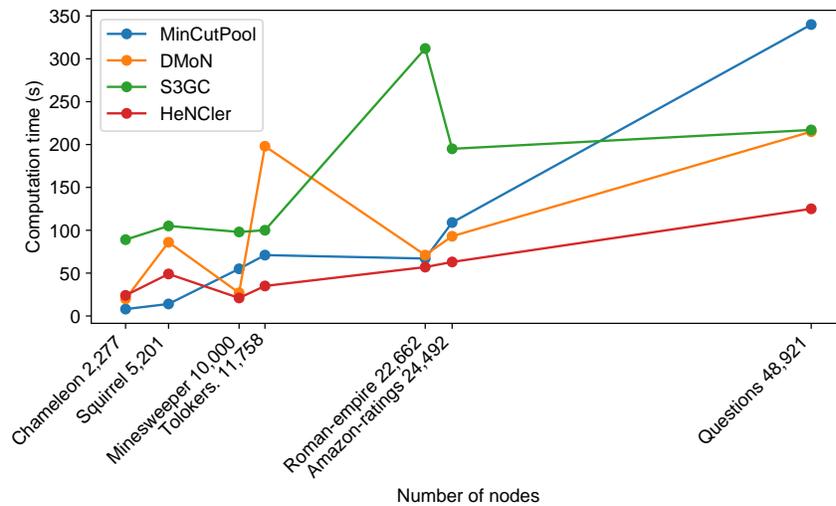


Figure 4: Computation times of MinCutPool, DMoN, S<sup>3</sup>GC, and HeNcler w.r.t. the number of nodes of the datasets. We observe that HeNcler scales linearly with the number of nodes, and that it is not sensitive to the number of edges, as opposed to DMoN, showing a significant peak for the Tolokers dataset due the large number of edges in this graph.