Label-Wise Graph Convolutional Network for Heterophilic Graphs

Anonymous Author(s) Anonymous Affiliation Anonymous Email

Abstract

Graph Neural Networks (GNNs) have achieved remarkable performance in mod-2 eling graphs for various applications. However, most existing GNNs assume the 3 graphs exhibit strong homophily in node labels, i.e., nodes with similar labels 4 are connected in the graphs. They fail to generalize to heterophilic graphs where 5 linked nodes may have dissimilar labels and attributes. Therefore, in this pa-6 per, we investigate a novel framework that performs well on graphs with either 7 homophily or heterophily. More specifically, we propose a label-wise message 8 passing mechanism to avoid the negative effects caused by aggregating dissimilar 9 node representations and preserve the heterophilic contexts for representation learn-10 ing. We further propose a bi-level optimization method to automatically select the 11 model for graphs with homophily/heterophily. Theoretical analysis and extensive 12 experiments demonstrate the effectiveness of our proposed framework for node 13 classification on both homophilic and heterophilic graphs. 14

15 **1** Introduction

Graph-structured data is very pervasive in the real-world such as knowledge graphs, traffic networks, and social networks. Therefore, it is important to model graphs for downstream tasks such as traffic prediction [39], recommendation system [14] and drug generation [3]. To capture the topology information in graph-structured data, Graph Neural Networks (GNNs) [36] adopt a message-passing mechanism which learns a node's representation by iteratively aggregating the representations of its neighbors. This can enrich the node features and preserve the node attributes and local topology for various downstream tasks.

Despite the great success of GNNs in modeling graphs, there is a concern in processing heterophilic 23 graphs where edges often link nodes dissimilar in attributes or labels. Specifically, existing works [42, 24 25 9] find that GNNs could fail to generalize to graphs with heterophily due to their implicit/explicit homophily assumption. For example, Graph Convolutional Network (GCNs) is even outperformed 26 by MLP that ignores the graph structure on heterophilic website datasets [42]. However, a recent 27 work [26] argues that homophily assumption is not a necessity for GNNs. They show that GCN 28 can work well on dense heterophilic graphs whose neighborhood patterns of different classes are 29 distinguishable. But their analysis and conclusion is limited to the heterophilic graphs under strict 30 conditions, and fails to show the relation between heterophily levels and performance of GNNs. Thus, 31 in Sec. 3, we conduct thoroughly theoretical and empirical analysis on GCN to investigate the impacts 32 of heterophily levels, which cover all the aforementioned observations. As the Theorem 1 and Fig. 1 33 show, the performance of GCN will firstly decrease then increase with the increment of heterophily 34 levels. And the aggregation in GCN could even lead to non-discriminative representations under 35 certain conditions. 36

Though heterophilic graphs challenge existing GNNs, the heterophilic neighborhood context itself provides useful information [26, 6]. Generally, two nodes of the same class tend to have similar heterophilic neighborhood contexts; while two nodes of different classes are more likely to have different heterophilic neighborhood contexts, which is verified in Appendix F. Thus, a heterophilic context-preserving mechanism can lead to more discriminative representations. One promising way

Submitted to the First Learning on Graphs Conference (LoG 2022). Do not distribute.

to preserve the heterophilic context is to conduct label-wise aggregation, i.e., separately aggregate 42 neighbors in each class. In this way, we can summarize the heterophilic neighbors belonging to 43 each class to an embedding to preserve the local context information for representation learning. As 44 shown in the example in Fig. 2, for node v_A , with label-wise aggregation, v_A will be represented 45 as [1.0, 5.5, 2.0, non-existence], in the order of v_A 's attribute, blue, green, and orange neighbors, 46 respectively. Compared with v_B , v_A 's representations of central node and neighborhood context 47 differ significantly with v_B . While for the aggregation in GCN, the obtained representations are 48 rather similar for two nodes. In other words, we obtain more discriminative features on heterophilic 49 graphs with label-wise aggregation, which is also verified by our analysis in Theorem 2. Though 50 promising, there is no existing work on exploring label-wise message passing to address the challenge 51 of heterophilic graphs. 52

Therefore, in this paper, we investigate novel label-wise aggregation for graph convolution to facilitate 53 the node classification on heterophilic graphs. In essence, we are faced with two challenges: (i) 54 the label-wise aggregation needs the label of each node; while for node classification, we are only 55 given a small set of labeled nodes. How to adopt label-wise graph convolution on sparsely labeled 56 heterophilic graphs to facilitate node classification? (ii) In practice, the homophily levels of the 57 given graphs can be various and are often unknown. For homophily graphs, the label-wise graph 58 convolution might not work as well as previous GNNs embedded with homophily assumption. How 59 to ensure the performance on both heterophilic and homophilic graphs? In an attempt to address these 60 challenges, we propose a novel framework Label-Wise GCN (LW-GCN). LW-GCN adopts a pseudo 61 label predictor to predict pseudo labels and designs a novel label-wise message passing to preserve 62 63 the heterophilic contexts with pseudo labels. To handle both heterophilic and homophilic graphs, apart from label-wise message passing GNN, LW-GCN also utilizes a GNN for homophilic graphs, 64 and adopts bi-level optimization on the validation data to automatically select the better model for the 65 given graph. The main contributions are: 66

- We theoretically show impacts of heterophily levels to GCN and demonstrate the potential limitations of GCN in learning on heterophilic graphs;
- We design a label-wise graph convolution to preserve the local context in heterophilic graphs,
 which is also proven by our theoretical and empirical analysis;
- We propose a novel framework LW-GCN, which deploys a pseudo label predictor and an automatic
 model selection module to achieve label-wise aggregation on sparsely labeled graphs and ensure
 the performance on both heterophilic and homophilic graphs; and
- Extensive experiments on real-world graphs with heterophily and homophily are conducted to demonstrate the effectiveness of LW-GCN.

76 2 Related Work

Graph neural networks (GNNs) have shown great success for various applications such as social 77 networks [14, 10], financial transaction networks [34, 12] and traffic networks [39, 40]. Based on 78 the definition of the graph convolution, GNNs can be categorized into two categories, i.e., spectral-79 based [4, 11, 19, 21] and spatial-based [33, 37, 1]. Spectral-based GNN models are defined according 80 to spectral graph theory. Bruna et al. [4] firstly generalize convolution operation to graph-structured 81 data from spectral domain. GCN [19] simplifies the graph convolution by first-order approximation. 82 For spatial-based graph convolution, it aggregates the information of the neighbors nodes [27, 14, 7]. 83 For instance, spatial graph convolution that incorporates the attention mechanism is applied in 84 85 GAT [33] to facilitate the information aggregation. Recently, to learn better node representations, deep graph neural networks [8, 20, 22] and self-supervised learning methods [32, 18, 43, 30, 38] 86 have been investigated. 87

However, the aforementioned methods are generally designed based on the homophily assumption of 88 the graph. Low homophily level in some real-word graphs can largely degrade their performance [42]. 89 Some efforts [28, 2, 16, 42, 41, 9, 15, 24, 23] have been taken to address the problem of heterophilic 90 graphs. For example, H2GCN [42] investigates three key designs for GNNs on heterophilic graphs. 91 SimP-GCN [16] adopts a node similarity preserving mechanism to handle graphs with heterophiliy. 92 FAGCN [2] adaptively aggregates low-frequency and high-frequency signals from neighbors to learn 93 representations for graphs with heterophily. GPR-GNN [9] proposes a generalized PageRank GNN 94 architecture that can learn positive/negative weights for the representations after different steps of 95 propagation to mitigate the graph heterophily issue. Recently, BM-GCN [15] proposes to utilize 96

pseudo labels in the convolutional operation. Specifically, the pseudo labels are used to obtain a 97

block similarity matrix to re-weight the edges in heterophilic graphs. Then, node pairs belonging to 98

different label combinations could have different information exchange. Our LW-GCN is inherently 99 different from these methods: (i) We propose a novel label-wise graph convolution to better capture

100 the neighbors' information in heterophilic graphs; and (ii) Automatic model selection is deployed to

achieve state-of-the-art performance on both homophilic and heterophilic graphs.

3 **Preliminaries** 103

In this section, we first present the notations and definition followed by the introduction of the GCN's 104 design. We then conduct the theoretical analysis to investigate the impacts of heterophily to GCN. 105

Notations and Definition 3.1 106

Let $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathbf{X})$ be an attributed graph, where $\mathcal{V} = \{v_1, ..., v_N\}$ is the set of N nodes, $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$ 107 is the set of edges, and $\mathbf{X} = {\{\mathbf{x}_1, ..., \mathbf{x}_N\}}$ is the set of node attributes. $\mathbf{A} \in \mathbb{R}^{N \times N}$ represents 108 the adjacency matrix of the graph \mathcal{G} , where $\mathbf{A}_{ij} = 1$ indicates an edge between nodes v_i and v_j ; otherwise, $A_{ij} = 0$. In the node classification task, each node belongs to one of C classes. We use y_i 110 to denote label of node v_i . Graphs can be split into homophilic and heterophilic graphs based on how 111 likely edges link nodes in the same class. The homophily level is measured by the homophily ratio: 112

Definition 1 (Homophily Ratio) It is the fraction of edges in a graph that connect nodes of the 113 same class. The homophily ratio h is calculated as $h = \frac{|\{(v_i, v_j) \in \mathcal{E}: y_i = y_j\}|}{|\mathcal{E}|}$. 114

When the homophily ratio is small, most of the edges will link nodes from different classes, which 115

indicates a heterophilic graph. In homophilic graphs, connected nodes are more likely to belong to 116

the same class, which will lead to a homophily ratio close to 1. 117

3.2 How does the Heterophily Affect the GCN? 118

GCN [19] is one of the most widely used graph neural networks. The operation in each layer of GCN 119 can be written as: 120

$$\mathbf{H}^{(k+1)} = \sigma(\tilde{\mathbf{A}}\mathbf{H}^{(k)}\mathbf{W}^{(k)}),\tag{1}$$

where $\mathbf{H}^{(k)}$ is the node representation matrix of the output of the k-th layer and $\tilde{\mathbf{A}}$ is the normalized 121 adjacency matrix. Generally, the symmetric normalized form $\mathbf{D}^{-\frac{1}{2}}\mathbf{A}\mathbf{D}^{-\frac{1}{2}}$ or $\mathbf{D}^{-1}\mathbf{A}$ is used as $\tilde{\mathbf{A}}$, where **D** is a diagonal matrix with $\mathbf{D}_{ii} = \sum_{i} \mathbf{A}_{ij}$. The adjacency matrix can be augmented with 123 a self-loop. σ is an activation function such as ReLU. In a single layer of GCN, the process can 124 be split into two steps. First, GCN layer averages the neighbor features with $\mathbf{Z} = \mathbf{A}\mathbf{X}$. Then, a 125 non-linear transformation $\sigma(\mathbf{ZW})$ is applied to obtain intermediate features or final predictions. The 126 step of averaging the neighbor features can benefit the node classification when the neighbors have 127 similar features. However, for heterophilic graphs, mixing neighbors that possess different features 128 129 may result in poor representations for node classification. This could be justified by the following theorem, which thoroughly analyzes the impacts of the heterophily level to the linear separability of 130 the representations after one step aggregation in GCN. 131

Assumptions. We first discuss the assumptions of the heterophilic graphs: (i) Following previous 132 works [42], the graph G is considered as a *d*-regular graph, i.e., each node has *d* neighbors; For 133 each node v, the label distribution of its neighbor node $u \in \mathcal{N}(v)$ follows $P(y_u = y_v|y_v) = h, P(y_u = y|y_v) = \frac{1-h}{C-1}, \forall y \neq y_v$. (ii) For nodes in different classes, their heterophilic neighbors' features follow different distributions and dimensions of features are independent to each other. 134 135 136 Specifically, let $\mathcal{N}_k(v)$ denote node v's neighbors of class k. For two nodes v and s in classes i and j 137 $(i \neq j)$, the features of their heterophilic neighbors $\mathcal{N}_k(v)$ and $\mathcal{N}_k(s)$ in class $k \in \{1, ..., C\}$ follow 138 two different normal distributions $N(\mu_{ik}, \sigma_{ik})$ and $N(\mu_{jk}, \sigma_{jk})$, where μ_{ik} and μ_{jk} represent 139 the means, σ_{ik} and σ_{ik} denote the standard deviations. Intuitively, though nodes in $\mathcal{N}_k(v)$ and 140 $\mathcal{N}_k(s)$ belong to the same class k, they are connected to nodes of different classes because of 141 their different properties. For example, in the molecule, the atom in the same class will exhibit 142 different features, when they are linked to different atoms. Therefore, this assumption is valid. And 143 it is also verified by the empirical analysis on large real-world heterophilic graphs in Appendix F. 144 Let $\sigma_i = \sqrt{\frac{1}{C}\sum_{k=1}^{C}(\mu_{ik} - \bar{\mu}_i) \bigodot (\mu_{ik} - \bar{\mu}_i)}$, where $\bar{\mu}_i = \frac{1}{C}\sum_{k=1}^{C}\mu_{ik}$ and \bigcirc represents the 145

element-wise product. We can have the following theorem. 146

Theorem 1 For an attributed graph $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathbf{X})$ that follows the above assumptions in Sec. 3.2, if $|\boldsymbol{\mu}_{ii} - \boldsymbol{\mu}_{jj}| > |\boldsymbol{\mu}_{ik} - \boldsymbol{\mu}_{jk}|$ and $\boldsymbol{\sigma}_i > \boldsymbol{\sigma}_{ii}$, $\forall k \in \{1, ..., C\}$, as the decrease of homophily ratio h, the discriminability of representations obtained by the averaging process in GCN layer, i.e. $\mathbf{Z} = \mathbf{D}^{-1}\mathbf{A}\mathbf{X}$, will firstly decrease until $h = \frac{1}{C}$ then increase. When $h = \frac{1}{C}$ and $d < \frac{\boldsymbol{\sigma}_i^2}{|\boldsymbol{\mu}_{ik} - \boldsymbol{\mu}_{jk}|^2}$, the representations after averaging process will be nearly non-discriminative.

The detailed proof can be found in Appendix C. The conditions in this theorem generally hold. Since 152 the intra-class distance is often much smaller than inter-calss distance, $|\mu_{ii} - \mu_{jj}| > |\mu_{ik} - \mu_{jk}|$ is generally meet in real-world graphs. As for σ_i , it computes the standard deviations of mean neighbor 154 features in different classes. As a result, σ_i is usually much larger than the σ_{ii} and $|\mu_{ik} - \mu_{ik}|$. 155 Therefore, the Theorem 1 generally holds for the real-world graphs. And we can observe from 156 Theorem 1 that (i) heterophily level in a certain range will largely degrade the performance of 157 GCN; (ii) GCN will be more negatively affected by the heterophilic graphs with lower node degrees. 158 Though our analysis is based on GCN, it can be easily extended to GNNs that average neighbor 159 representations in the aggregation (e.g. GraphSage [14], APPNP [20], and SGC [35]). For the 160 extension of the analysis on more complex message-passing mechanism, we leave it as future work. 161

To empirically verify the above theoretical analysis, we synthe-162 size graphs with different homophily ratios and node degrees 163 by deleting/adding edges in the crocodile graph following Ap-164 pendix E.1. The results of GCN and GAT [33] on graphs with 165 various node degrees are shown in Fig. 1. We can observe that 166 (i) as the homophily ratio decreases the performance of GCN 167 will keep decreasing until h is around 0.2 ($h \approx \frac{1}{C}$), then the 168 performance will start increase;(ii) when h is around $\frac{1}{C}$, the 169 performance can be very poor and even much worse than MLP on the graph with low node degrees. The observations are in 171 consistent with our Theorem 1, which further demonstrates the



general limitations of current GNN models in learning on graphs with heterophily. This trend has also be reported in [42, 26, 25]. Moreover, theoretical analysis is conducted in [26] to prove the effectiveness of GCN on heterophilic graphs with discriminative neighborhoods. However, it can only explain the observation when $h < \frac{1}{C}$. By contrast, our theoretical analysis can well explain the whole trend of GCN performance w.r.t the homophily ratio. A similar conclusion is made with the theoretical analysis in [25], but node features are not incorporated and are replaced by label embedding vectors in their analysis.

180 3.3 Problem Definition

Based on the analysis above, we can infer that current GNNs are effective on graphs with high homophily; while they are challenged by the graphs with heterophily. In real world, we are usually given graphs with various homophily levels. In addition, the graphs are often sparsely labeled. And due to the lack of labels, the homophily ratio of the given graph is generally unknown. Thus, we aim to develop a framework that works for semi-supervised node classification on graphs with any homophily level. The problem is defined as:

Problem 1 Given an attributed graph $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathbf{X})$ with a set of labels \mathcal{Y}_L for node set $\mathcal{V}_L \subset \mathcal{V}$, the homophily ratio h of \mathcal{G} is unknown, we aim to learn a GNN which accurately predicts the labels of the unlabeled nodes, i.e., $f(\mathcal{G}, \mathcal{Y}_L) \rightarrow \hat{\mathcal{Y}}_U$, where f is the function we aim to learn and $\hat{\mathcal{Y}}_U$ is the set of predicted labels for unlabeled nodes.

191 4 Methodology

As the analysis in Sec. 3 shows, the aggregation process in GCN will mix the neighbors in various 192 labels/distributions in heterophilic graphs, resulting in non-discriminative representations for local 193 context. Based on this motivation, we propose to adopt label-wise aggregation in graph convolution, 194 i.e., neighbors in the same class are separately aggregated, to preserve the heterophilic context. Next, 195 we give the details of the label-wise aggregation along with the theoretical analysis that verifies its 196 capability in obtaining distinguishable representations for heterophilic context. Then, we present how 197 to apply label-wise graph convolution on sparsely labeled graphs and how to ensure performance on 198 both heterophilic and homophilic graphs. 199



Figure 2: The illustration of label wise aggregation and overall framework of our LW-GCN.

200 4.1 Label-Wise Graph Convolution

In heterophilic graphs, we observe that the heterophilic neighbor context itself provides useful information. Let $\mathcal{N}_k(v)$ denote node v's neighbors of label class k. As shown in Appendix F, for two nodes u and v of the same class, i.e., $y_u = y_v$, the features of nodes in $\mathcal{N}_k(u)$ are likely to be similar to that of nodes in $\mathcal{N}_k(v)$; while for nodes u and s with $y_u \neq y_s$, the features of nodes in $\mathcal{N}_k(u)$ are likely to be different from that in $\mathcal{N}_k(s)$. Therefore, for each node $v \in \mathcal{V}$, we propose to summarize the information of $\mathcal{N}_k(v)$ by label-wise aggregation to capture the useful heterophilic context. Let $\mathbf{a}_{v,k}$ be the aggregated representation of neighbors in class k, the process of obtaining representation for heterophilic context with the label-wise aggregation can be formally written as:

$$\mathbf{a}_{v,k} = \sum_{u \in \mathcal{N}_k(v)} \frac{1}{|\mathcal{N}_k(v)|} \mathbf{x}_u, \quad \mathbf{h}_v^c = \texttt{CONCAT}(\mathbf{a}_{v,1}, \dots, \mathbf{a}_{v,C}), \tag{2}$$

where *C* is the number of classes. \mathbf{h}_v^c denotes the representation of the neighborhood context. As it is shown in Eq.(2), concatenation is applied to obtain representation of context to preserve the heterophilic context. When there is no neighbor of *v* belonging to class *k*, zero embedding is assigned for class *k*. We then can augment the representation of the centered node with the context representation as the general design of GNNs. Specifically, we concatenate the context representation \mathbf{h}_v^c and centered node representation \mathbf{x}_v followed by the non-linear transformation:

$$\mathbf{h}_{v} = \sigma(\mathbf{W} \cdot \text{CONCAT}(\mathbf{x}_{v}, \mathbf{h}_{v}^{c})), \tag{3}$$

where W denotes the learnable parameters in the label-wise graph convolution and σ denotes the activation function such as ReLU.

In this section, we further prove the superiority of label-wise graph convolution in learning discrimative representations for heterophilic context by the following theorem.

Theorem 2 We consider an attributed graph $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathbf{X})$ that follows the aforementioned assump-

tions in Sec. 3.2. If $|\boldsymbol{\mu}_{ik} - \boldsymbol{\mu}_{jk}| > \sqrt{\frac{C}{d}} \boldsymbol{\sigma}_{ik}, \forall k \in \{1, \dots, C\}$, the heterophilic context representation

h^c_v that is obtained by the label-wise aggregation with Eq.(2) will keep its discriminability regardless the value of homophily ratio h.

²²³ The detailed proof is presented in Appendix D. The difference between the groups of neighbors is

naturally larger than the intra-group variance. Since $\sqrt{\frac{C}{d}}$ is usually small (e.g. around 1.8 in the

Texas graph), the condition $|\mu_{ik} - \mu_{jk}| > \sqrt{\frac{C}{d}}\sigma_{ik}$ is generally satisfied in real-world scenarios. We also adopt the label-wise graph convolution on the synthetic graphs with different homophily ratios to empirically show its effectiveness. The results can be found in Appendix E.2.

4.2 LW-GCN: A Unified Framework for Graphs with Homophily or Heterophily

Though the analysis in Sec.3.1 proves the effectiveness of label-wise graph convolution in processing graphs with heterophily, there are still two major challenges for semi-supervised node classification on graphs with any heterophily levels: (i) how to conduct label-wise graph convolution on heterophilic graphs with a small number of labeled nodes; and (ii) how to make it work for both heterophilic and homophilic graphs. To address these challenges, we propose a novel framework LW-GCN, which is illustrated in Fig. 2. LW-GCN is composed of an MLP-based pseudo label predictor f_P , a GNN f_C using label-wise graph convolution, a GNN f_G for homophilic graph, and an automatic model selection module. The predictor f_P takes the node attributes as input to give pseudo labels. f_C utilizes the estimated pseudo labels from f_P to conduct label-wise graph convolution on \mathcal{G} for node classification. To ensure the performance on graphs with any homophily level, LW-GCN also trains f_G , i.e., a GNN for homophilic graphs, and can automatically select the model for graphs with unknown homophily ratios. For the model selection module, a bi-level optimization on validation set

is applied to learn the weights for model selection. Next, we give the details of each component.

242 4.2.1 Pseudo Label Prediction.

In label-wise graph convolution, neighbors in different classes are separately aggregated to update node representations. However, only a small number of nodes are provided with labels. Thus, a pseudo label predictor f_P is deployed to estimate labels for label-wise aggregation. Specifically, a MLP is utilized to obtain pseudo label of node v as $\hat{\mathbf{y}}_v^P = \text{MLP}(\mathbf{x}_v)$, where \mathbf{x}_v is the attributes of node v. Note that, we use MLP as the predictor because message passing of the GNNs may lead to poor predictions on heterophilic graphs. The loss function for training f_P is:

$$\min_{\theta_P} \mathcal{L}_P = \frac{1}{|\mathcal{V}_{train}|} \sum_{v \in \mathcal{V}_{train}} l(\hat{y}_v^P, y_v), \tag{4}$$

where \mathcal{V}_{train} is the set of labeled nodes in the training set, y_v denote the true label of node v, θ_P represents the parameters of the predictor f_P , and $l(\cdot)$ is the cross entropy loss.

4.2.2 Architecture of LW-GCN for Heterophilic Graphs.

With f_P , we can get pseudo labels $\hat{\mathcal{Y}}_U^P$ for unlabeled nodes $\mathcal{V}_U = \mathcal{V} \setminus \mathcal{V}_L$. Combining it with the provided \mathcal{Y}_L , we have labels $\mathcal{Y}^P \in (\hat{\mathcal{Y}}_U^P \cup \mathcal{Y}_L)$ necessary for label-wise aggregation in Eq.(2). Then, node representations can be updated with the heterophilic context by Eq.(3). Multiple layers of label-wise graph convolution can be applied to incorporate more hops of neighbors in representation learning. The process of one layer label-wise graph convolution with pseudo labels can be rewritten as:

$$\mathbf{a}_{v,k}^{(l)} = \sum_{u \in \mathcal{N}_{k}^{P}(v)} \frac{1}{|\mathcal{N}_{k}^{P}(v)|} \mathbf{h}_{u}^{(l)}, \quad \mathbf{h}_{v}^{l+1} = \sigma \Big(\mathbf{W}^{(l)} \cdot \texttt{CONCAT}(\mathbf{h}_{v}^{(l)}, a_{v,1}^{(l)}, \dots, a_{v,C}^{(l)}) \Big), \tag{5}$$

where $\mathcal{N}_{k}^{P}(v) = \{u : (v, u) \in \mathcal{E} \land \hat{y}_{u}^{P} = k\}$ stands for node v's neighbors with estimated label k. $\mathbf{h}_{v}^{(l)}$ is the representation of node v at the *l*-th layer label-wise graph convolution with $\mathbf{h}_{v}^{(0)} = \mathbf{x}_{v}$. In heterophilic graphs, different hops of neighbors may exhibit different distributions which can provide useful information for node classification. Therefore, the final node prediction can be conducted by combining the intermediate representations of the model with K layers:

$$\hat{\mathbf{y}}_{v}^{C} = \operatorname{softmax}\left(\mathbf{W}_{C} \cdot \operatorname{COMBINE}(\mathbf{h}_{v}^{(1)}, ..., \mathbf{h}_{v}^{(K)})\right), \tag{6}$$

where \mathbf{W}_C is a learnable weight matrix, $\hat{\mathbf{y}}_v^C$ is predicted label probabilities of node v. Various operations such as max-pooling and concatenation [37] can be applied as the COMBINE function.

265 4.2.3 Automatic Model Selection

In heterophilic graphs, the homophily ratio is very small and even can be around 0.2 [28]. With a 266 reasonable pseudo label predictor, the label-wise aggregation with pseudo labels will mix much less 267 noise than the general GNN aggregation. In contrast, for homophilic graphs such as citation networks, 268 their homophily ratios are close to 1. In this situation, directly aggregating all the neighbors may 269 introduce less noise in representations than aggregating label-wisely as the pseudo-labels contain noises. Therefore, it is necessary to determine whether to apply the label-wise graph convolution or 271 the state-of-the-art GNN for homophilic graphs. One straightforward way is to select the model based on the homophily ratio. However, graphs are generally sparsely labeled which makes it difficult to estimate the real homophily ratio. To address this problem, we propose to utilize the validation set to 274 automatically select the model.

In the model selection module, we combine predictions of the label-wise aggregation model for heterophilic graphs and traditional GNN models for homophilic graphs. Predictions from the GNN f_G for homophilic graphs are given by $\hat{\mathcal{Y}}^G = \text{GNN}(\mathbf{A}, \mathbf{X})$, where the GNN is flexible to various models for homophilic graphs. Here, we select GCNII [8] which achieves state-of-the-art results

on homophilic graphs. The model selection can be achieved by assigning higher weight to the 280 corresponding model prediction. The combined prediction is given as: 281

$$\hat{\mathbf{y}}_{v} = \frac{\exp(\phi_{1})}{\sum_{i=1}^{2} \exp(\phi_{i})} \hat{\mathbf{y}}_{v}^{C} + \frac{\exp(\phi_{2})}{\sum_{i=1}^{2} \exp(\phi_{i})} \hat{\mathbf{y}}_{v}^{G}, \tag{7}$$

where $\hat{\mathbf{y}}_v^G \in \hat{\mathcal{Y}}^G$ is the prediction of node v from f_G . ϕ_1 and ϕ_2 are the learnable weights to control 282 the contributions of two models in final prediction. ϕ_1 and ϕ_2 can be obtained by finding the values 283 that lead to good performance on validation set. More specifically, this goal can be formulated as the 284 following bi-level optimization problem: 285

$$\min_{\phi_1,\phi_2} \mathcal{L}_{val}(\theta_C^*(\phi_1,\phi_2),\theta_G^*(\phi_1,\phi_2),\phi_1,\phi_2) \quad s.t. \ \theta_C^*,\theta_G^* = \arg\min_{\theta_C,\theta_G} \mathcal{L}_{train}(\theta_C,\theta_G,\phi_1,\phi_2)$$

where \mathcal{L}_{val} and \mathcal{L}_{train} are the average cross entropy loss of the combined predictions $\{\hat{y}_v : v \in \mathcal{V}_{val}\}$ 286 and $\{\hat{y}_v : v \in \mathcal{V}_{train}\}$ on validation set and training set, respectively. 287

4.3 An Optimization Algorithm of LW-GCN 288

Computing the gradients for ϕ_1 and ϕ_2 is expensive in both computational cost and memory. To 289 alleviate this issue, we use an alternating optimization schema to iteratively update the model 290 parameters and the model selection weights. 291

Updating Lower Level θ_C and θ_G . Instead of calculating θ_C^* and θ_C^* per outer iteration, we fix ϕ_1 292 and ϕ_2 and update the mode parameters θ_G and θ_C for T steps by: 293

$$\theta_C^{t+1} = \theta_C^t - \alpha_C \nabla_{\theta_C} \mathcal{L}_{train}(\theta_C^t, \theta_G^t, \phi_1, \phi_2), \quad \theta_G^{t+1} = \theta_G^t - \alpha_G \nabla_{\theta_G} \mathcal{L}_{train}(\theta_C^t, \theta_G^t, \phi_1, \phi_2), \quad (9)$$

where θ_C^t and θ_G^t are model parameters after updating t steps. α_C and α_G are the learning rates for θ_C and θ_G . 294 295

296

Updating Upper Level ϕ_1 and ϕ_2 . Here, we use the updated model parameters θ_C^T and θ_G^T to approximate θ_C^* and θ_G^* . Moreover, to further speed up the optimization, we apply first-order 297 approximation [13] to compute the gradients of ϕ_1 and ϕ_2 : 298

$$\phi_1^{k+1} = \phi_1^k - \alpha_{\phi} \nabla_{\phi_1} \mathcal{L}_{val}(\bar{\theta}_C^T, \bar{\theta}_G^T, \phi_1^k, \phi_2^k), \quad \phi_2^{k+1} = \phi_2^k - \alpha_{\phi} \nabla_{\phi_2} \mathcal{L}_{val}(\bar{\theta}_C^T, \bar{\theta}_G^T, \phi_1^k, \phi_2^k), \quad (10)$$

where $\bar{\theta}_{C}^{T}$ and $\bar{\theta}_{G}^{T}$ means stopping the gradient. α_{ϕ} is the learning rate for ϕ_{1} and ϕ_{2} . 299

More details of the training algorithm are in Appendix A. 300

Experiments 5 301

In this section, we conduct experiments to demonstrate the effectiveness of LW-GCN. In particular, 302 we aim to answer the following research questions: 303

• **RQ1** Is our LW-GCN effective in node classification on both homophilic and heterophilic graphs? 304

• **RQ2** Can label-wise aggregation learn representations that well capture information for prediction? 305

• **RO3** How do the quality of pseudo labels and the automatic model selection affect LW-GCN? 306

5.1 Experimental Settings 307

Datasets. For homophilic graphs, we choose the widely used benchmark datasets, Cora, Citeseer, 308 and Pubmed [19]. The dataset splits of homophilic graphs are the same as the cited paper. As for 309 heterophilic graphs, we use three webpage datasets Texas, Cornell, and Wisconsin [28], and three 310 subgraphs of wiki, i.e., Squirrel, Chameleon, and Crocodile [31]. Following [42], 10 dataset splits 311 are used in each heterophilic graph for evaluation. In addition, we also use a large scale heterophilic 312 citation network, i.e., arxiv-year [24]. 5 public splits of arxiv-year are used for evaluation. The 313 statistics of the datasets are presented in Table 3 in the Appendix. 314

Compared Methods. We compare LW-GCN with state-of-the-art GNNs, which includes GCN [19]. 315

MixHop [18], SuperGAT [17], and GCNII [8]. We also compare with the following state-of-the-art 316

models designed for heterophilic graphs: FAGCN [2], SimP-GCN [16], H2GCN [42], GRP-GNN [9], 317

BM-GCN [15], ASGC [5], LINKX [24] and GloGNN++ [23]. In addition, the MLP are evaluated on 318

the datasets for reference. The details of these compared methods can be found in Appendix B.2. 319

Dataset	Wisconsin	Texas	Chameleon	Squirrel	Crocodile	arxiv-year	Cora	Pubmed
Ave. Degree Homo. Ratio	2.05 0.20	1.69 0.11	15.85 0.24	41.74 0.22	30.96 0.25	6.9 0.22	4.01 0.81	4.50 0.8
MLP GCN MixHop SuperGAT GCNII	$\begin{array}{c} 83.5 \pm 4.9 \\ 53.1 \pm 5.8 \\ 70.2 \pm 4.8 \\ 53.7 \pm 5.7 \\ 82.1 \pm 3.9 \end{array}$	$78.1 \pm 6.0 57.6 \pm 5.9 60.6 \pm 7.7 58.6 \pm 7.7 68.6 \pm 9.8$	$\begin{array}{c} 48.0 \pm 1.5 \\ 63.5 \pm 2.5 \\ 61.2 \pm 2.2 \\ 59.4 \pm 2.5 \\ 63.5 \pm 2.5 \end{array}$	$\begin{array}{c} 32.3 \pm 1.8 \\ 46.7 \pm 1.5 \\ 44.1 \pm 1.1 \\ 38.9 \pm 1.5 \\ 49.4 \pm 1.7 \end{array}$	$\begin{array}{c} 65.8 \pm 0.7 \\ 66.7 \pm 1.0 \\ 67.6 \pm 1.3 \\ 62.6 \pm 0.9 \\ 69.0 \pm 0.7 \end{array}$	$\begin{array}{c} 36.7 \pm 0.2 \\ 46.0 \pm 0.3 \\ 46.1 \pm 0.5 \\ 38.1 \pm 0.1 \\ 47.2 \pm 0.3 \end{array}$	$58.6 \pm 0.5 \\ 81.6 \pm 0.7 \\ 80.6 \pm 0.2 \\ 82.7 \pm 0.4 \\ 84.2 \pm 0.5$	$72.7 \pm 0.4 \\78.4 \pm 1.1 \\78.9 \pm 0.5 \\78.4 \pm 0.5 \\80.2 \pm 0.2$
FAGCN SimP-GCN H2GCN GPRGNN BM-GCN ASGC LINKX GloGNN++	$\begin{array}{c} 83.3 \pm 3.7 \\ 85.5 \pm 4.7 \\ 84.7 \pm 3.9 \\ 78.2 \pm 4.4 \\ 77.6 \pm 5.9 \\ 84.3 \pm 2.6 \\ 75.5 \pm 5.7 \\ \textbf{88.0} \pm \textbf{3.2} \end{array}$	$\begin{array}{c} 79.5 \pm 4.8 \\ 80.5 \pm 5.9 \\ 83.7 \pm 6.0 \\ 77.0 \pm 6.4 \\ 81.9 \pm 5.4 \\ 85.9 \pm 4.7 \\ 74.6 \pm 8.4 \\ 83.2 \pm 4.3 \end{array}$	$\begin{array}{c} 63.9 \pm 2.2 \\ 63.7 \pm 2.3 \\ 54.2 \pm 2.3 \\ 70.6 \pm 2.1 \\ 69.4 \pm 1.7 \\ 68.8 \pm 1.6 \\ 68.4 \pm 1.4 \\ \underline{71.2 \pm 2.5} \end{array}$	$\begin{array}{c} 43.3\pm\!2.5\\ 42.8\pm\!1.4\\ 36.0\pm\!1.1\\ 50.8\pm\!1.4\\ 53.1\pm\!1.8\\ 54.5\pm\!1.6\\ \underline{61.8\pm\!1.8}\\ 57.9\pm\!2.0\end{array}$	$\begin{array}{c} 67.1 \pm 0.9 \\ 63.7 \pm 2.3 \\ 66.7 \pm 0.5 \\ 65.6 \pm 0.9 \\ 64.3 \pm 1.1 \\ 66.4 \pm 0.7 \\ \hline 79.4 \pm 0.6 \\ \hline 78.4 \pm 0.9 \end{array}$	$\begin{array}{c} 40.6\pm 0.4\\ OOM\\ 49.1\pm 0.1\\ 45.1\pm 0.2\\ OOM\\ 39.2\pm 0.1\\ \underline{56.0\pm 1.3}\\ 54.8\pm 0.3\end{array}$	$\begin{array}{c} 83.1\pm 0.6\\ 82.8\pm 0.1\\ 81.6\pm 0.4\\ 83.8\pm 0.6\\ 81.5\pm 0.5\\ 76.8\pm 0.2\\ 64.7\pm 0.4\\ 66.7\pm 1.9\end{array}$	$\begin{array}{c} 78.8 \pm 0.3 \\ \underline{80.3 \pm 0.2} \\ \overline{79.5 \pm 0.2} \\ 79.9 \pm 0.1 \\ 77.9 \pm 0.4 \\ 74.4 \pm 0.1 \\ 70.4 \pm 0.7 \\ 78.1 \pm 0.2 \end{array}$
LW-GCN Weight for f_C	$\frac{86.9\pm2.2}{0.981}$	86.2 ± 5 .8 0.960	74.4 ±1.4 0.986	62.6 ±1.6 0.987	86.5 ±0.4 0.999	56.5 ±0.2 0.942	84.3 ±0.3 0.001	80.4 ±0.3 0.006

Table 1: Node classification results (Accuracy(%) \pm Std.) on homophilic/heterophilic graphs.

Settings of LW-GCN. For the label predictor f_P , we adopt a MLP with one-hidden layer. The dimension of the hidden layer in MLP is set as 64. As for the f_C , we adopt two layers of label-wise message passing on all the datasets. More discussion about the impacts of the depth on LW-GCN is given in Sec. I. The other hyperparameters such as hidden dimension and weight decay are tuned head on the validation set. See Appendix P.1 for more details

based on the validation set. See Appendix B.1 for more details.

325 5.2 Node Classification Performance

To answer **RQ1**, we conduct experiments on both heterophilic graphs and homophilic graphs. The average accuracy and standard deviations on homophilic/heterophilic graphs are reported in Table 1. Additional results on Cornell and Citeseer datasets are presented in Appendix H. The model selection weight for label-wise aggregation GNN f_C is shown along with the results of LW-GCN. Note that this model selection weight ranges from 0 to 1. When the weight is close to 1, the label-wise aggregation model is selected. When the weight for f_C is close to 0, the GNN f_G for homophilic graph is selected.

Performance on Heterophilic Graphs. We conduct experiments on 10 dataset splits on each
 heterophilic graph. From the results on heterophilic graphs, we can have following observations:

- MLP outperforms GCN and other GNNs for homophilic graphs by a large margin on Texas and Wisconsin; while GCN can achieve relatively good performance on dense heterophilic graphs such as Chameleon. This empirical result is consistent with our analysis in Theorem 1 that the heterophily will especially degrade the performance of GCN on graphs with low degrees.
- Though GCN and other GNNs designed for homophilic graphs can give relatively good performance
 on dense heterophilic graphs, our LW-GCN bring significant improvement by adopting label-wise
 aggregation. In addition, LW-GCN outperforms baselines on heterophilic graphs with low node
 degrees. This proves the superiority of label-wise aggregation in preserving heterophilic context.
- The model selection weight for f_C is close to 1 for heterophilic graphs, which verifies that the proposed LW-GCN can correctly select the label-wise aggregation GNN f_C for heterophilic graphs.

Compared with SimP-GCN which also aims to preserve node features, our LW-GCN performs significantly better on heterophilic graphs. This is because SimP-GCN only focuses on the similarity of central node attributes. In contrast, our label-wise aggregation can preserve both the central node features and the heterophilic local context for node classification. LW-GCN also outperforms the other GNNs that adopt message-passing mechanism designed for heterophilic graphs by a large margin. This further demonstrates the effectiveness of label-wise aggregation.

Performance on Homophilic Graphs. The average results and standard deviations of 5 runs on homophilic graphs, i.e., Cora, Citeseer, and Pubmed, are also reported in Table 1 and Appendix H. From the results, we can observe that existing GNNs for heterophilic graphs generally perform worse than state-of-the-art GNNs on homophilic graphs such as GCNII. In contrast, LW-GCN achieves comparable results with the the best model on homophilic graphs. This is because LW-GCN combines the GNN using label-wise message passing and a state-of-the-art GNN for homophilic graph. And it can automatically select the right model for the given homophilic graph.

Dataset	MLP	GCN	GCNII	LW-GCN\P	LW-GCN\G	${\rm LW}\text{-}{\rm GCN}_{GCN}$	LW-GCN
Cora Citeseer	$\begin{array}{c} 58.7 \pm \! 0.5 \\ 60.3 \pm \! 0.4 \end{array}$	$81.6 \pm 0.7 \\ 71.3 \pm 0.3$	$84.2 \pm 0.5 \\ 72.0 \pm 0.8$	84.2 ± 0.3 72.3 ± 0.5	$75.3 \pm 0.4 \\ 65.1 \pm 0.5$	$81.9 \pm 0.2 \\ 71.6 \pm 0.3$	84.3 ±0.3 72.3 ±0.4
Texas Crocodile	$78.1 \pm 6.0 \\ 65.8 \pm 0.7$	$57.6 \pm 5.9 \\ 66.7 \pm 1.0$	$68.6 \pm 9.8 \\ 69.0 \pm 0.7$	$82.4 \pm 5.2 \\ 84.6 \pm 2.4$	$85.9 \pm 5.6 \\ 85.8 \pm 0.9$	85.4 ±6.3 84.7 ±0.9	86.2 ±5.8 86.5 ±0.5

 Table 2: Ablation Study

357 5.3 Analysis of Node Representations

To answer **RO2**, we compare the representation similarity 358 of intra-class node pairs and inter-class node pairs on a 359 sparse heterophilic graph in Fig. 3. For both GCN and 360 LW-GCN, representations learned by the last layer are 361 used for analysis. we can observe that the learned rep-362 363 resentations of GCN are very similar for both intra-class pairs and inter-class pairs. This verifies that simply aggre-364 gating the neighbors will make the node representations 365 less discrimative. With label-wise aggregation, the sim-366 ilarity scores of intra-class pairs are significantly higher 367



Figure 3: Representation similarity distributions on Texas Graphs.

than inter-class node pairs. This demonstrates that the
 representations learned by label-wise message passing can well preserve the target nodes' features
 and their contextual information.

371 5.4 Ablation Study

To answer **RO3**, we conduct ablation studies to understand the contributions of each component to 372 LW-GCN. To investigate how the quality of pseudo labels can affect LW-GCN, we train a variant LW-GCN\P by replacing the MLP-based label predictor with a GCN model. To show the importance of 374 the automatic model selection, we train a variant LW-GCN\G which removes the GNN for homophilic 375 graphs and only uses label-wise aggregation GNN. Finally, we replace the GCNII backbone of f_G 376 to GCN, denoted as LW-GCN_{GCN}, to show LW-GCN is flexible to adopt various GNNs for f_G . 377 Experiments are conducted on both homophilic and heterophilic graphs. The results are shown in 378 Table 2 and ablation studies on the rest datasets are shown in Appendix G. We can observe that: 379 • On homophilic graphs, LW-GCN\P shows comparable results with LW-GCN, because GCNII 380 will be selected given a homophilic graph. On the heterophilic graph Texas, the performance of 381

LW-GCN\P is significantly worse than LW-GCN. This is because GNNs can produce poor pseudo labels on heterophilic graph, which degrades the label-wise message passing.

• LW-GCN\G performs much better than MLP. This shows label-wise graph convolution can capture structure information. However, LW-GCN\G performs worse than GCNII and LW-GCN on homophilic graphs, which indicates the necessity of combining GNN for homophilic graphs.

• LW-GCN_{GCN} achieves comparable results with GCN on homophilic graphs. On heterophilic graphs, LW-GCN_{GCN} performs similarly with LW-GCN. This shows the flexibility of LW-GCN in adopting traditional GNN models designed for homophilic graphs.

6 Conclusion and Future Work

In this paper, we analyze the impacts of the heterophily levels to GCN model and demonstrate its 391 limitations. We develop a novel label-wise graph convolution to learn representations that preserve the 392 node features and their heterophilic neighbors' information. An automatic model selection module 393 is applied to ensure the performance of the proposed framework on graphs with any homophily ratio. Theoretical and empirical analysis demonstrates the effectiveness of the label-wise aggregation. 395 Extensive experiments shows that our proposed LW-GCN can achieve sate-of-the-art results on both 396 homophilic and heterophilic graphs. There are several interesting directions need further investigation. 397 First, since better pseudo labels will benefit the label-wise message passing, it is promising to 398 incorporate the predictions of LW-GCN in label-wise message passing. Second, in some applications 399 such as link prediction, labels are not available. Therefore, we will investigate how to generate useful 400 pseudo labels for label-wise aggregation for applications where no labeled nodes are provided. 401

402 **References**

- [1] Sami Abu-El-Haija, Bryan Perozzi, Amol Kapoor, Nazanin Alipourfard, Kristina Lerman, Hrayr
 Harutyunyan, Greg Ver Steeg, and Aram Galstyan. Mixhop: Higher-order graph convolutional
 architectures via sparsified neighborhood mixing. In *international conference on machine learning*, pages 21–29. PMLR, 2019. 2, 14
- [2] Deyu Bo, Xiao Wang, Chuan Shi, and Huawei Shen. Beyond low-frequency information in graph convolutional networks. *arXiv preprint arXiv:2101.00797*, 2021. 2, 7, 14
- [3] Pietro Bongini, Monica Bianchini, and Franco Scarselli. Molecular generative graph neural
 networks for drug discovery. *Neurocomputing*, 450:242–252, 2021. 1
- [4] Joan Bruna, Wojciech Zaremba, Arthur Szlam, and Yann LeCun. Spectral networks and locally
 connected networks on graphs. *ICLR*, 2014. 2
- [5] Sudhanshu Chanpuriya and Cameron Musco. Simplified graph convolution with heterophily.
 arXiv preprint arXiv:2202.04139, 2022. 7, 14
- [6] Jie Chen, Shouzhen Chen, Zengfeng Huang, Junping Zhang, and Jian Pu. Exploiting neighbor effect: Conv-agnostic gnns framework for graphs with heterophily. *arXiv preprint arXiv:2203.11200*, 2022. 1
- [7] Jie Chen, Tengfei Ma, and Cao Xiao. Fastgen: fast learning with graph convolutional networks
 via importance sampling. *ICLR*, 2018. 2
- [8] Ming Chen, Zhewei Wei, Zengfeng Huang, Bolin Ding, and Yaliang Li. Simple and deep graph convolutional networks. In *ICML*, pages 1725–1735. PMLR, 2020. 2, 6, 7, 13, 14
- [9] Eli Chien, Jianhao Peng, Pan Li, and Olgica Milenkovic. Adaptive universal generalized
 pagerank graph neural network. *arXiv preprint arXiv:2006.07988*, 2020. 1, 2, 7, 14
- [10] Enyan Dai and Suhang Wang. Say no to the discrimination: Learning fair graph neural networks
 with limited sensitive attribute information. In WSDM, pages 680–688, 2021. 2
- [11] Michaël Defferrard, Xavier Bresson, and Pierre Vandergheynst. Convolutional neural networks
 on graphs with fast localized spectral filtering. In *NeurIPS*, pages 3844–3852, 2016. 2
- [12] Yingtong Dou, Zhiwei Liu, Li Sun, Yutong Deng, Hao Peng, and Philip S Yu. Enhancing graph
 neural network-based fraud detectors against camouflaged fraudsters. In *Proceedings of the* 29th ACM International Conference on Information & Knowledge Management, pages 315–324,
 2020. 2
- [13] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adapta tion of deep networks. In *International Conference on Machine Learning*, pages 1126–1135.
 PMLR, 2017. 7
- [14] Will Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large
 graphs. In *NeurIPS*, pages 1024–1034, 2017. 1, 2, 4
- [15] Dongxiao He, Chundong Liang, Huixin Liu, Mingxiang Wen, Pengfei Jiao, and Zhiyong Feng.
 Block modeling-guided graph convolutional neural networks. *AAAI*, 2022. 2, 7, 14
- [16] Wei Jin, Tyler Derr, Yiqi Wang, Yao Ma, Zitao Liu, and Jiliang Tang. Node similarity preserving
 graph convolutional networks. In *Proceedings of the 14th ACM International Conference on Web Search and Data Mining*, pages 148–156, 2021. 2, 7, 14
- [17] Dongkwan Kim and Alice Oh. How to find your friendly neighborhood: Graph attention design with self-supervision. In *International Conference on Learning Representations*, 2020. 7, 14
- [18] Dongkwan Kim and Alice Oh. How to find your friendly neighborhood: Graph attention design
 with self-supervision. In *International Conference on Learning Representations*, 2021. 2, 7
- [19] Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional
 networks. *arXiv preprint arXiv:1609.02907*, 2016. 2, 3, 7, 14
- [20] Johannes Klicpera, Aleksandar Bojchevski, and Stephan Günnemann. Predict then propagate:
 Graph neural networks meet personalized pagerank. *arXiv preprint arXiv:1810.05997*, 2018. 2,
 4
- [21] Ron Levie, Federico Monti, Xavier Bresson, and Michael M Bronstein. Cayleynets: Graph
 convolutional neural networks with complex rational spectral filters. *IEEE Transactions on Signal Processing*, 67(1):97–109, 2018. 2

- 454 [22] Guohao Li, Matthias Muller, Ali Thabet, and Bernard Ghanem. Deepgcns: Can gcns go as deep
 455 as cnns? In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages
 456 9267–9276, 2019. 2
- [23] Xiang Li, Renyu Zhu, Yao Cheng, Caihua Shan, Siqiang Luo, Dongsheng Li, and Weining Qian.
 Finding global homophily in graph neural networks when meeting heterophily. *arXiv preprint arXiv:2205.07308*, 2022. 2, 7, 14
- [24] Derek Lim, Felix Hohne, Xiuyu Li, Sijia Linda Huang, Vaishnavi Gupta, Omkar Bhalerao, and
 Ser Nam Lim. Large scale learning on non-homophilous graphs: New benchmarks and strong
 simple methods. Advances in Neural Information Processing Systems, 34:20887–20902, 2021.
 2, 7, 14
- [25] Sitao Luan, Chenqing Hua, Qincheng Lu, Jiaqi Zhu, Mingde Zhao, Shuyuan Zhang, Xiao-Wen
 Chang, and Doina Precup. Is heterophily a real nightmare for graph neural networks to do node
 classification? *arXiv preprint arXiv:2109.05641*, 2021. 4, 21
- Yao Ma, Xiaorui Liu, Neil Shah, and Jiliang Tang. Is homophily a necessity for graph neural networks? *arXiv preprint arXiv:2106.06134*, 2021. 1, 4
- [27] Mathias Niepert, Mohamed Ahmed, and Konstantin Kutzkov. Learning convolutional neural
 networks for graphs. In *ICML*, pages 2014–2023, 2016. 2
- [28] Hongbin Pei, Bingzhe Wei, Kevin Chen-Chuan Chang, Yu Lei, and Bo Yang. Geom-gcn:
 Geometric graph convolutional networks. *arXiv preprint arXiv:2002.05287*, 2020. 2, 6, 7
- [29] Oleg Platonov, Denis Kuznedelev, Artem Babenko, and Liudmila Prokhorenkova. Characteriz ing graph datasets for node classification: Beyond homophily-heterophily dichotomy. *arXiv preprint arXiv:2209.06177*, 2022. 21
- [30] Jiezhong Qiu, Qibin Chen, Yuxiao Dong, Jing Zhang, Hongxia Yang, Ming Ding, Kuansan
 Wang, and Jie Tang. Gcc: Graph contrastive coding for graph neural network pre-training. In
 Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, pages 1150–1160, 2020. 2
- [31] Benedek Rozemberczki, Carl Allen, and Rik Sarkar. Multi-scale attributed node embedding.
 Journal of Complex Networks, 9(2):cnab014, 2021. 7
- [32] Ke Sun, Zhanxing Zhu, and Zhouchen Lin. Multi-stage self-supervised learning for graph
 convolutional networks. AAAI, 2020. 2
- [33] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua
 Bengio. Graph attention networks. *ICLR*, 2018. 2, 4, 14
- [34] Daixin Wang, Jianbin Lin, Peng Cui, Quanhui Jia, Zhen Wang, Yanming Fang, Quan Yu, Jun
 Zhou, Shuang Yang, and Yuan Qi. A semi-supervised graph attentive network for financial
 fraud detection. *ICDM*, 2019. 2
- [35] Felix Wu, Amauri Souza, Tianyi Zhang, Christopher Fifty, Tao Yu, and Kilian Weinberger.
 Simplifying graph convolutional networks. In *International conference on machine learning*,
 pages 6861–6871. PMLR, 2019. 4, 14
- [36] Zonghan Wu, Shirui Pan, Fengwen Chen, Guodong Long, Chengqi Zhang, and S Yu Philip. A
 comprehensive survey on graph neural networks. *IEEE transactions on neural networks and learning systems*, 2020. 1
- [37] Keyulu Xu, Chengtao Li, Yonglong Tian, Tomohiro Sonobe, Ken-ichi Kawarabayashi, and
 Stefanie Jegelka. Representation learning on graphs with jumping knowledge networks. In
 International Conference on Machine Learning, pages 5453–5462. PMLR, 2018. 2, 6
- [38] Yuning You, Tianlong Chen, Yongduo Sui, Ting Chen, Zhangyang Wang, and Yang Shen.
 Graph contrastive learning with augmentations. *Advances in Neural Information Processing Systems*, 33, 2020. 2
- [39] Bing Yu, Haoteng Yin, and Zhanxing Zhu. Spatio-temporal graph convolutional networks: A deep learning framework for traffic forecasting. *arXiv preprint arXiv:1709.04875*, 2017. 1, 2
- [40] Tianxiang Zhao, Xianfeng Tang, Xiang Zhang, and Suhang Wang. Semi-supervised graph-tograph translation. In *CIKM*, pages 1863–1872, 2020. 2

[41] Jiong Zhu, Ryan A Rossi, Anup Rao, Tung Mai, Nedim Lipka, Nesreen K Ahmed, and Danai
 Koutra. Graph neural networks with heterophily. *arXiv preprint arXiv:2009.13566*, pages

507 11168–11176, 2020. 2

- [42] Jiong Zhu, Yujun Yan, Lingxiao Zhao, Mark Heimann, Leman Akoglu, and Danai Koutra.
 Beyond homophily in graph neural networks: Current limitations and effective designs. *arXiv* preprint arXiv:2006.11468, 2020. 1, 2, 3, 4, 7, 14, 21
- [43] Qikui Zhu, Bo Du, and Pingkun Yan. Self-supervised training of graph convolutional networks.
 arXiv preprint arXiv:2006.02380, 2020. 2

Algorithm 1 Training Algorithm of LW-GCN

Input: $\mathcal{G} = (\mathcal{V}, \mathcal{E}, X), \mathcal{Y}_L, p, \alpha_C, \alpha_G, \alpha_\phi$ and T **Output:** f_P , f_C , f_G , ϕ_1 and ϕ_2 1: Train f_P by optimizing Eq.(4) w.r.t θ_P 2: Obtain pseudo labels $\hat{\mathcal{Y}}^P$ with f_P 3: repeat Get combined predictions of f_C and f_G on \mathcal{V}_{val} 4: 5: Calculate the upper level loss \mathcal{L}_{val} Update ϕ_1 and ϕ_2 according to Eq.(10) 6: 7: for t = 1 to T do Obtain the lower level loss \mathcal{L}_{train} 8: 9: Update θ_C and θ_G by Eq.(9) 10: end for

11: **until** convergence

Dataset	Nodes	Edges	Classes	Hom. Ratio
Wisconsin	251	515	5	0.20
Texas	183	309	5	0.11
Cornell	183	280	5	0.30
Chameleon	2,277	36,101	5	0.24
Squirrel	5,201	217,073	5	0.22
Crocodile	11,631	360,040	5	0.25
arxiv-year	169,343	1,166,243	5	0.22
Cora	2,708	5,429	6	0.81
Citeseer	3,327	4,732	7	0.74
Pubmed	19,717	44,338	3	0.8

513 A Training Algorithm of LW-GCN

The training algorithm of LW-GCN is shown in Algorithm 1. In line 1 and 2, we firstly train the f_P to obtain the required pseudo labels for label-wise message passing. From line 4 to 6, we get the combined predictions from f_C and f_G and update the model selection weights with Eq.(10). From line 7 to 10, we update the model parameters θ_C and θ_G by minimizing \mathcal{L}_{train} with Eq.(9). The updating of model selection weights and model parameters are conducted iteratively until convenience.

519 **B** Additional Details of Experimental Settings

520 B.1 Implementation Details of LW-GCN

For experiments on each heterophilic graph, we report the results on the 10 public dataset splits. 521 For homophily graphs, we run each experiment 5 times on the provided public dataset split. The 522 hidden dimension of f_P is fixed as 64 for all graphs. For the f_C on Texas and Wisconsin, a linear 523 layer is firstly applied to transform the features followed by the label-wise graph convolutional layer. 524 As for the other graphs, the label-wise graph convolutional layer is directly applied to the node 525 features. The hidden layer dimension and weight decay rate are tuned based on the validation set 526 by grid search. Specifically, we vary the hidden dimension and weight decay in $\{32, 64, 128, 256\}$ 527 and $\{0.05, 0.005, 0.0005, 0.00005\}$, respectively. As for the f_G which deploys GCNII [8] as the 528 backbone, the hyperparameter settings are the same as the cited paper. During the training phase, the 529 learning rate is set as 0.01 for all the parameters and model selection weights. The inner iteration step 530 T is set as 1. Our machine uses an Intel i7-9700k CPU with 64GB RAM. A Nvidia 2080Ti GPU is 531 532 used to run all the experiments.

533 B.2 Implementation Details of Compared Methods

We adopt a two-layer MLP model on the datasets as baslines to show the effects of the graph structure and local context of the graphs. The hidden dimension is set the same as our LW-GCN. Apart

- from MLP, we compare LW-GCN with the following representative and state-of-the-art GNNs that
- originally designed for graphs with homophily:
- **GCN** [19]: This is a popular spectral-based Graph Convolutional Network, which aggregates the neighbor information and the centered node by averaging their representations. We apply the official code in https://github.com/tkipf/pygcn.
- **MixHop** [1]: It adopts a graph convolutional layer with powers of the adjacency matrix. The official code in https://github.com/samihaija/mixhop is implemented for comparsion.
- **SuperGAT** [17]: This is a GAT model augmented by the self-supervision. In SuperGAT, apart from the classification loss on provided labels, a self-supervised learning task is deployed to further guide the learning of attention for better information propagation based on GAT [33]. The official code from the authors in https://github.com/dongkwan-kim/SuperGAT is used.
- **GCNII** [8]: Based on GCN, residual connection and identity mapping are applied in GCNII to have a deep GNN for better performance. The experiments are run with the official implementation in https://github.com/chennnM/GCNII.
- ⁵⁵⁰ We also compare LW-GCN with the following baseline GNN models for heterophilic graphs:
- **FAGCN** [2]: FAGCN adaptively aggregates low-frequency and high-frequency signals from neighbors to improve the performance on heterophilic graphs. The implementation from authors in https://github.com/bdy9527/FAGCN is applied in our experiments.
- SimP-GCN [16]: A feature similarity preserving aggregation is applied to facilitate the representation learning on graphs with homophily and heterophily. We utilize the official code in https://github.com/ChandlerBang/SimP-GCN.
- **H2GCN** [42]: H2GCN investigates the limitations of GCN on graphs with heterophily. And it accordingly adopts three key designs for node classification on heterophilic graphs. We conduct experiments with the official code from authors in https://github.com/GemsLab/H2GCN.
- **GPR-GNN** [9]: This method introduces a new Generalized PageRank (GPR) GNN to adaptively learn the GPR weights that combine the aggregated representations in different orders. The learned GPR weights can be either positive or negative, which allows the GPR-GNN handle both heterophilic and homophilic graphs. We adopt the official code from authors in https: //github.com/jianhao2016/GPRGNN.
- **BM-GCN** [15]: This is one of the most recent methods designed for graphs with heterophily, which achieves state-of-the-art results on heterophilic graphs. A block-modeling is adopted to GCN to aggregate information from homophilic and heterophilic neighbors discrimatively. More specifically, the link between two nodes will be re-weighted based on the soft labels of two nodes and the block-similarity matrix. The training and evaluation process is based on the official code in https://github.com/hedongxiao-tju/BM-GCN.
- **ASGC** [5]: This method replaces the fixed feature propagation step of SGC [35] with an adaptive propagation, which can be effective for both homophilic graphs and heterophilic graphs. We use the official code released in https://openreview.net/forum?id=jRrpiqxtrWm.
- LINKX [24]: This methods separately embed the adjacency matrix and node features with multilayer perceptrions and simple transformations. We use the official code from authors in https://github.com/CUAI/Non-Homophily-Large-Scale.
- **GloGNN++** [23]: This method will learn a coefficient matrix to capture the correlations between nodes to aggregate information from global nodes in the graph. The values of the coefficient matrix can be signed and are derived from the optimization. In our experiments, we use the official code in https://github.com/recklessronan/glognn.
- The model architecture and hyperparameters of the baselines are set according to the experimental settings provided by the authors for reproduction. For datasets that are not given reproduction details, the hyperparameters of baselines will be tuned based on the performance on validation set to make a fair comparison.

585 C Proof of Theorem 1

Proof 1 In this proof, we focus on nodes in class i and class j, where $i \neq j$. Since dimensions of the node feature are independent to each other, without loss of generality, we consider one dimension of the feature and aggregated representation for node v, which is denoted as x_v and z_v . For node v in class *i*, the aggregated representation z_v in GCN layer is rewritten as:

$$z_v = \sum_{u \in \mathcal{N}(v)} \frac{1}{|\mathcal{N}(v)|} x_u.$$
(11)

With assumptions in Sec. 3.2, the expectation of aggregated representations of nodes in class *i* can be written as:

$$\mathbb{E}(z_v|y_v = i) = h \cdot \mu_{ii} + \frac{1-h}{C-1} \sum_{k=1, k \neq i}^C \mu_{ik},$$
(12)

Similarly, we can get the expectation of aggregated nodes representations in class j, i.e., $\mathbb{E}(z_v|y_v=j)$. Then, the difference between $\mathbb{E}(z_v|y_v=i)$ and $\mathbb{E}(z_v|y_v=j)$ is

$$\Delta_{i,j} = |\mathbb{E}(z_v | y_v = i) - \mathbb{E}(z_v | y_v = j)|$$

$$= |h \cdot (\mu_{ii} - \mu_{jj}) + \frac{1 - h}{C - 1} (\mu_{ij} - \mu_{ji}) + \frac{1 - h}{C - 1} \sum_{k=1, k \neq i, j}^C (\mu_{ik} - \mu_{jk})|$$

$$= |\frac{hC - 1}{C - 1} (\mu_{ii} - \mu_{jj}) + \frac{1 - h}{C - 1} (\sum_{k=1}^C (\mu_{ik} - \mu_{jk}))|$$
(13)

⁵⁹⁴ We firstly consider the situation of $h \ge \frac{1}{C}$. When $h \ge \frac{1}{C}$, we can infer the upper bound of $\Delta_{i,j}$ as:

$$\Delta_{i,j} \leq \frac{hC-1}{C-1} |\mu_{ii} - \mu_{jj}| + \frac{1-h}{C-1} \sum_{k=1}^{C} |\mu_{ik} - \mu_{jk}| = \frac{hC}{C-1} (|\mu_{ii} - \mu_{jj}| - \frac{1}{C} \sum_{k=1}^{C} |\mu_{ik} - \mu_{jk}|) + \frac{1}{C-1} (\sum_{k=1}^{C} |\mu_{ik} - \mu_{jk}| - |\mu_{ii} - \mu_{jj}|),$$
(14)

595 And the lower bound of $\Delta_{i,j}$ is:

$$\Delta_{i,j} \ge \frac{hC-1}{C-1} |\mu_{ii} - \mu_{jj}| - \frac{1-h}{C-1} \sum_{k=1}^{C} |\mu_{ik} - \mu_{jk}| = \frac{hC}{C-1} (|\mu_{ii} - \mu_{jj}| + \frac{1}{C} \sum_{k=1}^{C} |\mu_{ik} - \mu_{jk}|) - \frac{1}{C-1} (\sum_{k=1}^{C} |\mu_{ik} - \mu_{jk}| + |\mu_{ii} - \mu_{jj}|),$$
(15)

Thus, when $|\mu_{ii} - \mu_{jj}| > |\mu_{ik} - \mu_{jk}|, \forall k \in \{1, ...C\}$ and $h \ge \frac{1}{C}$, both the upper bound and lower bound of $\Delta_{i,j}$ will decrease with the decrease of h.

Next, we will show that lower h under the condition of $h \ge \frac{1}{C}$ will lead to higher variance of aggregated nodes. According to Eq.(11), the variance of $\{z_v : y_v = i\}$ can be written as:

$$Var(z_v|y_v=i) = Var(\sum_{u \in \mathcal{N}(v)} \frac{1}{|\mathcal{N}(v)|} x_u|y_v=i)$$

According to the assumption 1, the neighbor features are conditional independent to each other given the label of the center node. And for each neighbor node $u \in \mathcal{N}(v)$, we have $P(y_u = y_v|y_v) = h$, $P(y_u = y|y_v) = \frac{1-h}{C-1}$, $\forall y \neq y_v$. Therefore, for neighbor node $u \in \mathcal{N}(v)$ of node v whose label is *i*, its features follow a mixed distribution:

$$P(x_{u}|y_{v} = i)$$

$$= \sum_{k=1}^{C} P(y_{u} = k|y_{v} = i)P(x_{u}|y_{u} = k)$$

$$= h \cdot N(\mu_{ii}, \sigma_{ii}) + \frac{1-h}{C-1} \sum_{k=1, k \neq i} N(\mu_{ik}, \sigma_{ik})$$
(16)

 U_{504} Using the variance of mixture distribution, the variance of node v in class i can be derived as

$$Var(z_{v}|y_{v}=i) = \frac{1}{d}Var(x_{u}|y_{v}=i)$$

$$= \frac{1}{d}(\mathbb{E}[Var(x_{u}|y_{u}, y_{v}=i)] + Var[\mathbb{E}(x_{u}|y_{u}, y_{v}=i)])$$

$$= \frac{1}{d}\left(h\sigma_{ii}^{2} + \frac{1-h}{C-1}\sum_{k=1,k\neq i}^{C}\sigma_{ik}^{2} + h\mu_{ii}^{2} + \frac{1-h}{C-1}\sum_{k=1,k\neq i}^{C}\mu_{ik}^{2} - (h\mu_{ii} + \frac{1-h}{C-1}\sum_{k=1,k\neq i}^{C}\mu_{ik})^{2}\right)$$
(17)

Let $\bar{\mu}_i = \frac{1}{C} \sum_{k=1}^C \mu_{ik}$ and $\sigma_i^2 = \frac{1}{C} \sum_{k=1}^C (\mu_{ik} - \bar{\mu}_i)^2$. Then Eq.(17) can be rewritten as the following equation: $Var(z_v | y_v = i)$

$$= \frac{1}{d} \left(\frac{hC - 1}{C - 1} \sigma_{ii}^{2} + \frac{C - hC}{C - 1} \left(\frac{1}{C} \sum_{k=1}^{C} \sigma_{ik}^{2} + \sigma_{i}^{2} \right) + \frac{hC - 1}{C - 1} \mu_{ii}^{2} + \frac{C - hC}{C - 1} \bar{\mu}_{i}^{2} - \left(h\mu_{ii} + \frac{1 - h}{C - 1} \sum_{k=1, k \neq i}^{C} \mu_{ik} \right)^{2} \right)$$
(18)

As $h \ge \frac{1}{C}$, we can set $p = \frac{hC-1}{C-1}$, $0 \le p \le 1$ and $\frac{C-hC}{C-1} = 1 - p$. For the last three terms of Eq.(18), we have:

$$\frac{hC-1}{C-1}\mu_{ii}^2 + \frac{C-hC}{C-1}\bar{\mu}_i^2 - (h\mu_{ii} + \frac{1-h}{C-1}\sum_{k=1,k\neq i}^C \mu_{ik})^2$$

$$= p\mu_{ii}^2 + (1-p)\bar{\mu}_i^2 - (p\mu_{ii} + (1-p)\bar{\mu}_i)^2$$

$$= p(1-p)(\mu_{ii} - \bar{\mu}_i)^2 \ge 0$$
(19)

Combining Eq.(18) and Eq.(19), we are able to get the lower bound of the variance as: $Var(z_v|y_v = i)$

$$\geq \frac{hC-1}{d(C-1)}\sigma_{ii}^{2} + \frac{C-hC}{d(C-1)}\left(\frac{1}{C}\sum_{k=1}^{C}\sigma_{ik}^{2} + \sigma_{i}^{2}\right)$$

$$= \frac{hC}{d(C-1)}\left(\sigma_{ii}^{2} - \sigma_{i}^{2} - \frac{1}{C}\sum_{k=1}^{C}\sigma_{ik}^{2}\right) + \frac{1}{d(C-1)}\left(C\sigma_{i}^{2} + \sum_{k=1}^{C}\sigma_{ik}^{2} - \sigma_{ii}^{2}\right)$$
(20)

When $\sigma_i > \sigma_{ii}$, we know that with the decrease of h, the lower bound of $Var(z_v|y_v = i)$ will increase. Similarly, $Var(z_v|y_v = j)$ will also increase with a lower h. Combining with $|\mathbb{E}(z_v|y_v = i) - \mathbb{E}(z_v|y_v = j)|$ will decrease with the decrease of h, we can conclude that when $h \ge \frac{1}{C}$, the graph with lower h will lead to less discrimative aggregate representations.

We then prove when $h < \frac{1}{C}$, the decreasing of h will increase the discriminability of the aggregated representations by averaging. Specifically, with Eq.(13), we can infer that when $h < \frac{1}{C}$ the upper bound of $\Delta_{i,j}$ will be:

$$\Delta_{i,j} \leq \frac{1-hC}{C-1} |\mu_{ii} - \mu_{jj}| + \frac{1-h}{C-1} \sum_{k=1}^{C} |\mu_{ik} - \mu_{jk}| = \frac{-hC}{C-1} (|\mu_{ii} - \mu_{jj}| + \frac{1}{C} \sum_{k=1}^{C} |\mu_{ik} - \mu_{jk}|) + \frac{1}{C-1} (\sum_{k=1}^{C} |\mu_{ik} - \mu_{jk}| + |\mu_{ii} - \mu_{jj}|),$$
(21)

617 And the lower bound of $\Delta_{i,j}$ is:

$$\Delta_{i,j} \ge \frac{1-hC}{C-1} |\mu_{ii} - \mu_{jj}| - \frac{1-h}{C-1} \sum_{k=1}^{C} |\mu_{ik} - \mu_{jk}| = \frac{-hC}{C-1} (|\mu_{ii} - \mu_{jj}| - \frac{1}{C} \sum_{k=1}^{C} |\mu_{ik} - \mu_{jk}|) - \frac{1}{C-1} (\sum_{k=1}^{C} |\mu_{ik} - \mu_{jk}| - |\mu_{ii} - \mu_{jj}|),$$
(22)

Thus, when $h < \frac{1}{C}$ and $|\mu_{ii} - \mu_{jj}| > |\mu_{ik} - \mu_{jk}|, \forall k \in \{1, ...C\}$, both the upper bound and lower bound of $\Delta_{i,j}$ will increase with the decrease of h.

For the variance of aggregated representations when $h < \frac{1}{C}$, we can infer its following upper bound with Eq.(18):

$$Var(z_{v}|y_{v} = i)$$

$$\leq \frac{hC - 1}{d(C - 1)}\sigma_{ii}^{2} + \frac{C - hC}{d(C - 1)}(\frac{1}{C}\sum_{k=1}^{C}\sigma_{ik}^{2} + \sigma_{i}^{2})$$

$$= \frac{hC}{d(C - 1)}(\sigma_{ii}^{2} - \sigma_{i}^{2} - \frac{1}{C}\sum_{k=1}^{C}\sigma_{ik}^{2}) + \frac{1}{d(C - 1)}(C\sigma_{i}^{2} + \sum_{k=1}^{C}\sigma_{ik}^{2} - \sigma_{ii}^{2})$$
(23)

According to the assumption that $\sigma_i > \sigma_{ii}$, we know that with the decrease of h under the condition of $h < \frac{1}{C}$ the upper bound of the $Var(z_v|y_v = i)$ will decrease. We can have the same conclusion for $Var(z_v|y_v = j)$. Combining the trend that when $h < \frac{1}{C} |\mathbb{E}(z_v|y_v = i) - \mathbb{E}(z_v|y_v = j)|$ will increase with the decrease of h, we can conclude that when $h < \frac{1}{C}$, the graph with lower h will have more discriminative aggregate representations.

627 When $h = \frac{1}{C}$, we can get

$$\Delta_{i,j} = \frac{1}{C} |\sum_{k=1}^{C} (\mu_{ik} - \mu_{jk})|, \qquad (24)$$

628

$$Var(z_{v}|y_{v}=i) \ge \frac{1}{d} (\frac{1}{C} \sum_{k=1}^{C} \sigma_{ik}^{2} + \sigma_{i}^{2})|,$$
(25)

If $\sigma_i > \sqrt{d}|\mu_{ik} - \mu_{ik}|, \forall k \in \{1, \dots, C\}$, we can get $Var(z_v|y_v = i) > \Delta_{i,j}^2$. So when $h = \frac{1}{C}$ and $\sigma_i > \sqrt{d}|\mu_{ik} - \mu_{ik}|, \forall k \in \{1, \dots, C\}$, the representations after the averaging process will be non-discrimative.

632 **D** Proof of Theorem 2

Proof 2 In this proof, we also consider a center node v in class i. And we focus on one dimension of the node feature and aggregated representation. Specifically, for each dimension, the label-wise aggregation can be written as:

$$a_{v,k} = \sum_{u \in \mathcal{N}_k(v)} \frac{1}{|\mathcal{N}_k(v)|} x_u, \tag{26}$$

where $a_{v,k}$ denotes the aggregated feature of neighbors in class k. Since $u \in \mathcal{N}_k(v)$, we know node

$$u$$
's features x_u follows distribution as $x_u \sim N(\mu_{ik}, \sigma_{ik})$. The mean of $a_{v,k}$ in Eq.(26) is given as

$$\mathbb{E}(a_{v,k}|y_v=i) = \mu_{ik}.$$
(27)

Then the absolute difference between $\mathbb{E}(a_{v,k}|y_v = i)$ and $\mathbb{E}(a_{v,k}|y_v = j)$ will be:

$$\Delta_{i,j}^{k} = |\mathbb{E}(a_{v,k}|y_{v}=i) - \mathbb{E}(a_{v,k}|y_{v}=j)| = |\mu_{ik} - \mu_{jk}|.$$
(28)

Given the assumption that the features are conditionally independent given the label of center node,

640 the variance of $a_{v,k}$ can be written as:

$$Var(a_{v,k}|y_v=i) = \begin{cases} \frac{1}{dh}\sigma_{ik}^2 & \text{if } k=i;\\ \frac{C-1}{d(1-h)}\sigma_{ik}^2 & \text{else,} \end{cases}$$
(29)

In label-wise aggregation, we generally concatenate the $\{a_{v,k} : k \in \{1, ..., C\}\}$ for further classifi-

cation. Therefore, the lower bound of discriminability can be given by the representation of the class

643 that are most discriminative, which can be formally written as:

$$k^* = \arg\max_k \frac{(\Delta_{i,j}^k)^2}{Var(a_{v,k}|y_v=i)}$$
(30)

When $h \geq \frac{1}{C}$, we can get: 644

$$\frac{(\Delta_{i,j}^{k^*})^2}{Var(a_{v,k^*}|y_v=i)} \ge \frac{dh|\mu_{ii} - \mu_{ji}|^2}{\sigma_{ii}^2} \ge \frac{d|\mu_{ii} - \mu_{ji}|^2}{C\sigma_{ii}^2}$$
(31)

As for $h \leq \frac{1}{C}$, let $k \neq i$ we can infer that:

$$\frac{(\Delta_{i,j}^{k^*})^2}{Var(a_{v,k^*}|y_v=i)} \ge \frac{d(1-h)|\mu_{ik}-\mu_{jk}|^2}{(C-1)\sigma_{ik}^2} \ge \frac{d|\mu_{ik}-\mu_{jk}|^2}{C\sigma_{ik}^2}$$
(32)

Therefore, if the condition that $|\mu_{ik} - \mu_{jk}| > \sqrt{\frac{C}{d}}\sigma_{ik}, \forall k \in \{1, \dots, C\}$ is met, we can infer from 646

Eq.(31) and Eq.(32) that $\frac{(\Delta_{i,j}^{k^*})^2}{Var(a_{v,k^*}|y_v=i)} > 1$ regardless the value of the homophily ratio h. This 647

shows that label-wise aggregation can preserve the context and ensure the high discriminability 648

649 regardless the homophily ratio.

Additional Details and Experiments on Generated Graphs Е 650

Algorithm 2 Algorithm of Generating Graphs

Input: $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathbf{X}), \mathcal{Y}_L$, target homophily ratio h, and target node degree d Output: $\mathcal{G}' = (\mathcal{V}, \mathcal{E}', \mathbf{X})$ 1: Split the edges \mathcal{E} into heterophilic edges \mathcal{E}_n and homophilic edges \mathcal{E}_s . 2: if $|\mathcal{E}_s| \geq hd|\mathcal{V}|$ then Sample $hd|\mathcal{V}|$ edges from \mathcal{E}_s to get \mathcal{E}'_s 3: 4: else Obtain $hd|\mathcal{V}| - |\mathcal{E}_s|$ homophilic edges by randomly link nodes in the same class 5: Combine \mathcal{E}_s with added homophilic edges to obtain \mathcal{E}'_s 6: 7: end if 8: Randomly sample $d(1-h)|\mathcal{V}|$ edges from \mathcal{E}_n as \mathcal{E}'_n 9: Get \mathcal{E}' with $\mathcal{E}' = \mathcal{E}'_n \cup \mathcal{E}'_s$

E.1 Process of Graph Generation 651

To verify the conclusion in Theorem 1, we generate graphs with different homophily ratios and 652 average degrees on the large-scale crocodile graph. Specifically, the average node degree of the target 653 generated graphs is varied by $\{5, 10, 20\}$. For each node degree, we will sample the heterophilic 654 edges, i.e., edges linking nodes in different classes, and homophilic edges, i.e., edges linking nodes 655 in the same class from the original crocodile graph in different ratios to obtain realistic graphs with 656 different heterophily levels. The homophily ratios of the generated graphs range from 0 to 0.9 with 657 a step of 0.1. Since crocodile itself is a heterophilic graph that do not contain many homophilic 658 edges, there could be no enough homophilic edges to obtain a graph with high homophily and node 659 degrees. In this situation, we will randomly link nodes in the same class to get the required number of 660 homophilic edges for graph generation. For the train/validation/test splits of generated graphs, they 661 are the same as the original crocodile graph. The algorithm of the graph generation process can be 662 found in Algorithm 2. 663

More Experiments on Generated Graphs E.2 664

To verify our theoretical analysis that label-wise aggregation can lead to distinguishable representa-665 tions regardless the heterophily levels under mild conditions, we also compare LW-GCN with GCN 666 and GAT on the generated graphs with different homophily ratios and average node degrees. The 667 label-wise aggregation is conducted with the pseudo labels and provided ground-truth labels as it is 668 described in Sec.4.2.2. Since we only focus on the label-wise graph convolution in the experiments, 669 the model selection module is removed here. The other settings are the same as description in 670 Appendix B.1. The average results of 10 splits are shown in Fig. 4. From this figure, we can observe 671 that the performance of LW-GCN is much better than the GCN and GAT when the heterophily level is 672



Figure 4: Comparisons between GCN, GAT and our LW-GCN on generated graphs. Note that model selection module is not adopted in LW-GCN in these experiments.

high. For example, when $h \approx 0.2$, both GCN and GAT can hardly outperforms MLP. By contrast, the accuracy of LW-GCN outperform GCN and GAT by around 10%. This demonstrates the effectiveness of adopting label-wise aggregation in graph convolution. In addition, we can find that only adopting the model with label-wise graph convolution will give slightly worse performance than GCN/GAT when the homophily ratio is very high. This implies the necessity of deploying a model selection module.

679 F Analysis on Heterophilic Graphs



Figure 5: Similarity matrices of neighbors linked with centered nodes in different classes on Crocodile, Squirrel, and Chameleon.

In this section, we conduct empirical analysis to verify Assumption 2 in Sec.3.2. Specifically, we 680 aim to show (i) For nodes in the same class, features of their neighbors in the same class are similar; 681 (ii) For nodes in different classes, features of their neighbors in the same class follow different 682 distributions. Let $\mathcal{X}_{ik} = \{x_u : y_u = k, y_v = i, u \in \mathcal{N}(v), v \in \mathcal{V}\}$ be the set of neighbors which 683 belong to class k and are linked by the central node in class i. For neighbors in class k, we analyze 684 the average similarity scores between \mathcal{X}_{ik} and \mathcal{X}_{jk} to investigate whether neighbors in class k that 685 are linked by center nodes in different classes follow different distributions. Specifically, the average 686 similarity score between \mathcal{X}_{ik} and \mathcal{X}_{jk} is obtained by 687

$$s(\mathcal{X}_{ik}, \mathcal{X}_{jk}) = \frac{1}{|\mathcal{X}_{ik}| \times |\mathcal{X}_{jk}|} \sum_{v_i \in \mathcal{X}_{ik}} \sum_{v_j \in \mathcal{X}_{jk}} \frac{\mathbf{x}_i \cdot \mathbf{x}_j}{\|\mathbf{x}_i\| \|\mathbf{x}_j\|},$$
(33)

where \mathbf{x}_i and \mathbf{x}_j are features of node $v_i \in \mathcal{X}_{ik}$ and $v_j \in \mathcal{X}_{jk}$, respectively. The results on Crocodile, Chameleon, and Squirrel for representative neighbor classes are presented in Fig. 5, where (i, j)-th element in the similarity matrix denotes the average node feature cosine similarity between \mathcal{X}_{ik} and \mathcal{X}_{jk} . From this figure, we can observe that:

- For $\mathcal{X}_{ik}, \forall i \in 1, ..., C$, its intra-group similarity score is very high. This proves that the heterophilic neighbors' features are similar when the nodes are in the same class.
- The similarity scores between \mathcal{X}_{ik} and \mathcal{X}_{ik} are very small when $i \neq j$. This indicates that for nodes in different classes their heterophilic neighbors belonging to the same class still differs a lot.
- ⁶⁹⁶ With the above observations, Assumption 2 is justified.

697 G Additional Ablation Studies

Table 4 gives additional ablation studies on Pubmed, Chameleon, and Squirrel. The observations are similar to that of Table 2.

Dataset	MLP	GCN	GCNII	LW-GCN\P	LW-GCN\G	LW-GCN_{GCN}	LW-GCN
Pubmed	72.7 ± 0.4	78.4 ± 1.1	$80.2\pm\!0.2$	77.6 ± 0.7	$72.4\pm\!0.6$	$79.2\pm\!0.8$	$\textbf{80.3} \pm 0.3$
Chameleon Squirrel	$\begin{array}{c} 48.0 \pm 1.5 \\ 32.3 \pm 1.8 \end{array}$	$\begin{array}{c} 63.5 \pm \! 2.5 \\ 46.7 \pm \! 1.5 \end{array}$	$\begin{array}{c} 63.5 \pm \! 2.5 \\ 49.4 \pm \! 1.7 \end{array}$	74.7 ±1.4 62.3 ±2.3	$\begin{array}{c} 74.2 \pm 1.8 \\ 62.3 \pm 1.3 \end{array}$	$74.3 \pm 2.3 \\ 61.9 \pm 1.4$	74.4 ±1.2 62.6 ±1.6

Table 4: Ablation Study

699

700 H Additional Experimental Results

The additional experimental results on Cornell and Citeseer datasets are presented in Table 5 and Table 6. The observations are similar to that of Table 1.

Table 5: Additional comparisons with GNNs originally designed for graph with homophily.

Dataset	MLP	GCN	MixHop	SuperGAT	GCNII	LW-GCN
Cornell Citeseer	$\begin{array}{c} 79.2 \pm \! 5.7 \\ 60.3 \pm \! 0.4 \end{array}$	$\begin{array}{c} 57.3 \pm \! 5.8 \\ 71.3 \pm \! 0.3 \end{array}$	$\begin{array}{c} 79.5 \pm \! 6.3 \\ 68.7 \pm \! 0.3 \end{array}$	$\frac{57.3 \pm 4.3}{72.2 \pm 0.8}$	$\frac{80.3 \pm 5.3}{72.0 \pm 0.8}$	$\begin{array}{c} 84.3 \pm \! 5.2 \\ 72.3 \pm \! 0.4 \end{array}$

Table 6: Additional comparisons with GNNs designed for graph with heterophily.

Dataset	FAGCN	SimP-GCN	H2GCN	GPRGNN	BM-GCN	ASGC	LINKX	GloGNN+	LW-GCN
Cornell Citeseer	$\begin{array}{c} 78.3 \pm \! 4.5 \\ 71.7 \pm \! 0.6 \end{array}$	$\frac{81.4 \pm 7.4}{\underline{71.8 \pm 0.8}}$	$\begin{array}{c} 79.7 \pm \! 5.0 \\ 71.0 \pm \! 0.5 \end{array}$	$\begin{array}{c} 77.6 \pm \! 5.0 \\ 71.1 \pm \! 0.9 \end{array}$	$\begin{array}{c} 74.6 \pm \! 5.0 \\ 68.9 \pm \! 1.0 \end{array}$	$\begin{array}{c} 79.2 \pm \! 5.2 \\ 70.2 \pm \! 0.2 \end{array}$	$\begin{array}{c} 77.8 \pm \! 5.8 \\ 51.6 \pm \! 1.7 \end{array}$	$\begin{array}{c} \textbf{85.9} \pm \textbf{4.4} \\ \textbf{66.7} \pm \textbf{1.9} \end{array}$	$\frac{84.3\pm\!5.2}{72.3\pm\!0.4}$

I Impacts of Label-Wise Aggregation Layers

In this section, we explore the sensitivity of LW-GCN on the depth of f_C , i.e., the number of layers 704 of label-wise message passing. Since LW-GCN will not select f_C for homophilic graphs. We only 705 conduct the sensitivity analysis on heterophilic graphs. We vary the depth of f_C as $\{2, 3, \dots, 6\}$. The 706 other experimental settings are the same as that described in Sec. B.1. The results on Chameleon 707 and Squirrel are shown in Fig. 6. From the figure, we find that our LW-GCN is insensitive to the 708 number of layers, while the performance of GCN will drop with the increase of depth. This is because 709 aggregation of LW-GCN is performed label-wisely to capture the context information. Embeddings 710 of nodes in different classes will not be smoothed to similar values even after many iterations. 711



Figure 6: Classification accuracy with different model depth.

712 J Limitations of Our Work

In this paper, we conduct thoroughly theoretical and empirical analysis to show the impacts of heterophily levels to GCN. And we demonstrate the GCN model can be largely affected by heterophily and give poor prediction results. To alleviate the issue brought by heterophily, we develop a novel label-wise graph convolutional network to preserve the heterophilic context to facilitate the node

classification. However, there are some limitations of our work. First, node labels are required 717 for LW-GCN to obtain pseudo labels for label-wise graph convolution. However, in some tasks 718 such as link prediction, labels are not available. Therefore, we will investigate how to obtain useful 719 pseudo labels for applications that do not provide node labels. Second, in our theoretical analysis, 720 we make several assumptions for simplification. Concretely, we conduct analysis on the d-regular 721 graph. Following [42, 25], we also make an assumption on the label distribution of neighbor nodes. In our analysis, the node features are simplified to normal distribution and dimensions of features 723 are independent to each other. These assumptions may not hold for some real-world graphs. For 724 example, node degrees of the real network can be unbalanced which will contradict the assumption 725 of d-regularity. The label distributions and feature distributions of neighbor nodes can be much more 726 complex. Therefore, we will investigate the theoretical analysis on more flexible assumptions in the future. Third, recent studies [25, 29] show that the edge homophily ratio used in this paper could 728 have significant drawbacks especially when the distribution of classes is unbalanced. To address these 729 drawbacks, new measures such as adjusted homophily and label informativeness are proposed [29]. 730

⁷³¹ We leave the extension of our analysis on these new homophily measures as the future work.