
Revisiting the Message Passing in Heterophilous Graph Neural Networks

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Graph Neural Networks (GNNs) have demonstrated strong performance in graph
2 mining tasks due to their message-passing mechanism, which is aligned with the
3 homophily assumption that adjacent nodes exhibit similar behaviors. However,
4 in many real-world graphs, connected nodes may display contrasting behaviors,
5 termed as *heterophilous* patterns, which has attracted increased interest in het-
6 erophilous GNNs (HTGNNs). Although the message-passing mechanism seems
7 unsuitable for heterophilous graphs due to the propagation of class-irrelevant infor-
8 mation, it is still widely used in many existing HTGNNs and consistently achieves
9 notable success. This raises the question: *why does message passing remain effec-*
10 *tive on heterophilous graphs?* To answer this question, in this paper, we revisit the
11 message-passing mechanisms in heterophilous graph neural networks and reform-
12 ulate them into a unified heterophilous message-passing (HTMP) mechanism.
13 Based on HTMP and empirical analysis, we reveal that the success of message
14 passing in existing HTGNNs is attributed to implicitly enhancing the compatibility
15 matrix among classes. Moreover, we argue that the full potential of the compat-
16 ibility matrix is not completely achieved due to the existence of incomplete and
17 noisy semantic neighborhoods in real-world heterophilous graphs. To bridge this
18 gap, we introduce a new approach named CMGNN, which operates within the
19 HTMP mechanism to explicitly leverage and improve the compatibility matrix. A
20 thorough evaluation involving 10 benchmark datasets and comparative analysis
21 against 13 well-established baselines highlights the superior performance of the
22 HTMP mechanism and CMGNN method.

23 1 Introduction

24 Graph Neural Networks (GNNs) have shown remarkable performance in graph mining tasks, such
25 as social network analysis [1, 2] and recommender systems [3, 4]. The design principle of GNNs is
26 typically based on the homophily assumption [5], which assumes that nodes are inclined to exhibit
27 behaviors similar to their neighboring nodes [6]. However, this assumption does not always hold
28 in real-world graphs, where the connected nodes demonstrate a contrasting tendency known as the
29 *heterophily* [7]. In response to the challenges of heterophily in graphs, *heterophilous GNNs (HTGNNs)*
30 have attracted considerable research interest [6, 8–10], with numerous innovative approaches being
31 introduced recently [11–24]. However, the majority of these methods continue to employ a message-
32 passing mechanism, which was not originally designed for heterophilous graphs, as they tend to
33 incorporate excessive information from disparate classes. This naturally raises a question: *Why does*
34 *message passing remain effective on heterophilous graphs?*

35 Recently, a few efforts [6] have begun to investigate this question and reveal that vanilla message
36 passing can work on heterophilous graphs under certain conditions. However, the absence of a unified

37 and comprehensive understanding of message passing within existing HTGNNs has hindered the
 38 creation of innovative approaches. In this paper, we first revisit the message-passing mechanisms
 39 in existing HTGNNs and reformulate them into a unified heterophilous message-passing (HTMP)
 40 mechanism, which extends the definition of neighborhood in various ways and simultaneously utilizes
 41 the messages of multiple neighborhoods. Specifically, HTMP consists of three major steps namely
 42 aggregating messages with explicit guidance, combining messages from multiple neighborhoods, and
 43 fusing intermediate representations.

44 Equipped with HTMP, we further conduct empirical analysis on real-world graphs. The results reveal
 45 that the success of message passing in existing HTGNNs is attributed to **implicitly enhancing the**
 46 **compatibility matrix**, which exhibits the probabilities of observing edges among nodes from different
 47 classes. In particular, by increasing the distinctiveness between the rows of the compatibility matrix
 48 via different strategies, the node representations of different classes become more discriminative in
 49 heterophilous graphs.

50 Drawing from previous observations, we contend that nodes within real-world graphs might exhibit a
 51 semantic neighborhood that only reveals a fraction of the compatibility matrix, accompanied by noise.
 52 This could limit the effectiveness of enhancing the compatibility matrix and result in suboptimal
 53 representations. To fill this gap, we further propose a novel Compatibility Matrix-aware Graph Neural
 54 Network (CMGNN) under HTMP mechanism, which utilizes the compatibility matrix to construct
 55 desired neighborhood messages as supplementary for nodes and explicitly enhances the compatibility
 56 matrix by a targeted constraint. We build a benchmark to fairly evaluate CMGNN and existing
 57 methods, which encompasses 13 diverse baseline methods and 10 datasets that exhibit varying
 58 levels of heterophily. Extensive experimental results demonstrate the superiority of CMGNN and
 59 HTMP mechanism. The contributions of this paper are summarized as:

- 60 • We revisit the message-passing mechanisms in existing HTGNNs and reformulate them into a
 61 unified heterophilous message-passing mechanism (HTMP), which not only provides a macroscopic
 62 view of message passing in HTGNNs but also enables people to develop new methods flexibly.
- 63 • We reveal that the effectiveness of message passing on heterophilous graphs is attributed to
 64 implicitly enhancing the compatibility matrix among classes, which gives us a new perspective to
 65 understand the message passing in HTGNNs.
- 66 • Based on HTMP mechanism and empirical analysis, we propose CMGNN to unlock the potential
 67 of the compatibility matrix in HTGNNs. We further build a unified benchmark that overcomes the
 68 issues of current datasets for fair evaluation¹. Experiments show the superiority of CMGNN.

69 2 Preliminaries

70 Given a graph $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathbf{X}, \mathbf{A}, \mathbf{Y})$, \mathcal{V} is the node set and \mathcal{E} is the edge set. Nodes are characterized
 71 by the feature matrix $\mathbf{X} \in \mathbb{R}^{N \times d_f}$, where $N = |\mathcal{V}|$ denotes the number of nodes, d_f is the features
 72 dimension. $\mathbf{Y} \in \mathbb{R}^{N \times 1}$ is the node labels with the one-hot version $\mathbf{C} \in \mathbb{R}^{N \times K}$, where K is
 73 the number of node classes. The neighborhood of node v_i is denoted as \mathcal{N}_i . $\mathbf{A} \in \mathbb{R}^{N \times N}$ is
 74 the adjacency matrix, and $\mathbf{D} = \text{diag}(\mathbf{d}_1, \dots, \mathbf{d}_n)$ represents the diagonal degree matrix, where
 75 $\mathbf{d}_i = \sum_j \mathbf{A}_{ij}$. $\tilde{\mathbf{A}} = \mathbf{A} + \mathbf{I}$ represents the adjacency matrix with self-loops. Let $\mathbf{Z} \in \mathbb{R}^{N \times d_r}$ be the
 76 node representations with dimension d_r learned by the models. We use $\mathbf{1}$ to represent a matrix with
 77 all elements equal to 1, and $\mathbf{0}$ for a matrix with all elements equal to 0.

78 **Homophily and Heterophily.** High homophily is observed in graphs where a substantial portion of
 79 connected nodes shares identical labels, while high heterophily corresponds to the opposite situation.
 80 For measuring the homophily level, two widely used metrics are edge homophily h^e [12] and node
 81 homophily h^n [15], defined as $h^e = \frac{|\{e_{u,v} | e_{u,v} \in \mathcal{E}, \mathbf{Y}_u = \mathbf{Y}_v\}|}{|\mathcal{E}|}$ and $h^n = \frac{1}{|\mathcal{V}|} \sum_{v \in \mathcal{V}} \frac{|\{u | u \in \mathcal{N}_v, \mathbf{Y}_u = \mathbf{Y}_v\}|}{d_v}$.
 82 Both metrics have a range of $[0, 1]$, where higher values indicate stronger homophily and lower values
 83 indicate stronger heterophily.

84 **Vanilla Message Passing (VMP).** The vanilla message-passing mechanism plays a pivotal role in
 85 transforming and updating node representations based on the neighborhood [25]. Typically, the

¹Codebase is available at the supplementary material.

Table 1: Revisiting the message passing in representative heterophilous GNNs under the perspective of HTMP mechanism.

Method	Neighborhood Indicators		Aggregation Guidance		COMBINE	FUSE
	Type	\mathcal{A}	Type	\mathcal{B}		
GCN [1]	Raw	$[\tilde{\mathbf{A}}]$	DegAvg	$[\tilde{\mathbf{B}}^d]$	/	$\mathbf{Z} = \mathbf{Z}^L$
APPNP [26]		$[\mathbf{I}, \tilde{\mathbf{A}}]$		$[\mathbf{I}, \tilde{\mathbf{B}}^d]$	WeightedAdd	$\mathbf{Z} = \mathbf{Z}^L$
GCNII [27]		$[\mathbf{I}, \tilde{\mathbf{A}}]$		$[\mathbf{I}, \tilde{\mathbf{B}}^d]$	WeightedAdd	$\mathbf{Z} = \mathbf{Z}^L$
GAT [28]		$[\tilde{\mathbf{A}}]$	AdaWeight	$[\mathbf{B}^{aw}]$	/	$\mathbf{Z} = \mathbf{Z}^L$
GPR-GCN [20]		$[\tilde{\mathbf{A}}]$	DegAvg	$[\tilde{\mathbf{B}}^d]$	/	AdaAdd
OrderedGNN [21]		$[\mathbf{I}, \mathbf{A}]$		$[\mathbf{I}, \mathbf{B}^d]$	AdaCat	$\mathbf{Z} = \mathbf{Z}^L$
ACM-GCN [18]		$[\mathbf{I}, \mathbf{A}, \tilde{\mathbf{A}}]$		$[\mathbf{I}, \mathbf{B}^d, \mathbf{I} - \mathbf{B}^d]$	AdaAdd	$\mathbf{Z} = \mathbf{Z}^L$
FAGCN [11]		$[\mathbf{I}, \mathbf{A}]$		$[\mathbf{I}, \mathbf{B}^{naw}]$	WeightedAdd	$\mathbf{Z} = \mathbf{Z}^L \mathbf{W}$
GBK-GNN [24]		$[\mathbf{I}, \mathbf{A}, \mathbf{A}]$	AdaWeight	$[\mathbf{I}, \mathbf{B}^{aw}, \mathbf{1} - \mathbf{B}^{aw}]$	Add	$\mathbf{Z} = \mathbf{Z}^L$
SimP-GCN [14]		$[\mathbf{I}, \tilde{\mathbf{A}}, \mathbf{A}_f]$	DegAvg	$[\mathbf{I}, \tilde{\mathbf{B}}^d, \mathbf{B}_f^d]$	AdaAdd	$\mathbf{Z} = \mathbf{Z}^L$
H2GCN [12]		$[\mathbf{A}, \mathbf{A}_{h2}]$		$[\mathbf{B}^d, \mathbf{B}_{h2}^d]$	Cat	Cat
Geom-GCN [15]		$[\mathbf{A}_{c1}, \dots, \mathbf{A}_{cr}, \dots, \mathbf{A}_{cR}]$		$[\mathbf{B}_{c1}^d, \dots, \mathbf{B}_{cr}^d, \dots, \mathbf{B}_{cR}^d]$	Cat	$\mathbf{Z} = \mathbf{Z}^L$
MixHop [16]		$[\mathbf{I}, \mathbf{A}, \mathbf{A}_{h2}, \dots, \mathbf{A}_{hk}]$	AdaWeight	$[\mathbf{I}, \mathbf{B}^d, \mathbf{B}_{h2}^d, \dots, \mathbf{B}_{hk}^d]$	Cat	$\mathbf{Z} = \mathbf{Z}^L$
UGCN [13]		$[\tilde{\mathbf{A}}, \tilde{\mathbf{A}}_{h2}, \mathbf{A}_f]$		$[\tilde{\mathbf{B}}^{aw}, \tilde{\mathbf{B}}_{h2}^{aw}, \mathbf{B}_f^{aw}]$	AdaAdd	$\mathbf{Z} = \mathbf{Z}^L$
WRGNN [22]	$[\mathbf{A}_{c1}, \dots, \mathbf{A}_{cr}, \dots, \mathbf{A}_{cR}]$	$[\mathbf{B}_{c1}^{aw}, \dots, \mathbf{B}_{cr}^{aw}, \dots, \mathbf{B}_{cR}^{aw}]$		Add	$\mathbf{Z} = \mathbf{Z}^L$	
HOG-GCN [17]	$[\mathbf{I}, \mathbf{A}_{hk}]$	RelaEst		$[\mathbf{I}, \mathbf{B}^{re}]$	WeightedAdd	$\mathbf{Z} = \mathbf{Z}^L$
GloGNN [19]	$[\mathbf{I}, \mathbf{1}]$		$[\mathbf{I}, \mathbf{B}^{re}]$	WeightedAdd	$\mathbf{Z} = \mathbf{Z}^L$	
GGCN [23]	Dis		$[\mathbf{I}, \mathbf{A}_p, \mathbf{A}_n]$	$[\mathbf{I}, \mathbf{B}_p^{re}, \mathbf{B}_n^{re}]$	AdaAdd	$\mathbf{Z} = \mathbf{Z}^L$

* The correspondence between the full form and the abbreviation: Raw Neighborhood (Raw), Neighborhood Redefine (ReDef), Neighborhood Discrimination (Dis), Degree-based Averaging (DegAvg), Adaptive Weights (AdaWeight), Relation Estimation (RelaEst), Addition (Add), Weighted Addition (WeightAdd), Adaptive Weighted Addition (AdaAdd), Concatenation (Cat), Adaptive Dimension Concatenation (AdaCat).
 * More details about the notations are available in Appendix A.1.

86 mechanism operates iteratively and comprises two stages:

$$\tilde{\mathbf{Z}}^l = \text{AGGREGATE}(\mathbf{A}, \mathbf{Z}^{l-1}), \quad \mathbf{Z}^l = \text{COMBINE}(\mathbf{Z}^{l-1}, \tilde{\mathbf{Z}}^l), \quad (1)$$

87 where the AGGREGATE function first aggregates the input messages \mathbf{Z}^{l-1} from neighborhood \mathbf{A}
 88 into the aggregated one $\tilde{\mathbf{Z}}^l$, and subsequently, the COMBINE function combines the messages of
 89 node ego and neighborhood aggregation, resulting in updated representations \mathbf{Z}^l .

90 3 Revisiting Message Passing in Heterophilous GNNs.

91 To gain a thorough and unified insight into the effectiveness of message passing in HTGNNs, we
 92 revisit message passing in various notable HTGNNs [11–24] and propose a unified heterophilous
 93 message passing (HTMP) mechanism, structured as follows:

$$\tilde{\mathbf{Z}}_r^l = \text{AGGREGATE}(\mathbf{A}_r, \mathbf{B}_r, \mathbf{Z}^{l-1}), \quad \mathbf{Z}^l = \text{COMBINE}(\{\tilde{\mathbf{Z}}_r^l\}_{r=1}^R), \quad \mathbf{Z} = \text{FUSE}(\{\mathbf{Z}^l\}_{l=0}^L). \quad (2)$$

94 Generally, HTMP extends the definition of neighborhood in various ways and simultaneously utilize
 95 the messages of multiple neighborhoods, which is the key for better adapting to heterophily. We
 96 use R to denote the number of neighborhoods used by the model. In each message passing layer l ,
 97 HTMP separately aggregates messages within R neighborhoods and combines them. The method-
 98 ological analysis of some representative HTGNNs and more details can be seen in Appendix A.
 99 Compared to the VMP mechanism, HTMP mechanism has advances in the following functions:

100 (i) To characterize different neighborhoods, the AGGREGATE function in HTMP includes the **neigh-**
 101 **borhood indicator** \mathbf{A}_r to indicate the neighbors within a specific neighborhood r . The adjacency
 102 matrix \mathbf{A} in VMP is a special neighborhood indicator that marks the neighbors in the raw neigh-
 103 borhood. To further characterize the aggregation of different neighborhoods, HTMP introduces the
 104 **aggregation guidance** \mathbf{B}_r for each neighborhood r . In VMP, the aggregation guidance is an implicit
 105 parameter of the AGGREGATE function since it only works for the raw neighborhood. A commonly
 106 used form of the AGGREGATE function is $\text{AGGREGATE}(\mathbf{A}_r, \mathbf{B}_r, \mathbf{Z}^{l-1}) = (\mathbf{A}_r \odot \mathbf{B}_r) \mathbf{Z}^{l-1} \mathbf{W}_r^l$,
 107 where \odot is the Hadamard product and \mathbf{W}_r^l is a weight matrix for message transformation. We take

108 this as the general form of the AGGREGATE function and only analyze the neighborhood indicators
 109 and the aggregation guidance in the following.

110 The *neighborhood indicator* $\mathbf{A}_r \in \{0, 1\}^{N \times N}$ indicates neighbors associated with central nodes
 111 within neighborhood r . To describe the multiple neighborhoods in HTGNNs, neighborhood indicators
 112 can be formed as a list $\mathcal{A} = [\mathbf{A}_1, \dots, \mathbf{A}_r, \dots, \mathbf{A}_R]$. For the sake of simplicity, we consider the identity
 113 matrix $\mathbf{I} \in \mathbb{R}^{N \times N}$ as a special neighborhood indicator for acquiring the nodes' ego messages. The
 114 *aggregation guidance* $\mathbf{B}_r \in \mathbb{R}^{N \times N}$ can be viewed as pairwise aggregation weights in most cases,
 115 which has the multiple form $\mathcal{B} = [\mathbf{B}_1, \dots, \mathbf{B}_r, \dots, \mathbf{B}_R]$. Table 1 illustrates the connection between
 116 message passing in various HTGNNs and HTMP mechanism.

117 (ii) Considering the existence of multiple neighborhoods, the **COMBINE** function in HTMP need to
 118 integrate multiple messages instead of only the ego node and the raw neighborhood. Thus, the input
 119 of the COMBINE function is a set of messages $\tilde{\mathbf{Z}}_r^l$ aggregated from the corresponding neighborhoods.
 120 In HTGNNs, addition and concatenation are two common approaches, each of which has variants.
 121 An effective COMBINE function is capable of simultaneously processing messages from various
 122 neighborhoods while preserving their distinct features, thereby reducing the effects of heterophily.

123 (iii) In VMP, the final output representations are usually the one of the final layer: $\mathbf{Z} = \mathbf{Z}^L$. Some
 124 HTGNNs utilize the combination of intermediate representations to leverage messages from different
 125 localities, adapting to the heterophilous structural properties in different graphs. Thus, we introduce
 126 an additional **FUSE** function in HTMP which integrates multiple representations \mathbf{Z}^l of different
 127 layers l into the final \mathbf{Z} . Similarly, the FUSE function is based on addition and concatenation.

128 4 Why Does Message Passing Still Remain Effective in Heterophilous 129 Graphs?

130 Based on HTMP mechanism, we further dive into the motivation behind the message passing of
 131 existing HTGNNs. Our discussion begins by examining the difference between homophilous and
 132 heterophilous graphs. Initially, we consider the homophily ratios h^e and h^n , as outlined in Section 2.
 133 However, a single number is not able to indicate enough conditions of a graph. Ma et al. [6] propose
 134 the existence of a special case of heterophily, named "good" heterophily, where the VMP mechanism
 135 can achieve strong performance and the homophily ratio shows no difference. Thus, to better study
 136 the heterophily property, here we introduce the *Compatibility Matrix* [7] to describe graphs:

137 **Definition 1 Compatibility Matrix (CM):** *The potential connection preference among classes within*
 138 *a graph. It's formatted as a matrix $\mathbf{M} \in \mathbb{R}^{K \times K}$, where the i -th row \mathbf{M}_i denotes the connection*
 139 *probabilities between class i and all classes. It can be estimated empirically by the statistics among*
 140 *nodes as follows:*

$$140 \quad \mathbf{M} = \text{Norm}(\mathbf{C}^T \mathbf{C}^{nb}), \quad \mathbf{C}^{nb} = \hat{\mathbf{A}} \mathbf{C}, \quad (3)$$

141 *where $\text{Norm}(\cdot)$ denotes the L1 normalization and T is the matrix transpose operation. $\mathbf{C}^{nb} \in \mathbb{R}^{N \times K}$*
 142 *is the **semantic neighborhoods** of nodes, which indicates the proportion of neighbors from each class*
 143 *in nodes' neighborhoods.*

144 We visualize the CM of a homophilous graph Photo [29] and a heterophilous graph Amazon-
 145 Ratings [30] in Figure 1(a) and 1(b). The CM in Photo displays an identity-like matrix, where the
 146 diagonal elements can be viewed as the homophily level of each class. With this type of CM, the VMP
 147 mechanism learns representations comprised mostly of messages from same the class, while messages
 148 of other classes are diluted. *Then how does HTMP mechanism work on heterophilous graphs without*
 149 *an identity-like CM?* The "good" heterophily inspires us, which we believe corresponds to a CM with
 150 enough discriminability among classes. We conduct experiments on synthetic graphs to confirm this
 151 idea, with details available in Appendix C. Also, we find "good" heterophily in real-world graphs
 152 though it's not as significant as imagined. Thus, we have the following observation:

153 **Observation 1 (Connection between CM and VMP).** *When enough (depends on data) discriminabil-*
 154 *ity exists among classes in CM, vanilla message passing can work well in heterophilous graphs.*

155 With this observation, we have a conjecture: *Is HTMP mechanism trying to enhance the discriminabil-*
 156 *ity of CM?* Some special designs in HTMP intuitively meet this. For example, *feature-similarity-based*
 157 *neighborhood indicators* and *neighborhood discrimination* are designed to construct neighborhoods

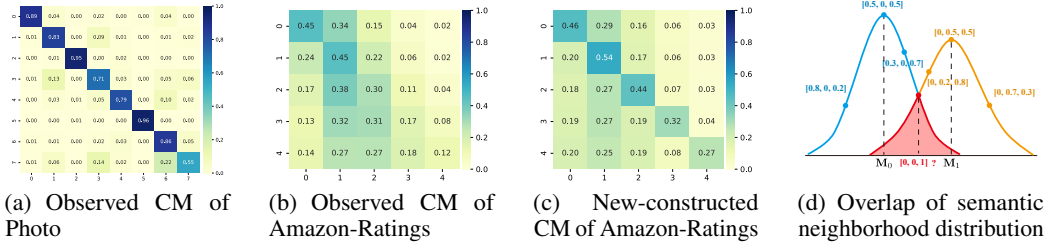


Figure 1: Visualizations of the compatibility matrix and the example of distribution overlap.

158 with high homophily, that is, an identity-like CM with high discriminability. We plot the CM of
 159 feature-similarity-based neighborhood on Amazon-Ratings in Figure 1(c) to confirm it. Moreover,
 160 we investigate two representative methods ACM-GCN [18] and GPRGNN [20], showing that they
 161 also meet this conjecture with the posterior proof in Appendix D. ACM-GCN combines the messages
 162 of node ego, low-frequency and high-frequency with adaptive weights, which actually motifs the
 163 edge weights and node weights to build a new CM. GPRGNN has a FUSE function with adaptive
 164 weights while other settings are the same as GCN. It actually integrates the CMs of multiple-order
 165 neighborhoods with adaptive weights to form a more discriminative CM. These lead to the answer to
 166 the aforementioned question:

167 **Observation 2** (Connection between CM and HTMP). *The unified goal of various message passing*
 168 *in existing HTGNNs is to utilize and enhance the discriminability of CM on heterophilous graphs.*
 169 *In other words, the success of message passing in existing HTGNNs benefits from utilizing and*
 170 *enhancing the discriminability of CM.*

171 Furthermore, we notice that the power of CM is not fully released due to the incomplete and noisy
 172 semantic neighborhoods in real-world heterophilous graphs. We use the perspective of distribution
 173 to describe the issue more intuitively: The semantic neighborhoods of nodes from the same class
 174 collectively form a distribution, whose mean value indicates the connection preference of that class,
 175 i.e. M_i for class i . Influenced by factors such as degree and randomness, the semantic neighborhood
 176 of nodes in real-world graphs may display only a fraction of CM accompanied by noise. It can
 177 lead to the overlap between different distributions as shown in Figure 1(d), where the existence of
 178 overlapping parts means nodes from different classes may have the same semantic neighborhood.
 179 This brings a great challenge since the overlapping semantic neighborhood may become redundant
 180 information during message passing.

181 5 Method

182 To fill this gap, we further propose a method named Compatibility Matrix-Aware GNN (CMGNN),
 183 which leverages the CM to construct desired neighborhood messages as supplementary, providing
 184 valuable neighborhood information for nodes to mitigate the impact of incomplete and noisy se-
 185 mantic neighborhoods. The desired neighborhood message denotes the averaging message within
 186 a neighborhood when a node’s semantic neighborhoods meet the CM of the corresponding class,
 187 which converts the discriminability from CM into messages. CMGNN follows the HTMP mechanism
 188 and constructs a supplementary neighborhood indicator along with the corresponding aggregation
 189 guidance to introduce supplementary messages. Further, CMGNN introduces a simple constraint to
 190 explicitly enhance the discriminability of CM.

191 **Message Passing in CMGNN.** CMGNN aggregates messages from three neighborhoods for
 192 each node, including the ego neighborhood, raw neighborhood, and supplementary neighborhood.
 193 Following the HTMP mechanism, the message passing of CMGNN can be described as follows:

$$\begin{aligned}
 \tilde{\mathbf{Z}}_r^l &= \text{AGGREGATE}(\mathbf{A}_r, \mathbf{B}_r, \mathbf{Z}^{l-1}) = (\mathbf{A}_r \odot \mathbf{B}_r) \mathbf{Z}^{l-1} \mathbf{W}_r^l, \\
 \mathbf{Z}^l &= \text{COMBINE}(\{\tilde{\mathbf{Z}}_r^l\}_{r=1}^3) = \text{AdaWeight}(\{\tilde{\mathbf{Z}}_r^l\}_{r=1}^3), \\
 \mathbf{Z} &= \text{FUSE}(\{\mathbf{Z}^l\}_{l=0}^L) = \parallel_{l=0}^L \mathbf{Z}^l,
 \end{aligned}
 \tag{4}$$

194 where AdaWeight is the adaptive weighted addition implemented by an MLP with Softmax, \parallel denotes
 195 the concatenation. The neighborhood indicators and aggregation guidance of the three neighborhoods
 196 are formatted as follows:

$$\mathbf{A}_1^l = \mathbf{I}, \mathbf{B}_1^l = \mathbf{I}, \quad \mathbf{A}_2^l = \mathbf{A}, \mathbf{B}_2^l = \mathbf{D}^{-1}\mathbf{1}, \quad \mathbf{A}_3^l = \mathbf{A}^{sup}, \mathbf{B}_3^l = \mathbf{B}^{sup}, \quad (5)$$

197 where \mathbf{A}^{sup} and \mathbf{B}^{sup} are described below.

198 The supplementary neighborhood indicator \mathbf{A}^{sup} assigns K additional virtual neighbors for each
 199 node: $\mathbf{A}^{sup} = \mathbf{1} \in \mathbb{R}^{N \times K}$. Specifically, these additional neighbors are K virtual nodes, constructed
 200 as the prototypes of classes based on the labels of the training set. The attributes $\mathbf{X}^{ptt} \in \mathbb{R}^{K \times d_f}$,
 201 neighborhoods $\mathbf{A}^{ptt} \in \mathbb{R}^{K \times N}$ and labels $\mathbf{Y}^{ptt} \in \mathbb{R}^{K \times K}$ of prototypes are defined as follows:

$$\mathbf{X}^{ptt} = \text{Norm}(\mathbf{C}_{train}^T \mathbf{X}_{train}), \mathbf{A}^{ptt} = \mathbf{0}, \mathbf{Y}^{ptt} = \mathbf{I}, \quad (6)$$

202 where \mathbf{C}_{train} and \mathbf{X}_{train} are the one-hot labels and attributes of nodes in the training set. Utilizing
 203 class prototypes as supplementary neighborhoods can provide each node with representative messages
 204 of classes, which builds the basis for desired neighborhood messages.

205 The supplementary aggregation guidance $\mathbf{B}^{sup} = \hat{\mathbf{C}}\hat{\mathbf{M}}$ indicates the desired semantic neighborhood
 206 of nodes, i.e. the desired proportion of neighbors from each class in nodes' neighborhoods according
 207 to the probability that nodes belong to each class. $\hat{\mathbf{M}}$ is the estimated compatibility matrix described
 208 in below. Using soft logits instead of one-hot pseudo labels preserves the real characteristics of nodes
 209 and reduces the impact of wrong predictions. During the message aggregation in the supplementary
 210 neighborhoods, the input representations \mathbf{Z}^{l-1} are replaced by the representations of virtual prototype
 211 nodes \mathbf{Z}_{ptt}^{l-1} , which are obtained by the same message-passing mechanism as real nodes.

212 Similar to existing methods [18, 19], we also regard topology structure as a kind of additional
 213 available node features. Thus, the input representation of the first layer can be obtained in two ways:

$$\mathbf{Z}^0 = [\mathbf{X}\mathbf{W}^X \parallel \hat{\mathbf{A}}\mathbf{W}^A]\mathbf{W}^0, \text{ or } \mathbf{Z}^0 = \mathbf{X}\mathbf{W}^0. \quad (7)$$

214 Note that in practice, we use ReLU as the activation function between layers. From the perspective of
 215 HTMP mechanism, our special design is to introduce an additional neighborhood indicator \mathbf{A}^{sup} by
 216 neighborhood redefining and aggregation guidance \mathbf{B}^{sup} , which can be seen as a form of relation
 217 estimation along with good interpretability. Meanwhile, these designs greatly reduce the time and
 218 space cost via the $N \times K$ form.

219 **Compatibility Matrix Estimation.** The CM can be directly calculated via Eq 3 with full-available
 220 labels. However, the label information is not entirely available in semi-supervised settings. Thus, we
 221 try to estimate the CM with the help of semi-supervised and pseudo labels. Since the pseudo labels
 222 predicted by the model might be wrong, which can lead to low-quality estimation, we introduce the
 223 confidence $\mathbf{g} \in \mathbb{R}^{N \times 1}$ based on the information entropy to reduce the impact of wrong predictions,
 224 where a high entropy means low confidence:

$$\mathbf{g}_i = \log K - \mathbf{H}(\hat{\mathbf{C}}_i) \in [0, \log K], \quad (8)$$

225 where $\hat{\mathbf{C}} \in \mathbb{R}^{N \times K}$ is the soft pseudo labels composed of labels from the training set and model
 226 predictions. Then the nodes' semantic neighborhoods $\mathbf{C}^{nb} = \text{Norm}(\mathbf{A}(\mathbf{g} \cdot \hat{\mathbf{C}})) \in \mathbb{R}^{N \times K}$ are
 227 calculated considering the confidence.

228 Further, the degrees of nodes also influence the estimation. As we mentioned in Section 4, the
 229 semantic neighborhood of low-degree nodes may display incomplete CM, leading to a significant gap
 230 between semantic neighborhoods and corresponding CM. Thus, they deserve low weights during the
 231 estimation. We manually set up two fixed thresholds and a weighting function range in $[0, 1]$:

$$\mathbf{w}_i^d = \begin{cases} \mathbf{d}_i/2K, & \mathbf{d}_i \leq K, \\ 0.25 + \mathbf{d}_i/4K, & K < \mathbf{d}_i \leq 3K, \\ 1, & \text{otherwise.} \end{cases} \quad (9)$$

232 When a node's degree \mathbf{d}_i is smaller than the number of classes K , its semantic neighborhood is
 233 unlikely to display complete CM, corresponding to a low weight. And when the node degree is
 234 greater than $3K$, we believe it can display near-complete CM, corresponding to the maximum weight.
 235 Finally, we can estimate the compatibility matrix $\hat{\mathbf{M}} \in \mathbb{R}^{K \times K}$ as follows:

$$\hat{\mathbf{M}} = \text{Norm}((\mathbf{w}^d \cdot \mathbf{g} \cdot \hat{\mathbf{C}})^T) \mathbf{C}^{nb}. \quad (10)$$

Table 2: Node classification accuracy comparison (%). The error bar (\pm) denotes the standard deviation of results over 10 trial runs. The best and second-best results in each column are highlighted in **bold** font and underlined. OOM denotes out-of-memory error during the model training.

Dataset	Roman-Empire	Amazon-Ratings	Chameleon-F	Squirrel-F	Actor	Flickr	BlogCatalog	Wikics	Pubmed	Photo	Avg. Rank
Homo.	0.05	0.38	0.25	0.22	0.22	0.24	0.4	0.65	0.8	0.83	
Nodes	22,662	24,492	890	2,223	7,600	7,575	5,196	11,701	19,717	7,650	
Edges	65,854	186,100	13,584	65,718	30,019	479,476	343,486	431,206	88,651	238,162	
Classes	18	5	5	5	5	9	6	10	3	8	
MLP	62.29 \pm 1.03	42.66 \pm 0.84	38.66 \pm 4.02	36.74 \pm 1.80	36.70 \pm 0.85	89.82 \pm 0.63	93.57 \pm 0.55	78.94 \pm 1.22	87.48 \pm 0.46	89.96 \pm 1.22	11
GCN	38.58 \pm 2.35	45.16 \pm 0.49	42.12 \pm 3.82	38.47 \pm 1.82	30.11 \pm 0.74	68.25 \pm 2.75	78.15 \pm 0.95	77.53 \pm 1.41	87.70 \pm 0.32	94.31 \pm 0.33	10.8
GAT	59.55 \pm 1.45	46.90 \pm 0.47	40.89 \pm 3.50	38.22 \pm 1.71	30.94 \pm 0.95	57.22 \pm 3.04	88.36 \pm 1.37	76.69 \pm 0.87	87.45 \pm 0.53	94.59 \pm 0.48	11.4
GCNII	82.53 \pm 0.37	47.53 \pm 0.72	41.56 \pm 4.15	40.70 \pm 1.80	37.51 \pm 0.92	<u>91.64 \pm 0.67</u>	<u>96.48 \pm 0.62</u>	84.63 \pm 0.66	89.96 \pm 0.43	95.18 \pm 0.39	4.1
H2GCN	68.61 \pm 1.05	37.20 \pm 0.67	42.29 \pm 4.57	35.82 \pm 2.20	33.32 \pm 0.90	91.25 \pm 0.58	96.24 \pm 0.39	78.34 \pm 2.01	89.32 \pm 0.37	<u>95.66 \pm 0.26</u>	8.2
MixHop	79.16 \pm 0.70	47.95 \pm 0.65	44.97 \pm 3.12	40.43 \pm 1.40	36.97 \pm 0.90	91.10 \pm 0.46	96.21 \pm 0.42	84.19 \pm 0.61	89.42 \pm 0.37	95.63 \pm 0.30	4.7
GBK-GNN	66.05 \pm 1.44	40.20 \pm 1.96	42.01 \pm 4.89	36.52 \pm 1.45	35.70 \pm 1.12	OOM	OOM	81.07 \pm 0.83	88.18 \pm 0.45	93.48 \pm 0.42	10.7
GGCN	OOM	OOM	41.23 \pm 4.08	36.76 \pm 2.19	35.68 \pm 0.87	90.84 \pm 0.65	95.58 \pm 0.44	<u>84.76 \pm 0.65</u>	89.04 \pm 0.40	95.18 \pm 0.44	8.5
GlöGNN	68.63 \pm 0.63	48.62 \pm 0.59	40.95 \pm 5.95	36.85 \pm 1.97	36.66 \pm 0.81	90.47 \pm 0.77	94.51 \pm 0.49	82.83 \pm 0.52	89.60 \pm 0.34	95.09 \pm 0.46	8.2
HöGCGN	OOM	OOM	43.35 \pm 3.66	38.63 \pm 1.95	36.47 \pm 0.83	90.94 \pm 0.72	94.75 \pm 0.65	83.74 \pm 0.69	OOM	94.79 \pm 0.26	7.3
GPR-GNN	71.19 \pm 0.75	46.64 \pm 0.52	41.84 \pm 4.68	38.04 \pm 1.98	36.21 \pm 0.98	91.19 \pm 0.47	96.37 \pm 0.44	84.07 \pm 0.54	89.28 \pm 0.37	95.48 \pm 0.24	6.7
ACM-GCN	71.15 \pm 0.73	50.64 \pm 0.61	45.20 \pm 4.14	40.90 \pm 1.74	35.88 \pm 1.40	91.43 \pm 0.65	96.19 \pm 0.45	84.39 \pm 0.43	89.99 \pm 0.40	95.52 \pm 0.40	4.3
OrderedGNN	83.10 \pm 0.75	51.30 \pm 0.61	42.07 \pm 4.24	<u>37.75 \pm 2.53</u>	<u>37.22 \pm 0.62</u>	91.42 \pm 0.79	96.27 \pm 0.73	85.50 \pm 0.80	90.09 \pm 0.37	95.73 \pm 0.33	3.3
CMGNN	84.35 \pm 1.27	52.13 \pm 0.55	45.70 \pm 4.92	41.89 \pm 2.34	36.82 \pm 0.78	92.66 \pm 0.46	97.00 \pm 0.52	84.50 \pm 0.73	89.99 \pm 0.32	95.48 \pm 0.29	2.1

236 **Objective Function.** As mentioned in Sec 4, the CMs in real-world graphs don't always have
 237 significant discriminability, which may lead to low effectiveness of supplementary messages. Thus, we
 238 introduce an additional discrimination loss \mathcal{L}_{dis} to reduce the similarity of the desired neighborhood
 239 message among different classes, which enhances the discriminability among classes in CM. The
 240 overall loss consists of a CrossEntropy loss \mathcal{L}_{ce} and the discrimination loss \mathcal{L}_{dis} :

$$\mathcal{L} = \mathcal{L}_{ce} + \lambda \mathcal{L}_{dis}, \quad \mathcal{L}_{dis} = \sum_{i \neq j} \text{Sim}(\hat{\mathbf{M}}_i \mathbf{Z}_{ptt}, \hat{\mathbf{M}}_j \mathbf{Z}_{ptt}), \quad (11)$$

241 where $\mathbf{Z}_{ptt} \in \mathbb{R}^{K \times d_r}$ is the representation of virtual prototypes nodes. More details about the
 242 implementation of CMGNN is available in Appendix E.

243 6 Benchmarks and Experiments

244 In this section, we conduct comprehensive experiments to demonstrate the effectiveness of the
 245 proposed CMGNN with a newly organized benchmark for fair comparisons.

246 6.1 New Benchmark

247 As reported in [30], some widely adopted datasets in existing works have critical drawbacks, which
 248 lead to unreliable results. Therefore, with a comprehensive review of existing benchmark evaluation,
 249 we construct a new benchmark to fairly perform experimental validation. Specifically, we integrate 13
 250 representative homophilous and heterophilous GNNs, construct a unified codebase, and evaluate their
 251 node classification performances on 10 unified organized datasets with various heterophily levels.

252 **Drawbacks of Existing Datasets.** Existing works mostly follow the settings and datasets used
 253 in [15], including 6 heterophilous datasets (Cornell, Texas, Wisconsin, Actor, Chameleon, and
 254 Squirrel) and 3 homophilous datasets (Cora, Citeseer, and Pubmed). Platonov et al. [30] pointed out
 255 that there are serious data leakages in Chameleon and Squirrel, while Cornell, Texas, and Wisconsin
 256 are too small with very imbalanced classes. Further, we revisit other datasets and discover new
 257 drawbacks: (i) In the ten splits of Citeseer, there are two inconsistent ones, which have smaller
 258 training, validation, and test sets that could cause issues with statistical results; (ii) The data split
 259 ratios for Cora are not consistent with the expected ones. These drawbacks may lead to certain issues
 260 with the conclusions of previous works. The detailed descriptions of dataset drawbacks are listed in
 261 Appendix F.1.

262 **Newly Organized Datasets.** The datasets used in the benchmark include Roman-Empire, Amazon-
 263 Ratings, Chameleon-F, Squirrel-F, Actor, Flickr, BlogCatalog, Wikics, Pubmed, and Photo. Their
 264 statistics are summarized in Table 2, with details in Appendix F.2. For consistency with existing meth-
 265 ods, we randomly construct 10 splits with predefined proportions (48%/32%/20% for train/valid/test)
 266 for each dataset and report the mean performance and standard deviation of 10 splits.

Table 3: Ablation study results (%) between CMGNN and three ablation variants, where SM denotes supplementary messages of the desired neighborhoods and DL denotes the discrimination loss.

Variants	Roman-Empire	Amazon-Ratings	Chameleon-F	Squirrel-F	Actor	Flickr	BlogCatalog	Wikics	Pubmed	Photo
CMGNN	84.35 ± 1.27	52.13 ± 0.55	45.70 ± 4.92	41.89 ± 2.34	36.82 ± 0.78	92.66 ± 0.46	97.00 ± 0.52	84.50 ± 0.73	89.99 ± 0.32	95.48 ± 0.29
W/O SM	83.84 ± 1.09	51.98 ± 0.61	42.35 ± 4.21	40.79 ± 1.89	36.02 ± 1.21	92.32 ± 0.83	96.52 ± 0.63	83.97 ± 0.83	89.70 ± 0.44	95.41 ± 0.40
W/O DL	83.68 ± 1.24	52.04 ± 0.37	44.97 ± 3.99	41.60 ± 2.43	36.28 ± 1.12	92.66 ± 0.46	97.00 ± 0.52	83.29 ± 1.83	89.99 ± 0.32	95.26 ± 0.35
W/O SM and DL	83.52 ± 1.91	51.58 ± 1.04	41.12 ± 2.93	40.07 ± 2.41	35.61 ± 1.48	92.32 ± 0.83	96.52 ± 0.63	81.62 ± 1.67	89.70 ± 0.44	94.66 ± 0.42

267 **Baseline Methods.** As baseline methods, we choose 13 representative homophilous and het-
 268 erophilous GNNs, including (i) shallow base model: MLP; (ii) homophilous GNNs: GCN [1],
 269 GAT [28], GCNII [27]; (iii) heterophilous GNNs: H2GCN [12], MixHop [16], GBK-GNN [24],
 270 GGCN [23], GloGNN [19], HOGGCN [17], GPR-GNN [20], ACM-GCN [18] and OrderedGNN [21].
 271 For each method, we integrate its official/reproduced code into a unified codebase and search for
 272 parameters in the space suggested by the original papers. More experimental settings can be found in
 273 Appendix F.4 and G.1.

274 6.2 Main Results

275 Following the constructed benchmark, we evaluate methods and report the performance in Table 2.

276 **Performance of Baseline Methods.** With the new benchmarks, some interesting observations and
 277 conclusions can be found when analyzing the performance of baseline methods. First, comparing the
 278 performance of MLP and GCN, we can find "good" heterophily in Amazon-Ratings, Chameleon-F,
 279 and Squirrel-F. Meanwhile, when the homophily level is not high enough, "bad" homophily may also
 280 exist as shown in BlogCatalog and Wikics. These results once again support the observations about
 281 CMs. Therefore, **homophilous GNNs** can also work well in heterophilous graphs as GCNII has
 282 an average rank of 4.1, which is better than most HTGNNs. This is attributed to the initial residual
 283 connection in GCNII actually playing the role of ego/neighbor separation, which is suitable in
 284 heterophilous graphs. As for **heterophilous GNNs**, they are usually designed for both homophilous
 285 and heterophilous graphs. Surprisingly, MixHop, as an early method, demonstrated quite good
 286 performance. In fact, from the perspective of HTMP, it can be considered a degenerate version
 287 of OrderedGNN with no learnable dimensions. As previous SOTA methods, OrderedGNN and
 288 ACM-GCN prove their strong capabilities again.

289 **Performance of CMGNN.** CMGNN achieves the best performance in 6 datasets and an average
 290 rank of 2.1, which outperforms baseline methods. This demonstrates the superiority of utilizing
 291 and enhancing the CM to handle incomplete and noisy semantic neighborhoods, especially in
 292 heterophilous graphs. Regarding the suboptimal performance in Actor, we believe that this is due
 293 to the CM in this dataset are not discriminative enough to provide valuable information via the
 294 supplementary messages and hard to enhance. In homophilous graphs, due to the identity-like CMs,
 295 the overlap between distributions is relatively less, leading to a minor contribution from supplement
 296 messages. Yet CMGNN still achieves top-level performances.

297 6.3 Ablation Study

298 We conduct an ablation study on two key designs of CMGNN, including the supplementary messages
 299 of the desired neighborhood (SM) and the discrimination loss (DL). The results are shown in Table 3.
 300 *First of all*, both SM and DL have indispensable contributions except for Flickr, BlogCatalog, and
 301 Pubmed, in which the discrimination loss has no effect. This may be due to the discriminability of
 302 desired neighborhood messages reaching the bottlenecks and can not be further improved by DL.
 303 *Meanwhile*, the extent of their contributions varies across datasets. SM plays a more important role in
 304 most datasets except Roman-Empire, Wikics, and Photo, in which the number of nodes that need
 305 supplementary messages is relatively small and DL has great effects. **Further**, we notice that with
 306 SM and DL, CMGNN can reach a smaller standard deviation most of the time. This illustrates
 307 that CMGNN achieves more stable results by handling nodes with incomplete and noisy semantic
 308 neighborhoods. As for the opposite result on Chameleon-F, this may attributed to the small size of
 309 this dataset (890 nodes), which can lead to naturally unstable results.

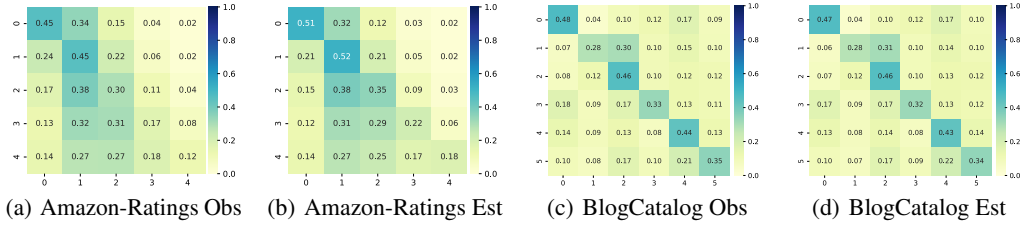


Figure 2: The visualization of observed (Obs) and estimated (Est) compatibility matrixes.

Table 4: Node classification accuracy (%) comparison among nodes with different degrees.

Dataset Deg. Prop.(%)	Amazon-Ratings					Flickr					BlogCatalog				
	0~20	20~40	40~60	60~80	80~100	0~20	20~40	40~60	60~80	80~100	0~20	20~40	40~60	60~80	80~100
CMGNN	59.78	58.36	53.08	41.74	<u>47.86</u>	92.56	91.19	<u>92.71</u>	93.24	<u>93.65</u>	94.13	97.17	98.29	97.99	97.47
ACM-GCN	<u>57.35</u>	<u>56.21</u>	<u>51.74</u>	41.55	46.47	<u>90.44</u>	<u>91.17</u>	92.85	<u>93.19</u>	89.50	92.17	96.68	<u>97.83</u>	97.84	96.51
OrderedGNN	56.32	56.16	51.20	41.85	50.26	86.48	90.07	92.40	92.79	93.40	92.19	96.09	97.48	97.36	96.27
GCNII	50.61	49.94	47.49	41.85	47.76	87.49	90.54	92.29	92.68	95.09	<u>92.81</u>	<u>96.73</u>	97.58	<u>97.90</u>	<u>97.43</u>

310 6.4 Visualization of Compatibility Matrix Estimation

311 We visualize the observed and estimated CMs by CMGNN in Figure 2 with heat maps. Obviously,
 312 CMGNN estimates CMs that are very close to those existing in graphs. This shows that even
 313 with incomplete node labels, CMGNN can estimate high-quality CMs which provides valuable
 314 neighborhood information to nodes. Meanwhile, it can adapt to graphs with various levels of
 315 heterophily. More results can be seen in Appendix G.2.1.

316 6.5 Performance on Nodes with Various Levels of Degrees

317 To verify the effect of CMGNN on nodes with incomplete and noisy semantic neighborhoods, we
 318 divide the test set nodes into 5 parts according to their degrees and report the classification accuracy
 319 respectively. We compare CMGNN with 3 top-performance methods and show the results in Table 4.
 320 In general, nodes with low degrees tend to have incomplete and noisy semantic neighborhoods.
 321 Thus, our outstanding performances on the top 20% nodes with the least degree demonstrate the
 322 effectiveness of CMGNN for providing desired neighborhood messages. Further, we can find that
 323 OrderedGNN and GCNII are good at dealing with nodes with high degrees, while ACM-GCN is
 324 relatively good at nodes with low degrees. And CMGNN, to a certain extent, can be adapted to both
 325 situations at the same time.

326 7 Conclusion and Limitations

327 In this paper, we revisit the message passing mechanism in existing heterophilous GNNs and
 328 reformulate them into a unified heterophilous message passing (HTMP) mechanism. Based on the
 329 HTMP mechanism and empirical analysis, we reveal that the reason for message passing remaining
 330 effective is attributed to implicitly enhancing the compatibility matrix among classes. Further, we
 331 propose a novel method CMGNN to unlock the potential of the compatibility matrix by handling the
 332 incomplete and noisy semantic neighborhoods. The experimental results show the effectiveness of
 333 CMGNN and the feasibility of designing a new method following HTMP mechanism. We hope the
 334 HTMP mechanism and benchmark can further provide convenience to the community.

335 This work mainly focuses on the message passing mechanism in existing HTGNNs under the
 336 semi-supervised setting. Thus, the other designs in HTGNNs such as objective functions are not
 337 analyzed in this paper. The proposed HTMP mechanism is suitable for only a large part of existing
 338 HTGNNs which still follow the message passing mechanism.

References

- 339
- 340 [1] Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional
341 networks. In International Conference on Learning Representations, 2017.
- 342 [2] Yanfu Zhang, Hongchang Gao, Jian Pei, and Heng Huang. Robust self-supervised struc-
343 tural graph neural network for social network prediction. In Proceedings of the ACM Web
344 Conference 2022, pages 1352–1361, 2022.
- 345 [3] Xiang Wang, Xiangnan He, Meng Wang, Fuli Feng, and Tat-Seng Chua. Neural graph collabo-
346 rative filtering. In Proceedings of the 42nd international ACM SIGIR conference on Research
347 and development in Information Retrieval, pages 165–174, 2019.
- 348 [4] Xiangnan He, Kuan Deng, Xiang Wang, Yan Li, Yongdong Zhang, and Meng Wang. Lightgcn:
349 Simplifying and powering graph convolution network for recommendation. In Proceedings
350 of the 43rd International ACM SIGIR conference on research and development in Information
351 Retrieval, pages 639–648, 2020.
- 352 [5] Miller McPherson, Lynn Smith-Lovin, and James M Cook. Birds of a feather: Homophily in
353 social networks. Annual review of sociology, 27(1):415–444, 2001.
- 354 [6] Yao Ma, Xiaorui Liu, Neil Shah, and Jiliang Tang. Is homophily a necessity for graph neural
355 networks? In International Conference on Learning Representations, 2022.
- 356 [7] Jiong Zhu, Ryan A Rossi, Anup Rao, Tung Mai, Nedim Lipka, Nesreen K Ahmed, and Danai
357 Koutra. Graph neural networks with heterophily. In Proceedings of the AAAI conference on
358 artificial intelligence, volume 35, pages 11168–11176, 2021.
- 359 [8] Xin Zheng, Yixin Liu, Shirui Pan, Miao Zhang, Di Jin, and Philip S Yu. Graph neural networks
360 for graphs with heterophily: A survey. arXiv preprint arXiv:2202.07082, 2022.
- 361 [9] Jiong Zhu, Yujun Yan, Mark Heimann, Lingxiao Zhao, Leman Akoglu, and Danai Koutra.
362 Heterophily and graph neural networks: Past, present and future. IEEE Data Engineering
363 Bulletin, 2023.
- 364 [10] Chenghua Gong, Yao Cheng, Xiang Li, Caihua Shan, Siqiang Luo, and Chuan Shi. Towards
365 learning from graphs with heterophily: Progress and future. arXiv preprint arXiv:2401.09769,
366 2024.
- 367 [11] Deyu Bo, Xiao Wang, Chuan Shi, and Huawei Shen. Beyond low-frequency information in graph
368 convolutional networks. In Proceedings of the AAAI Conference on Artificial Intelligence,
369 volume 35, pages 3950–3957, 2021.
- 370 [12] Jiong Zhu, Yujun Yan, Lingxiao Zhao, Mark Heimann, Leman Akoglu, and Danai Koutra. Be-
371 yond homophily in graph neural networks: Current limitations and effective designs. Advances
372 in neural information processing systems, 33:7793–7804, 2020.
- 373 [13] Di Jin, Zhizhi Yu, Cuiying Huo, Rui Wang, Xiao Wang, Dongxiao He, and Jiawei Han.
374 Universal graph convolutional networks. Advances in Neural Information Processing Systems,
375 34:10654–10664, 2021.
- 376 [14] Wei Jin, Tyler Derr, Yiqi Wang, Yao Ma, Zitao Liu, and Jiliang Tang. Node similarity preserving
377 graph convolutional networks. In Proceedings of the 14th ACM international conference on
378 web search and data mining, pages 148–156, 2021.
- 379 [15] Hongbin Pei, Bingzhe Wei, Kevin Chen-Chuan Chang, Yu Lei, and Bo Yang. Geom-gcn: Geo-
380 metric graph convolutional networks. In International Conference on Learning Representations,
381 2020.
- 382 [16] Sami Abu-El-Haija, Bryan Perozzi, Amol Kapoor, Nazanin Alipourfard, Kristina Lerman, Hrayr
383 Harutyunyan, Greg Ver Steeg, and Aram Galstyan. Mixhop: Higher-order graph convolutional
384 architectures via sparsified neighborhood mixing. In international conference on machine
385 learning, pages 21–29. PMLR, 2019.

- 386 [17] Tao Wang, Di Jin, Rui Wang, Dongxiao He, and Yuxiao Huang. Powerful graph convolutional
387 networks with adaptive propagation mechanism for homophily and heterophily. In Proceedings
388 of the AAAI conference on artificial intelligence, volume 36, pages 4210–4218, 2022.
- 389 [18] Sitao Luan, Chenqing Hua, Qincheng Lu, Jiaqi Zhu, Mingde Zhao, Shuyuan Zhang, Xiao-Wen
390 Chang, and Doina Precup. Revisiting heterophily for graph neural networks. Advances in
391 neural information processing systems, 35:1362–1375, 2022.
- 392 [19] Xiang Li, Renyu Zhu, Yao Cheng, Caihua Shan, Siqiang Luo, Dongsheng Li, and Weining Qian.
393 Finding global homophily in graph neural networks when meeting heterophily. In International
394 Conference on Machine Learning, pages 13242–13256. PMLR, 2022.
- 395 [20] Eli Chien, Jianhao Peng, Pan Li, and Olgica Milenkovic. Adaptive universal generalized
396 pagerank graph neural network. In International Conference on Learning Representations,
397 2021.
- 398 [21] Yunchong Song, Chenghu Zhou, Xinbing Wang, and Zhouhan Lin. Ordered GNN: Ordering
399 message passing to deal with heterophily and over-smoothing. In The Eleventh International
400 Conference on Learning Representations, 2023.
- 401 [22] Susheel Suresh, Vinith Budde, Jennifer Neville, Pan Li, and Jianzhu Ma. Breaking the limit of
402 graph neural networks by improving the assortativity of graphs with local mixing patterns. In
403 Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining,
404 pages 1541–1551, 2021.
- 405 [23] Yujun Yan, Milad Hashemi, Kevin Swersky, Yaoqing Yang, and Danai Koutra. Two sides of the
406 same coin: Heterophily and oversmoothing in graph convolutional neural networks. In 2022
407 IEEE International Conference on Data Mining (ICDM), pages 1287–1292. IEEE, 2022.
- 408 [24] Lun Du, Xiaozhou Shi, Qiang Fu, Xiaojun Ma, Hengyu Liu, Shi Han, and Dongmei Zhang.
409 Gbk-gnn: Gated bi-kernel graph neural networks for modeling both homophily and heterophily.
410 In Proceedings of the ACM Web Conference 2022, pages 1550–1558, 2022.
- 411 [25] Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. Neural
412 message passing for quantum chemistry. In International conference on machine learning,
413 pages 1263–1272. PMLR, 2017.
- 414 [26] Johannes Gasteiger, Aleksandar Bojchevski, and Stephan Günnemann. Predict then propagate:
415 Graph neural networks meet personalized pagerank. In International Conference on Learning
416 Representations, 2019.
- 417 [27] Ming Chen, Zhewei Wei, Zengfeng Huang, Bolin Ding, and Yaliang Li. Simple and deep graph
418 convolutional networks. In International conference on machine learning, pages 1725–1735.
419 PMLR, 2020.
- 420 [28] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and
421 Yoshua Bengio. Graph attention networks. In The International Conference on Learning
422 Representations, 2018.
- 423 [29] Oleksandr Shchur, Maximilian Mumme, Aleksandar Bojchevski, and Stephan Günnemann.
424 Pitfalls of graph neural network evaluation. arXiv preprint arXiv:1811.05868, 2018.
- 425 [30] Oleg Platonov, Denis Kuznedelev, Michael Diskin, Artem Babenko, and Liudmila
426 Prokhorenkova. A critical look at the evaluation of GNNs under heterophily: Are we re-
427 ally making progress? In The Eleventh International Conference on Learning Representations,
428 2023.
- 429 [31] Keyulu Xu, Chengtao Li, Yonglong Tian, Tomohiro Sonobe, Ken-ichi Kawarabayashi, and
430 Stefanie Jegelka. Representation learning on graphs with jumping knowledge networks. In
431 International conference on machine learning, pages 5453–5462. PMLR, 2018.
- 432 [32] Michaël Defferrard, Xavier Bresson, and Pierre Vandergheynst. Convolutional neural networks
433 on graphs with fast localized spectral filtering. Advances in neural information processing
434 systems, 29, 2016.

- 435 [33] Bingbing Xu, Huawei Shen, Qi Cao, Yunqi Qiu, and Xueqi Cheng. Graph wavelet neural
436 network. In International Conference on Learning Representations, 2018.
- 437 [34] Will Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large
438 graphs. Advances in neural information processing systems, 30, 2017.
- 439 [35] Yao Ma, Xiaorui Liu, Tong Zhao, Yozen Liu, Jiliang Tang, and Neil Shah. A unified view on
440 graph neural networks as graph signal denoising. In Proceedings of the 30th ACM International
441 Conference on Information & Knowledge Management, pages 1202–1211, 2021.
- 442 [36] Meiqi Zhu, Xiao Wang, Chuan Shi, Houye Ji, and Peng Cui. Interpreting and unifying graph
443 neural networks with an optimization framework. In Proceedings of the Web Conference 2021,
444 pages 1215–1226, 2021.
- 445 [37] Yixin Liu, Zhao Li, Shirui Pan, Chen Gong, Chuan Zhou, and George Karypis. Anomaly
446 detection on attributed networks via contrastive self-supervised learning. IEEE Transactions on
447 Neural Networks and Learning Systems, 2021.
- 448 [38] Péter Mernyei and Cătălina Cangea. Wiki-cs: A wikipedia-based benchmark for graph neural
449 networks. arXiv preprint arXiv:2007.02901, 2020.
- 450 [39] Prithviraj Sen, Galileo Namata, Mustafa Bilgic, Lise Getoor, Brian Galligher, and Tina Eliassi-
451 Rad. Collective classification in network data. AI magazine, 29(3):93–93, 2008.
- 452 [40] Édouard Grave, Piotr Bojanowski, Prakhar Gupta, Armand Joulin, and Tomáš Mikolov. Learning
453 word vectors for 157 languages. In Proceedings of the Eleventh International Conference on
454 Language Resources and Evaluation (LREC 2018), 2018.
- 455 [41] Leskovec Jure. Snap datasets: Stanford large network dataset collection. Retrieved December
456 2021 from <http://snap.stanford.edu/data>, 2014.
- 457 [42] Benedek Rozemberczki, Carl Allen, and Rik Sarkar. Multi-scale attributed node embedding.
458 Journal of Complex Networks, 9(2):cnab014, 2021.
- 459 [43] Jie Tang, Jimeng Sun, Chi Wang, and Zi Yang. Social influence analysis in large-scale networks.
460 In Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery
461 and data mining, pages 807–816, 2009.
- 462 [44] Jeffrey Pennington, Richard Socher, and Christopher D Manning. Glove: Global vectors for
463 word representation. In Proceedings of the 2014 conference on empirical methods in natural
464 language processing (EMNLP), pages 1532–1543, 2014.

465 A More Details of HTMP Mechanism

466 In this part, we list more details about the HTMP mechanism, including additional analysis about
467 HTMP, method-wise analysis and overall analysis.

468 A.1 Additional Analysis of HTMP Mechanism

469 A.1.1 Neighborhood Indicators

470 The neighborhood indicator explicitly marks the neighbors of all nodes within a specific neighbor-
471 hood. In existing heterophilous GNNs, neighborhood indicators typically take one of the following
472 forms: (i) Raw Neighborhood (Raw); (ii) Neighborhood Redefining (ReDef); and (3) Neighborhood
473 Discrimination (Dis).

474 **Raw Neighborhood.** Raw neighborhood, including \mathbf{A} and $\tilde{\mathbf{A}}$, provides the basic neighborhood
475 information. The only difference between them lies in whether there is differential treatment of the
476 node’s ego messages. For example, APPNP [26] applies additional weighting to the nodes’ ego
477 messages compared with GCN [1]. For the sake of simplicity, we consider the identity matrix $\mathbf{I} \in$
478 $\mathbb{R}^{N \times N}$ as a special neighborhood indicator for acquiring the nodes’ ego messages. In heterophilous
479 GNNs, ego/neighbor separation is a common strategy that can mitigate the confusion of ego messages
480 with neighbor messages.

481 **Neighborhood Redefining.** Neighborhood redefining is the most commonly used technique in
482 heterophilous GNNs, aiming to capture additional information from new neighborhoods. As a repre-
483 sentative example, *high-order neighborhood* \mathbf{A}_h can provide long-distance connection information
484 but also result in additional computational costs. *Feature-similarity-based neighborhood* \mathbf{A}_f is often
485 defined by the k-NN relationships within the feature space. Fundamentally, it only utilizes node
486 features and thus needs to be used in conjunction with other neighborhood indicators. Otherwise,
487 the model will be limited by the amount of information in node features. GloGNN [19] introduces
488 *fully-connected neighborhood* $\mathbf{1} \in \mathbb{R}^{N \times N}$, which can capture global neighbor information from all
489 nodes. However, it can also cause significant time and space consumption. Additionally, there are
490 some *custom-defined neighborhood* \mathbf{A}_c . For example, Geom-GCN [15] redefines neighborhoods
491 based on the geometric relationships between node pairs. These neighborhood indicators may have
492 limited generality, and the effectiveness is reliant on the specific method.

493 **Neighborhood Discrimination.** Neighborhood discrimination aims to mark whether neighbors
494 share the same label with central nodes. The neighborhoods are partitioned into positive \mathbf{A}_p and
495 negative ones \mathbf{A}_n , which include homophilous and heterophilous neighbors respectively. GGCN [23]
496 divides the raw neighborhood based on the similarity of node representations with a threshold
497 of 0. Explicitly distinguishing neighbors allows for targeted processing, making the model more
498 interpretable. However, its performance is influenced by the accuracy of the discrimination, which
499 may lead to the accumulation of errors.

500 A.1.2 Aggregation Guidance

501 After identifying the neighborhood, the aggregation guidance controls what type of messages to
502 gather from the corresponding neighbors. The existing aggregation guidance mainly includes three
503 kinds of approaches: (1) Degree Averaging (DegAvg), (2) Adaptive Weights (AdaWeight), and (3)
504 Relationship Estimation (RelaEst).

505 **Degree Averaging.** Degree averaging, formatted as $\mathbf{B}^d = \mathbf{D}^{-\frac{1}{2}} \mathbf{1} \mathbf{D}^{-\frac{1}{2}}$ or $\mathbf{B}^d = \mathbf{D}^{-1} \mathbf{1}$, is the most
506 common aggregation guidance, which plays the role of a low-pass filter to capture the smooth signals
507 and is fixed during model training. Further, combining negative degree averaging with an identity
508 aggregation guidance $\mathbf{I} \in \mathbb{R}^{N \times N}$ can capture the difference between central nodes and neighbors, as
509 used in ACM-GCN [18]. Degree averaging is simple and efficient but depends on the discriminability
510 of corresponding neighborhoods.

511 **Adaptive Weights.** Another common strategy is allowing the model to learn the appropriate aggrega-
512 tion guidances \mathbf{B}^{aw} . GAT [28] proposes an attention mechanism to learn aggregate weights, which
513 guides many subsequent heterophilous methods. To better handle heterophilous graphs, FAGCN [11]
514 introduces negative-available attention weights \mathbf{B}^{naw} to capture the difference between central nodes

515 and heterophilous neighbors. Adaptive weights can personalize message aggregation for different
 516 neighbors, yet it’s difficult for models to attain the desired effect.

517 **Relationship Estimation.** Recently, some methods have tried to estimate the pair-wise relationships
 518 \mathbf{B}^{r^e} between nodes and use them to guide message aggregation. HOG-GCN [17] estimates the
 519 pair-wise homophily levels between nodes as aggregation guidances based on both attribute and
 520 topology space. GloGNN [19] treats all nodes as neighbors and estimates a coefficient matrix
 521 as aggregation guidance based on the idea of linear subspace expression. GGCN [23] estimates
 522 appropriate weights for message aggregation with the degrees of nodes and the similarities between
 523 node representations. Relationship estimation usually has theoretical guidance, which brings strong
 524 interpretability. However, it may also result in significant temporal and spatial complexity when
 525 estimating pair-wise relations.

526 A.1.3 COMBINE Function

527 After message aggregation, the COMBINE functions integrate messages from multiple neighborhoods
 528 into layer representations. COMBINE functions in heterophilous GNNs are commonly based on
 529 two operations: addition and concatenation, each of which has variants. To merge several messages
 530 together, addition (Add) is a naive idea. Further, to control the weight of messages from different
 531 neighborhoods, weighted addition (WeightedAdd) is applied. However, it is a global setting and
 532 cannot adapt to the differences between nodes. Thus, adaptive weighted addition (AdaAdd) is
 533 proposed, which can learn personalized message combination weights for each node, but it will result
 534 in additional time consumption. Although the addition is simple and efficient, some methods [12, 16]
 535 believe that it may blur messages from different neighborhoods, which can be harmful in heterophilous
 536 GNNs, so they employ a concatenation operation (Cat) to separate the messages. Nevertheless, such
 537 an approach not only increases the space cost but may also retain additional redundant messages. To
 538 address these issues, OrderedGNN [21] proposes an adaptive concatenation mechanism (AdaCat)
 539 that can combine multiple messages with learnable dimensions. This is an innovative and worthy
 540 further exploration practice, but the difficulty of model learning should also be considered.

541 A.1.4 FUSE Function

542 Further, the FUSE functions integrate messages from multiple layers into the final representation. For
 543 the FUSE function, utilizing the representation of the last layer as the final representation is widely
 544 accepted: $\mathbf{Z} = \mathbf{Z}^L$. JKNet [31] proposes that the combination of representations from intermediate
 545 layers can capture both local and global information. H2GCN [12] applies it in heterophilous graphs,
 546 preserving messages from different localities with concatenation. Similarly, GPRGNN [20] combines
 547 the representations of multiple layers into the final representation through adaptive weighted addition.

548 A.1.5 AGGREGATE function

549 The most commonly used AGGREGATE function is $\text{AGGREGATE}(\mathbf{A}_r, \mathbf{B}_r, \mathbf{Z}_r^{l-1}) = (\mathbf{A}_r \odot$
 550 $\mathbf{B}_r) \mathbf{Z}_r^{l-1} \mathbf{W}_r^l$. We take this as the fixed form of the AGGREGATE function following. Actually,
 551 the input representations \mathbf{Z}_r^{l-1} and weight matrixes \mathbf{W}_r^l also can be specially designed. Taking
 552 the initial node representations \mathbf{Z}^0 as input is a relatively common approach as in APPNP [26],
 553 GCNII [27], FAGCN [11] and GloGNN [19]. Further, GCNII [27] adds an identity matrix \mathbf{I}_w to the
 554 weight matrixes to keep more original messages. However, the methods that specially design these
 555 components are few and with a similar form. Thus, we don’t discuss them too much, but leave it for
 556 future extensions.

557 A.2 Revisiting Representative GNNs with HTMP Mechanism

558 In this part, we utilize HTMP mechanism to revisit the representative GNNs. We start from ho-
 559 mophilous GNNs as simple examples and further extend to heterophilous GNNs.

560 A.2.1 GCN

561 Graph Convolutional Networks (GCN) [1] utilizes a low-pass filter to gather messages from neighbors
 562 as follows:

$$\mathbf{Z}^l = \hat{\mathbf{A}} \mathbf{Z}^{l-1} \mathbf{W}^l. \quad (12)$$

563 It can be revisited by HTMP with the following components:

$$\begin{aligned} \mathbf{A}_0 &= \tilde{\mathbf{A}}, \quad \mathbf{B}_0 = \mathbf{B}^d = \tilde{\mathbf{D}}^{-\frac{1}{2}} \mathbf{1} \tilde{\mathbf{D}}^{-\frac{1}{2}}, \\ \mathbf{Z}^l &= \mathbf{Z}_0^l = (\mathbf{A}_0 \odot \mathbf{B}_0) \mathbf{Z}^{l-1} \mathbf{W}^l = \hat{\mathbf{A}} \mathbf{Z}^{l-1} \mathbf{W}^l. \end{aligned} \quad (13)$$

564 Specifically, GCN has a raw neighborhood indicator $\tilde{\mathbf{A}}$ and a degree averaging aggregation guidance
565 \mathbf{B}^d . Since there is only one neighborhood, the COMBINE function is meaningless in GCN. GCN
566 utilizes a naive way to fuse messages about the original neighborhood and central nodes. However, it
567 may confuse the representations in heterophilous graphs.

568 A.2.2 APPNP

569 PPNP [26] is also a general method whose message passing is based on Personalized PageRank
570 (PPR). To avoid massive consumption, APPNP is introduced as the approximate version of PPNP
571 with an iterative message-passing mechanism:

$$\mathbf{Z}^l = \mu \mathbf{Z}^0 + (1 - \mu) \hat{\mathbf{A}} \mathbf{Z}^{l-1}. \quad (14)$$

572 It can be revisited by $_$ with the following components:

$$\begin{aligned} \mathcal{A} &= [\mathbf{A}_0, \mathbf{A}_1], \quad \mathcal{B} = [\mathbf{B}_0, \mathbf{B}_1], \\ \mathbf{A}_0 &= \mathbf{I}, \quad \mathbf{B}_0 = \mathbf{I}, \quad \mathbf{W}_0^l = \mathbf{I}, \\ \tilde{\mathbf{Z}}_0^l &= (\mathbf{A}_0 \odot \mathbf{B}_0) \mathbf{Z}^0 \mathbf{W}_0^l = \mathbf{Z}^0, \\ \mathbf{A}_1 &= \mathbf{A}, \quad \mathbf{B}_1 = \mathbf{D}^{-\frac{1}{2}} \mathbf{1} \mathbf{D}^{-\frac{1}{2}}, \quad \mathbf{W}_1^l = \mathbf{I}, \\ \tilde{\mathbf{Z}}_1^l &= (\mathbf{A}_1 \odot \mathbf{B}_1) \mathbf{Z}^{l-1} \mathbf{W}_1^l = \hat{\mathbf{A}} \mathbf{Z}^{l-1}. \end{aligned} \quad (15)$$

573 Specifically, APPNP aggregates messages from node ego and neighborhoods separately and combines
574 them with a weighted addition. Compared with GCN, APPNP assigns adjustable weights to nodes,
575 for controlling the proportion of ego and neighbor messages during message-passing, which becomes
576 a worthy design in heterophilous graphs.

577 A.2.3 GAT

578 Going a step further, Graph Attention Networks (GAT) [28] allows learnable weights for each
579 neighbor:

$$\mathbf{z}_i^l = \sum_{j \in \tilde{\mathcal{N}}(i)} \alpha_{ij} \mathbf{z}_j^{l-1} \mathbf{W}^l, \quad (16)$$

580 where α_{ij} is the weight for aggregating neighbor node j to center node i , whose construction process
581 is as follows:

$$\begin{aligned} \alpha_{ij} &= \frac{\exp(e_{ij})}{\sum_{k \in \tilde{\mathcal{N}}(i)} \exp(e_{ik})}, \\ e_{ij} &= \text{LeakyReLU}([\mathbf{z}_i^{l-1} | \mathbf{z}_j^{l-1}] \mathbf{a}). \end{aligned} \quad (17)$$

582 Let \mathbf{P}^{GAT} be the matrix of aggregation weights in GAT:

$$\mathbf{P}_{ij}^{GAT} = \begin{cases} \alpha_{ij}, & \tilde{\mathbf{A}}_{ij} = 1, \\ 0, & \tilde{\mathbf{A}}_{ij} = 0. \end{cases} \quad (18)$$

583 HTMP can revisit GAT with the following components:

$$\begin{aligned} \mathbf{A}_0 &= \tilde{\mathbf{A}}, \quad \mathbf{B}_0 = \mathbf{B}^{aw} = \mathbf{P}^{GAT}, \\ \mathbf{Z}^l &= \mathbf{Z}_0^l = (\mathbf{A}_0 \odot \mathbf{B}_0) \mathbf{Z}^{l-1} \mathbf{W}^l = \mathbf{P}^{GAT} \mathbf{Z}^{l-1} \mathbf{W}^l, \end{aligned} \quad (19)$$

584 which is the matrix version of Eq 16. Specifically, GAT aggregate messages from raw neighborhood
585 $\tilde{\mathbf{A}}$ with adaptive weights \mathbf{B}^{aw} . Aggregation guidance with adaptive weights is a nice idea, but simple
586 constraints are not enough for the model to learn ideal results.

587 **A.2.4 GCNII**

588 GCNII [27] is a novel homophilous GNN with two key designs: initial residual connection and
589 identity mapping, which can be formatted as follows:

$$\mathbf{Z}^l = \left(\alpha \mathbf{Z}^0 + (1 - \alpha) \tilde{\mathbf{D}}^{-\frac{1}{2}} \tilde{\mathbf{A}} \tilde{\mathbf{D}}^{-\frac{1}{2}} \mathbf{Z}^{l-1} \right) (\beta \mathbf{W}^l + (1 - \beta) \mathbf{I}_w), \quad (20)$$

590 where α and β are two predefined parameters and $\mathbf{I}_w \in \mathbb{R}^{d_r \times d_r}$ is an identity matrix.

591 From the perspective of HTMP, it can be viewed as follows:

$$\begin{aligned} \mathcal{A} &= [\mathbf{I}, \tilde{\mathbf{A}}], \quad \mathcal{B} = [\mathbf{I}, \tilde{\mathbf{B}}^d], \quad \mathbf{W}_0^l = \mathbf{W}_1^l = (\beta \mathbf{W}^l + (1 - \beta) \mathbf{I}_w), \\ \tilde{\mathbf{Z}}_0^l &= (\mathbf{I} \odot \mathbf{I}) \mathbf{Z}^0 (\beta \mathbf{W}^l + (1 - \beta) \mathbf{I}_w) = \mathbf{Z}^0 (\beta \mathbf{W}^l + (1 - \beta) \mathbf{I}_w), \\ \tilde{\mathbf{Z}}_1^l &= (\tilde{\mathbf{A}} \odot \tilde{\mathbf{B}}^d) \mathbf{Z}^{l-1} (\beta \mathbf{W}^l + (1 - \beta) \mathbf{I}_w) = \hat{\mathbf{A}} \mathbf{Z}^{l-1} (\beta \mathbf{W}^l + (1 - \beta) \mathbf{I}_w), \end{aligned} \quad (21)$$

592 where the COMBINE function is weighted addition. Specifically, the first design of GCNII is a form
593 of ego/neighbor separation, and the second design is a novel transformation weights matrix. This can
594 also be specially designed, but only GCNII does this, so we won't analyze it too much and leave it as
595 a future extension.

596 **A.2.5 Geom-GCN**

597 Geom-GCN [15] is one of the most influential heterophilous GNNs, which employs the geometric
598 relationships of nodes within two kinds of neighborhoods to aggregate the messages through bi-level
599 aggregation:

$$\begin{aligned} \mathbf{Z}^l &= \left(\parallel_{i \in \{g, s\}} \parallel_{r \in R} \mathbf{Z}_{i,r}^l \right) \mathbf{W}^l, \\ \mathbf{Z}_{i,r}^l &= \mathbf{D}_{i,r}^{-\frac{1}{2}} \mathbf{A}_{i,r} \mathbf{D}_{i,r}^{-\frac{1}{2}} \mathbf{Z}^{l-1}, \end{aligned} \quad (22)$$

600 where \parallel denotes the concatenate operator, $\{g, s\}$ is the set of neighborhoods including the original
601 graph and the latent space. R is the set of geometric relationships. $\mathbf{A}_{i,r}$ is the corresponding adjacency
602 matrix in neighborhood i and relationship r .

603 It can be revisited by HTMP with the following components:

$$\begin{aligned} \mathcal{A} &= [\mathbf{A}_{i,r} | i \in \{g, s\}, r \in R], \quad \mathcal{B} = [\mathbf{B}_{i,r}^d | i \in \{g, s\}, r \in R], \\ \tilde{\mathbf{Z}}_{i,r}^l &= (\mathbf{A}_{i,r} \odot \mathbf{B}_{i,r}^d) \mathbf{Z}_{i,r}^{l-1} \mathbf{W}_{i,r}^l = \mathbf{D}_{i,r}^{-\frac{1}{2}} \mathbf{A}_{i,r} \mathbf{D}_{i,r}^{-\frac{1}{2}} \mathbf{Z}^{l-1} \mathbf{W}_{i,r}^l, \end{aligned} \quad (23)$$

604 where the COMBINE function is concatenation and the weight matrix \mathbf{W}^l in Eq 22 can be viewed as
605 the combination of multiple $\mathbf{W}_{i,r}^l$. Specifically, Geom-GCN redefines multiple neighborhoods based
606 on the customized geometric relations in both raw and latent space. The messages are aggregated
607 from each neighborhood and combined by a concatenation. This approach may be applicable to some
608 datasets, yet it has weak universality.

609 **A.2.6 H2GCN**

610 H2GCN [12] is also an influential method with three key designs: ego- and neighbor-message
611 separation, higher-order neighborhoods, and the combination of intermediate representations. Its
612 single-layer representations are constructed as follows:

$$\mathbf{Z}^l = \left[\hat{\mathbf{A}} \mathbf{Z}^{l-1} \parallel \hat{\mathbf{A}}_{h2} \mathbf{Z}^{l-1} \right], \quad (24)$$

613 where $\hat{\mathbf{A}}_{h2}$ denotes the 2-order adjacency matrix with normalization.

614 It can be revisited by HTMP with the following components:

$$\begin{aligned} \mathcal{A} &= [\mathbf{A}, \mathbf{A}_{h2}], \quad \mathcal{B} = [\mathbf{B}^d, \mathbf{B}_{h2}^d], \quad \mathbf{W}_0^l = \mathbf{W}_1^l = \mathbf{I}, \\ \tilde{\mathbf{Z}}_0^l &= (\mathbf{A} \odot \mathbf{B}^d) \mathbf{Z}^{l-1} \mathbf{I} = \hat{\mathbf{A}} \mathbf{Z}^{l-1}, \\ \tilde{\mathbf{Z}}_1^l &= (\mathbf{A}_{h2} \odot \mathbf{B}_{h2}^d) \mathbf{Z}^{l-1} \mathbf{I} = \hat{\mathbf{A}}_{h2} \mathbf{Z}^{l-1}, \end{aligned} \quad (25)$$

615 where the COMBINE function is concatenation. Meanwhile, H2GCN also uses the concatenation
616 as the FUSE function. Specifically, H2GCN aggregates messages from the raw and 2-order neigh-
617 borhoods in a layer of message passing and keeps them apart in the representations. The design
618 of ego/neighbor separation is first introduced by H2GCN and gradually becomes a necessity for
619 subsequent methods.

620 A.2.7 SimP-GCN

621 SimP-GCN [14] constructs an additional graph based on the feature similarity. It has two key concepts:
622 (1) the information from the original graph and feature kNN graph should be balanced, and (2) each
623 node can adjust the contribution of its node features. Specifically, the message passing in SimP-GCN
624 is as follows:

$$\mathbf{Z}^l = \left(\text{diag}(\mathbf{s}^l) \hat{\mathbf{A}} + \text{diag}(1 - \mathbf{s}^l) \hat{\mathbf{A}}_f + \gamma \mathbf{D}_K^l \right) \mathbf{Z}^{l-1} \mathbf{W}^l, \quad (26)$$

625 where $\mathbf{s}^l \in \mathbb{R}^n$ is a learnable score vector that balances the effect of the original and feature graphs,
626 $\mathbf{D}_K^l = \text{diag}(K_1^l, K_2^l, \dots, K_n^l)$ is a learnable diagonal matrix.

627 It can be revisited by HTMP with the following components:

$$\begin{aligned} \mathcal{A} &= [\mathbf{I}, \tilde{\mathbf{A}}, \mathbf{A}_f], \quad \mathcal{B} = [\mathbf{I}, \tilde{\mathbf{B}}^d, \mathbf{B}_f^d], \\ \tilde{\mathbf{Z}}_0^l &= (\mathbf{I} \odot \mathbf{I}) \mathbf{Z}^{l-1} \mathbf{W}^l = \mathbf{Z}^{l-1} \mathbf{W}^l, \\ \tilde{\mathbf{Z}}_1^l &= (\tilde{\mathbf{A}} \odot \tilde{\mathbf{B}}^d) \mathbf{Z}^{l-1} \mathbf{W}^l = \hat{\mathbf{A}} \mathbf{Z}^{l-1} \mathbf{W}^l, \\ \tilde{\mathbf{Z}}_2^l &= (\mathbf{A}_f \odot \mathbf{B}_f^d) \mathbf{Z}^{l-1} \mathbf{W}^l = \hat{\mathbf{A}}_f \mathbf{Z}^{l-1} \mathbf{W}^l, \end{aligned} \quad (27)$$

628 where the COMBINE function is adaptive weighted addition. Specifically, SimP-GCN aggregates
629 messages from ego, raw and feature-similarity-based neighborhoods, and combines them with
630 node-specific learnable weights. The feature-similarity-based neighborhoods can provide more
631 homophilous messages to enhance the discriminability of the compatibility matrix. However, it's still
632 limited by the amount of information on node features.

633 A.2.8 FAGCN

634 FAGCN [11] proposes considering both low-frequency and high-frequency information simultane-
635 ously, and transferring them into the negative-allowable weights during message passing:

$$\mathbf{Z}_i^l = \mu \mathbf{Z}_i^0 + \sum_{j \in \mathcal{N}_i} \frac{\alpha_{ij}^G}{\sqrt{d_i d_j}} \mathbf{Z}_j^{l-1}, \quad (28)$$

636 where α_{ij}^G can be negative as follows:

$$\alpha_{ij}^G = \tanh(\mathbf{g}^T [\mathbf{X}_i \| \mathbf{X}_j]), \quad (29)$$

637 which can form a weight matrix:

$$\mathbf{P}_{ij}^{FAG} = \begin{cases} \alpha_{ij}^G, & \mathbf{A}_{ij} = 1, \\ 0, & \mathbf{A}_{ij} = 0. \end{cases} \quad (30)$$

638 It can be revisited by HTMP with the following components:

$$\begin{aligned} \mathcal{A} &= [\mathbf{I}, \mathbf{A}], \quad \mathcal{B} = [\mathbf{I}, \mathbf{D}^{-\frac{1}{2}} \mathbf{P}^{FAG} \mathbf{D}^{-\frac{1}{2}}], \quad \mathbf{W}_0^l = \mathbf{W}_1^l = \mathbf{I}, \\ \tilde{\mathbf{Z}}_0^l &= (\mathbf{I} \odot \mathbf{I}) \mathbf{Z}^0 \mathbf{I} = \mathbf{Z}^0, \\ \tilde{\mathbf{Z}}_1^l &= (\mathbf{A} \odot \mathbf{D}^{-\frac{1}{2}} \mathbf{P}^{FAG} \mathbf{D}^{-\frac{1}{2}}) \mathbf{Z}^{l-1} \mathbf{I} = \mathbf{D}^{-\frac{1}{2}} \mathbf{P}^{FAG} \mathbf{D}^{-\frac{1}{2}} \mathbf{Z}^{l-1}, \end{aligned} \quad (31)$$

639 where the COMBINE function is weighted addition, same as the matrix form of Eq 28. Specifically,
640 FAGCN aggregates messages from node ego and raw neighborhood with negative-allowable weights.
641 It has a similar form to GAT but allows for ego/neighbor separation and negative weights, which
642 means the model can capture the difference between center nodes and neighbors.

643 A.2.9 GGCN

644 GGCN [23] explicitly distinguishes between homophilous and heterophilous neighbors based on
645 node similarities, and assigns corresponding positive and negative weights:

$$\mathbf{Z}^l = \alpha^l \left(\beta_0^l \hat{\mathbf{Z}}^l + \beta_1^l (\mathbf{S}_{pos}^l \odot \tilde{\mathbf{A}}_{\mathcal{T}}^l) \hat{\mathbf{Z}}^l + \beta_2^l (\mathbf{S}_{neg}^l \odot \tilde{\mathbf{A}}_{\mathcal{T}}^l) \hat{\mathbf{Z}}^l \right), \quad (32)$$

646 where $\hat{\mathbf{Z}}^l = \mathbf{Z}^{l-1} \mathbf{W}^l + b^l$, $\tilde{\mathbf{A}}_{\mathcal{T}}^l = \tilde{\mathbf{A}} \odot \mathcal{T}^l$ is an adjacency matrix weighted by the structure property,
647 β_0^l , β_1^l and β_2^l are learnable scalars. The neighbors are distinguished by the cosine similarity of node
648 representations with a threshold of 0:

$$\begin{aligned} \mathbf{S}_{ij}^l &= \begin{cases} \text{Cosine}(\mathbf{Z}_i, \mathbf{Z}_j), & i \neq j \ \& \ \mathbf{A}_{ij} = 1, \\ 0, & \text{otherwise.} \end{cases}, \\ \mathbf{S}_{pos, ij}^l &= \begin{cases} \mathbf{S}_{ij}^l, & \mathbf{S}_{ij}^l > 0, \\ 0, & \text{otherwise.} \end{cases}, \\ \mathbf{S}_{neg, ij}^l &= \begin{cases} \mathbf{S}_{ij}^l, & \mathbf{S}_{ij}^l < 0, \\ 0, & \text{otherwise.} \end{cases}. \end{aligned} \quad (33)$$

649 It can be revisited by HTMP with the following components:

$$\begin{aligned} \mathcal{A} &= [\mathbf{I}, \mathbf{A}_p, \mathbf{A}_n], \quad \mathcal{B} = [\mathbf{I}, \mathbf{S}_{pos}^l \odot \mathcal{T}^l, \mathbf{S}_{neg}^l \odot (\mathcal{T}^l)^l], \\ \tilde{\mathbf{Z}}_0^l &= (\mathbf{I} \odot \mathbf{I}) \mathbf{Z}^{l-1} \mathbf{W}^l = \mathbf{Z}^{l-1} \mathbf{W}^l, \\ \tilde{\mathbf{Z}}_1^l &= (\mathbf{A}_p \odot \mathbf{S}_{pos}^l \odot \mathcal{T}^l) \mathbf{Z}^{l-1} \mathbf{W}^l = (\mathbf{S}_{pos}^l \odot \mathcal{T}^l) \mathbf{Z}^{l-1} \mathbf{W}^l, \\ \tilde{\mathbf{Z}}_2^l &= (\mathbf{A}_n \odot \mathbf{S}_{neg}^l \odot \mathcal{T}^l) \mathbf{Z}^{l-1} \mathbf{W}^l = (\mathbf{S}_{neg}^l \odot \mathcal{T}^l) \mathbf{Z}^{l-1} \mathbf{W}^l, \end{aligned} \quad (34)$$

650 where \mathbf{A}_p and \mathbf{A}_n are discriminated by the representation similarities:

$$\begin{aligned} \mathbf{A}_{p, ij} &= \begin{cases} 1, & \mathbf{S}_{pos, ij}^l > 0 \ \& \ \mathbf{A}_{ij} = 1, \\ 0, & \text{otherwise.} \end{cases}, \\ \mathbf{A}_{n, ij} &= \begin{cases} 1, & \mathbf{S}_{neg, ij}^l < 0 \ \& \ \mathbf{A}_{ij} = 1, \\ 0, & \text{otherwise.} \end{cases}. \end{aligned} \quad (35)$$

651 The COMBINE function is an adaptive weighted addition. Specifically, GGCN divides the raw
652 neighborhood into positive and negative ones based on the similarities among node presentations.
653 On this basis, it aggregates messages from node ego, positive and negative neighborhoods, and
654 combines them with node-specific learnable weights. This approach allows for targeted processing
655 for homophilous and heterophilous neighbors, yet can suffer from the accuracy of discrimination,
656 which may lead to the accumulation of errors.

657 A.2.10 ACM-GCN

658 ACM-GCN [18] introduces 3 channels (identity, low pass and high pass) to capture different informa-
659 tion and mixes them with node-wise adaptive weights:

$$\mathbf{Z}^l = \text{diag}(\alpha_I^l) \mathbf{Z}^{l-1} \mathbf{W}_I^l + \text{diag}(\alpha_L^l) \hat{\mathbf{A}} \mathbf{Z}^{l-1} \mathbf{W}_L^l + \text{diag}(\alpha_H^l) (\mathbf{I} - \hat{\mathbf{A}}) \mathbf{Z}^{l-1} \mathbf{W}_H^l, \quad (36)$$

660 where $\text{diag}(\alpha_I^l)$, $\text{diag}(\alpha_L^l)$, $\text{diag}(\alpha_H^l) \in \mathbb{R}^{N \times 1}$ are learnable weight vectors.

661 It can be revisited by HTMP with the following components:

$$\begin{aligned} \mathcal{A} &= [\mathbf{I}, \mathbf{A}, \mathbf{A}], \quad \mathcal{B} = [\mathbf{I}, \mathbf{B}^d, \mathbf{I} - \mathbf{B}^d], \\ \tilde{\mathbf{Z}}_0^l &= (\mathbf{I} \odot \mathbf{I}) \mathbf{Z}^{l-1} \mathbf{W}_I^l = \mathbf{Z}^{l-1} \mathbf{W}_I^l, \\ \tilde{\mathbf{Z}}_1^l &= (\mathbf{A} \odot \mathbf{B}^d) \mathbf{Z}^{l-1} \mathbf{W}_L^l = \hat{\mathbf{A}} \mathbf{Z}^{l-1} \mathbf{W}_L^l, \\ \tilde{\mathbf{Z}}_2^l &= (\mathbf{A} \odot (\mathbf{I} - \mathbf{B}^d)) \mathbf{Z}^{l-1} \mathbf{W}_H^l = (\mathbf{I} - \hat{\mathbf{A}}) \mathbf{Z}^{l-1} \mathbf{W}_H^l, \end{aligned} \quad (37)$$

662 where the COMBINE function is adaptive weighted addition. Specifically, ACM-GCN aggregates
663 node ego, low-frequency, and high-frequency messages from ego and raw neighborhoods, and
664 combines them with node-wise adaptive weights. With simple but effective designs, ACM-GCN
665 achieves outstanding performance, which shows that complicated designs are not necessary.

666 **A.2.11 OrderedGNN**

667 OrderedGNN [21] is a SOTA method that introduces a node-wise adaptive dimension concatenation
 668 function to combine messages from neighbors of different hops:

$$\mathbf{Z}^l = \mathbf{P}_d^l \odot \mathbf{Z}^{l-1} + (1 - \mathbf{P}_d^l) \odot (\hat{\mathbf{A}}\mathbf{Z}^{l-1}), \quad (38)$$

669 where $\mathbf{P}_d \in \mathbb{R}^{N \times d_r}$ is designed to be matrix with each line $\mathbf{P}_{d,i}^l$ being a dimension indicate vector,
 670 which starts with continuous 1s while the others be 0s. In practice, to keep the differentiability, it's
 671 "soften" as follows:

$$\begin{aligned} \hat{\mathbf{P}}_d^l &= \text{cumsum}_{\leftarrow} \left(\text{softmax} \left(f_{\xi}^l \left(\mathbf{Z}^{l-1}, \hat{\mathbf{A}}\mathbf{Z}^{l-1} \right) \right) \right), \\ \mathbf{P}_d^l &= \text{SOFTOR}(\mathbf{P}_d^{l-1}, \hat{\mathbf{P}}_d^l), \end{aligned} \quad (39)$$

672 where f_{ξ}^l is a learnable layer that fuses two messages.

673 It can be revisited by HTMP with the following components:

$$\begin{aligned} \mathcal{A} &= [\mathbf{I}, \mathbf{A}], \quad \mathcal{B} = [\mathbf{I}, \mathcal{B}^d], \quad \mathbf{W}_0^l = \mathbf{W}_1^l = \mathbf{I}, \\ \tilde{\mathbf{Z}}_0^l &= (\mathbf{I} \odot \mathbf{I})\mathbf{Z}^{l-1} = \mathbf{Z}^{l-1}, \\ \tilde{\mathbf{Z}}_1^l &= (\mathbf{A} \odot \mathcal{B}^d)\mathbf{Z}^{l-1} = \hat{\mathbf{A}}\mathbf{Z}^{l-1}, \end{aligned} \quad (40)$$

674 where the COMBINE function is concatenation with node-wise adaptive dimensions. Specifically, in
 675 each layer, OrderedGNN aggregates messages from node ego and raw neighborhood and concatenates
 676 them with learnable dimensions. Combined with the multi-layer architecture, this approach can
 677 aggregate messages from neighbors of different hops and combine them not only with adaptive
 678 contributions but also as separately as possible.

679 **A.3 Analysis and Advice for Designing Models**

680 The HTMP mechanism splits the message-passing mechanism of HTGNNs into multiple modules,
 681 establishing connections among methods. For instance, most message passing in HTGNNs have
 682 personalized processing for nodes. Some methods [24, 11, 13, 22] utilize the learnable aggregation
 683 guidance and some others [14, 18, 21, 23] count on learnable COMBINE functions. Though
 684 neighborhood redefining is commonly used in HTGNNs, there are also many methods [24, 11, 18,
 685 20, 21] using only raw neighborhoods to handle heterophily and achieve good performance. Degree
 686 averaging, which plays the role of a low-pass filter to capture the smooth signals, can still work well
 687 in many HTGNNs [12, 14–16, 20]. High-order neighbor information may be helpful in heterophilous
 688 graphs. Existing HTGNNs utilize it in two ways: directly defining high-order [12, 13, 16, 17] or
 689 even full-connected [19] neighborhood indicators and by the multi-layer architecture of message
 690 passing [20, 21].

691 With the aid of HTMP, we can revisit existing methods from a unified and comprehensible perspective.
 692 An obvious observation is that *the coordination among designs is important while good combinations*
 693 *with easy designs can also achieve wonderful results*. For instance, in ACM-GCN [18], the separation
 694 and adaptive addition of ego, low-frequency, and high-frequency messages can accommodate the
 695 personalized conditions of each node. OrderedGNN's design [21], which includes an adaptive
 696 connection mechanism, ego/neighbor separation, and multi-layer architecture, allows discrete and
 697 adaptive combinations of messages from multi-hop neighborhoods. This advises us to *take into*
 698 *account all components simultaneously* when designing models. As an illustration, please be cautious
 699 about using multiple learnable components. Also, here are some additional model design tips and
 700 considerations. Please *separate the messages from node ego and neighbors*. When combining them
 701 afterward, whether by weighted addition or concatenation, this approach is at least harmless if not
 702 beneficial, especially when dealing with heterophilous graphs. Last but not least, try to design a
 703 model capable of *personalized handling different nodes*. Available components include but are not
 704 limited to, custom-defined neighborhood indicators, aggregation guidance with adaptive weights or
 705 estimated relationships, and learnable COMBINE functions. This is to accommodate the diversity
 706 and sparsity of neighborhoods that nodes in real-world graphs may have.

707 B Related Works

708 **Homophilous Graph Neural Networks.** Graph Neural Networks (GNNs) have showcased impres-
709 sive capabilities in handling graph-structured data. Traditional GNNs are predominantly founded
710 on the assumption of homophily, broadly categorized into two classes: spectral-based GNNs and
711 spatial-based GNNs. **Firstly**, spectral-based GNNs acquire node representations through graph
712 convolution operations employing diverse graph filters [1, 32, 33]. **Secondly**, spatial-based meth-
713 ods gather information from neighbors and update the representation of central nodes through the
714 message-passing mechanism [26, 28, 34]. Moreover, for a more **comprehensive understanding of**
715 **existing homophilous GNNs**, several unified frameworks [35, 36] have been proposed. Ma et al. [35]
716 propose that the aggregation process in some representative homophilous GNNs can be regarded
717 as solving a graph denoising problem with a smoothness assumption. Zhu et al. [36] establishes
718 a connection between various message-passing mechanisms and a unified optimization problem.
719 However, these methods have limitations, as the aggregated representations may lose discriminability
720 when heterophilous neighbors dominate [11, 12].

721 **Heterophilous Graph Neural Networks.** Recently, some heterophilous GNNs have emerged to
722 tackle the heterophily problem [11–23]. **Firstly**, a commonly adopted strategy involves *expanding the*
723 *neighborhood with higher homophily or richer messages*, such as high order neighborhoods [12, 13],
724 feature-similarity-based neighborhoods [13, 14], and custom-defined neighborhoods [15, 22]. **Sec-**
725 **ondly**, some approaches [11, 17–19, 23] aim to *leverage information from heterophilous neighbors*,
726 considering that not all heterophily is detrimental et al.[6]. **Thirdly**, some methods [12, 16, 20, 21]
727 adapt to heterophily by extending the combine function in message passing, creating variations for
728 addition and concatenation. On this basis, several works have **reviewed existing heterophilous**
729 **methods**. Zheng et al. [8] and Zhu et al. [9] identifies effective designs in heterophilous GNNs and
730 analyzes the relationship between heterophily and graph-related issues. Gong et al. [10] provide
731 a higher-level perspective on learning heterophilous graphs, summarizing and classifying existing
732 methods based on learning strategies, architectures, and applications. However, *these reviews merely*
733 *classify and list methods hierarchically, lacking unified understandings and not exploring the reason*
734 *behind the effectiveness of message passing in heterophilous graphs*.

735 C The Detail of Experiments on Synthetic Datasets

736 To explore the performance impact of homophily level, node degrees and compatibility matrix (CMs)
737 on simple GNNs, we conduct some experiments on synthetic datasets.

738 C.1 Synthetic Datasets

739 We construct synthetic graphs considering the factors of homophily, CMs and degrees. For homophily,
740 we set 3 levels including Lowh (0.2), Midh (0.5), and Highh (0.8). For CMs, we set two levels of
741 discriminability, including Easy and Hard. For degrees, we set two levels including Lowdeg (4)
742 and Highdeg (18). Note that with a certain homophily level, we can only control the non-diagonal
743 elements of CMs. Thus, there are a total of 12 synthetic graphs following the above settings. These
744 synthetic graphs are based on the Cora dataset, which provides node features and labels, which means,
745 only the edges are constructed. We visualize the CMs of these graphs in Figure 3. Since there is no
746 significant difference in CMs between low-degree and high-degree, we only plot the high-degree
747 ones. Further, the edges are randomly constructed under the guidance of these CMs and degrees to
748 form the synthetic graphs.

749 C.2 Experiments on Synthetic Datasets

750 We use GCN to analyze the performance impact of the above factors. The semi-supervised node
751 classification performance of GCN is shown in Table 5 while the baseline performance of MLP (72.54
752 ± 2.18) is the same among these datasets since their difference is only on edges. From these results,
753 we have some observations: (1) high homophily is not necessary, GCN can also work well on low
754 homophily but discriminative CM; (2) low degrees have a negative impact on performance, especially
755 when the CMs are relatively weak discriminative, this also indicates that nodes with lower degrees
756 are more likely to have confused neighborhoods; and (3) when dealing with nodes with confused
757 neighborhoods, GCN may contaminate central nodes with their neighborhoods’ messages, which

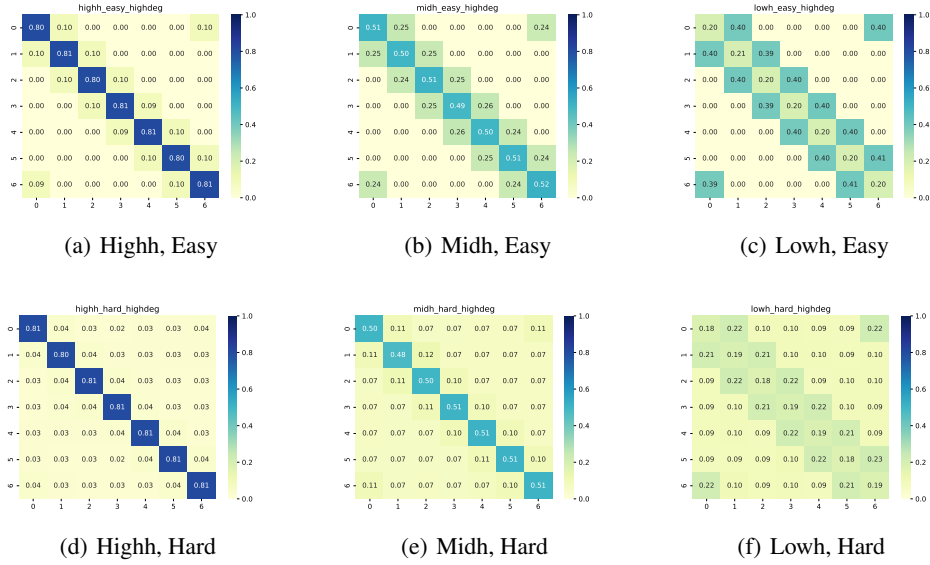


Figure 3: The visualization of compatibility matrix on synthetic graphs.

Table 5: Node classification accuracy of GCN on Synthetic Datasets.

Factors	Highh, Esay	Highh, Hard	Midh, Easy	Midh, Hard	Lowh, easy	Lowh, Hard
Highd	99.15 ± 0.35	99.48 ± 0.24	86.42 ± 4.13	90.52 ± 1.05	89.34 ± 2.19	39.22 ± 2.34
Lowd	89.98 ± 1.59	91.25 ± 0.85	70.85 ± 1.59	70.20 ± 1.41	56.46 ± 2.63	40.91 ± 1.75

758 leads to performance worse than MLP. This once again remind us the importance of ego/neighbor
 759 separation.

760 D Empirical Evidence for the Conjecture about CM

761 In this part, we show the empirical evidence for the conjecture about CM as mentioned in Sec 4.
 762 Specifically, we plot the observed and desired CM of ACM-GCN and GPRGNN in Figure 4. The
 763 results show that ACM-GCN and GPRGNN have enhanced the discriminability of CM, which can be
 764 empirical evidence for the conjecture.

765 The desired CMs are obtained as follows: For ACM-GCN, we leverage the learned weights in the
 766 COMBINE function to rebuild a weighted adjacency matrix \mathbf{A}^{acm} based on the low-pass filter $\hat{\mathbf{A}}$
 767 and high-pass filter $\mathbf{I} - \hat{\mathbf{A}}$, then regard \mathbf{A}^{acm} as the neighborhood and calculate the desired CM.

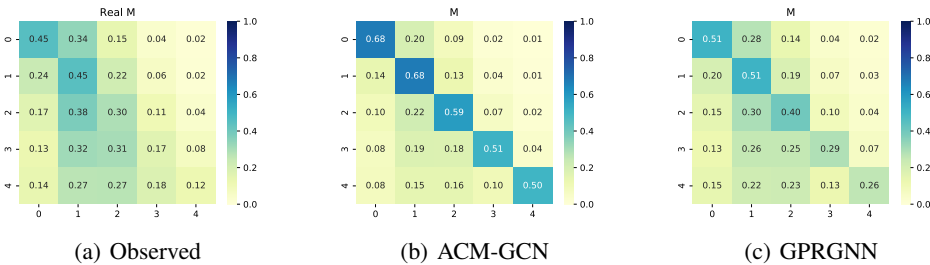


Figure 4: The visualization of compatibility matrix on Amazon-Ratings.

768 For GPRGNN, we utilize the leaned weights in the FUSE function to rebuild a weighted adjacency
 769 matrix \mathbf{A}^{gpr} based on the multi-hop adjacency matrixes $[\mathbf{I}, \mathbf{A}, \mathbf{A}^2, \dots, \mathbf{A}^k]$ then regard \mathbf{A}^{gpr} as the
 770 neighborhood and calculate the desired CM.

771 E Additional Detailed Implementation of CMGNN

772 **Overall Message Passing Mechanism.** The overall message passing mechanism in CMGNN is
 773 formatted as follows:

$$\begin{aligned} \mathbf{Z}^l &= \text{diag}(\alpha_0^l) \mathbf{Z}^{l-1} \mathbf{W}_0^l + \text{diag}(\alpha_1^l) \hat{\mathbf{A}} \mathbf{Z}^{l-1} \mathbf{W}_1^l + \text{diag}(\alpha_2^l) (\mathbf{A}^{sup} \odot \mathbf{B}^{sup}) \mathbf{Z}^{l-1} \mathbf{W}_2^l, \\ \mathbf{Z} &= \parallel_{l=0}^L \mathbf{Z}^l, \end{aligned} \quad (41)$$

774 where $\text{diag}(\alpha_0^l), \text{diag}(\alpha_1^l), \text{diag}(\alpha_2^l) \mathbb{R}^{N \times 1}$ are the learned combination weights introduced below.

775 **COMBNIE Function with Adaptive Weights.** Firstly, we list the aggregated messages $\tilde{\mathbf{Z}}_r^l$ from 3
 776 neighborhoods:

$$\begin{aligned} \tilde{\mathbf{Z}}_0^l &= \mathbf{Z}^{l-1} \mathbf{W}_0^l, \quad \tilde{\mathbf{Z}}_1^l = \hat{\mathbf{A}} \mathbf{Z}^{l-1} \mathbf{W}_1^l, \\ \tilde{\mathbf{Z}}_2^l &= (\mathbf{A}^{sup} \odot \mathbf{B}^{sup}) \mathbf{Z}^{l-1} \mathbf{W}_2^l. \end{aligned} \quad (42)$$

777 The combination weights are learned by an MLP with Softmax:

$$[\alpha_0^l, \alpha_1^l, \alpha_2^l] = \text{Softmax}(\text{Sigmoid}([\mathbf{Z}_0^l \parallel \mathbf{Z}_1^l \parallel \mathbf{Z}_2^l \parallel \mathbf{d}] \mathbf{W}_{att}^l) \mathbf{W}_{mix}^l), \quad (43)$$

778 where $\mathbf{W}_{att}^l \in \mathbb{R}^{(3d_r+1) \times 3}$ and $\mathbf{W}_{mix}^l \in \mathbb{R}^{3 \times 3}$ are two learnable weight matrixes, \mathbf{d} is the node
 779 degrees which may be helpful to weights learning.

780 **The Message Passing of Supplementary Prototypes.** In practice, the virtual prototype nodes are
 781 viewed as additional nodes, which have the same message passing mechanism as real nodes:

$$\begin{aligned} \mathbf{Z}^{ppt,l} &= \text{diag}(\alpha_0^{ppt,l}) \mathbf{Z}^{ppt,l-1} \mathbf{W}_0^l + \text{diag}(\alpha_1^{ppt,l}) \hat{\mathbf{A}}^{ppt} \mathbf{Z}^{ppt,l-1} \mathbf{W}_1^l \\ &\quad + \text{diag}(\alpha_2^{ppt,l}) (\mathbf{A}^{ppt,sup} \odot \mathbf{B}^{ppt,sup}) \mathbf{Z}^{ppt,l-1} \mathbf{W}_2^l, \\ \mathbf{Z}^{ppt} &= \parallel_{l=0}^L \mathbf{Z}^{ppt,l}, \end{aligned} \quad (44)$$

782 where $\mathbf{A}^{sup,ppt} = \mathbf{1} \in \mathbb{R}^{K \times K}$ and $\mathbf{B}^{sup,ppt} = \hat{\mathbf{C}}^{ppt} \hat{\mathbf{M}}$ are similar with those of real nodes.

783 **Update Strategy for the Estimation of the Compatibility Matrix.** For the sake of efficiency, we do
 784 not estimate the compatibility matrix in each epoch. Instead, we save it as fixed parameters and only
 785 update it when the evaluation performance is improved during the training.

786 **Predition of CMGNN.** CMGNN leverages the prediction of the model during message passing. For
 787 initialization, nodes have the same probabilities belonging to each class. During the message passing,
 788 the prediction soft label $\hat{\mathbf{C}}$ is replaced by the output of CMGNN, formatted as follow:

$$\hat{\mathbf{C}} = \text{CLA}(\mathbf{Z}), \quad (45)$$

789 where CLA is a classifier implemented by an MLP and \mathbf{Z} is the final node representations.

790 F More Detail about the Benchmark

791 In this section, we describe the details of the new benchmarks, including (i) the reason why we need
 792 a new benchmark: drawbacks of existing datasets; (ii) detailed descriptions of new datasets; (iii)
 793 baseline methods and the codebase; and (iv) details of obtaining benchmark performance.

794 F.1 Drawbacks in Existing Datasets

795 As mentioned in [30], the widely used datasets Cornell, Texas, and Wisconsin² have a too small scale
 796 for evaluation. Further, the original datasets Chameleon and Squirrel have an issue of data leakage,

²<https://www.cs.cmu.edu/afs/cs.cmu.edu/project/theo-11/www/wkwb>

797 where some nodes may occur simultaneously in both training and testing sets. Then, the splitting
 798 ratio of training, validation, and testing sets are different across various datasets, which is ignored in
 799 previous works.

800 Therefore, to build a comprehensive and fair benchmark for model effectiveness evaluation, we
 801 will newly organize 10 datasets with unified splitting across various homophily values in the next
 802 Subsection F.2.

803 F.2 New Datasets

804 In our benchmark, we adopt ten different types of publicly available datasets with a unified splitting
 805 setting (48%/32%/20% for training/validation/testing) for fair model comparison, including **Roman-**
 806 **Empire** [30], **Amazon-Ratings** [30], **Chameleon-F** [30], **Squirrel-F** [30], **Actor** [15], **Flickr** [37],
 807 **BlogCatalog** [37], **Wikics** [38], **Pubmed** [39], and **Photo** [29]. The datasets have a variety of
 808 homophily values from low to high. The statistics and splitting of these datasets are shown in Table 6.
 809 The detailed description of the datasets is as follows:

Table 6: Statistics and splitting of the experimental benchmark datasets.

Dataset	Nodes	Edges	Attributes	Classes	Avg. Degree	Undirected	Homophily	Train / Valid / Test
Roman-Empire	22,662	65,854	300	18	2.9	✓	0.05	10,877 / 7,251 / 4,534
Amazon-Ratings	24,492	186,100	300	5	7.6	✓	0.38	11,756 / 7,837 / 4,899
Chameleon-F	890	13,584	2,325	5	15.3	✗	0.25	427 / 284 / 179
Squirrel-F	2,223	65,718	2,089	5	29.6	✗	0.22	1,067 / 711 / 445
Actor	7,600	30,019	932	5	3.9	✗	0.22	3,648 / 2,432 / 1,520
Flickr	7,575	479,476	12,047	9	63.3	✓	0.24	3,636 / 2,424 / 1,515
BlogCatalog	5,196	343,486	8,189	6	66.1	✓	0.40	2,494 / 1,662 / 1,040
Wikics	11,701	431,206	300	10	36.9	✓	0.65	5,616 / 3,744 / 2,341
Pubmed	19,717	88,651	500	3	4.5	✓	0.80	9,463 / 6,310 / 3,944
Photo	7,650	238,162	745	8	31.1	✓	0.83	3,672 / 2,448 / 1,530

810 • **Roman-Empire**³ [30] is derived from the extensive article on the Roman Empire found on the
 811 English Wikipedia, chosen for its status as one of the most comprehensive entries on the platform.
 812 It contains 22,662 nodes and 65,854 edges between nodes. Each node represents an individual word
 813 from the text, with the total number of nodes mirroring the length of the article. An edge between
 814 two nodes is established under one of two conditions: the words are sequential in the text or they
 815 are linked in the sentence’s dependency tree, indicating a grammatical relationship where one word
 816 is syntactically dependent on the other. Consequently, the graph is structured as a chain graph,
 817 enriched with additional edges that represent these syntactic dependencies. The graph encompasses
 818 a total of 18 distinct node classes, with each node being equipped with 300-dimensional attributes
 819 obtained by fastText word embeddings [40].

820 • **Amazon-Ratings**³ [30] is sourced from the Amazon product co-purchasing network metadata
 821 dataset [41]. It contains 24,492 nodes and 186,100 edges between nodes. The nodes within
 822 this graph represent products, encompassing a variety of categories such as books, music CDs,
 823 DVDs, and VHS video tapes. An edge between nodes signifies that the respective products are
 824 often purchased together. The objection is to forecast the average rating assigned to a product by
 825 reviewers, with the ratings being categorized into five distinct classes. For the purpose of node
 826 feature representation, we have utilized the 300-dimensional mean values derived from fastText
 827 word embeddings [40], extracted from the textual descriptions of the products.

828 • **Chameleon-F** and **Squirrel-F**³ [30] are specialized collections of Wikipedia page-to-page net-
 829 works [42], of which the data leakage nodes are filtered out by [30]. Within these datasets, each
 830 node symbolizes a web page, and edges denote the mutual hyperlinks that connect them. The
 831 node features are derived from a selection of informative nouns extracted directly from Wikipedia
 832 articles. For the purpose of classification, nodes are categorized into five distinct groups based
 833 on the average monthly web traffic they receive. Specifically, Chameleon-F contains 890 nodes
 834 and 13,584 edges between nodes, with each node being equipped with 2,325-dimensional features.
 835 Squirrel-F contains 2,223 nodes and 65,718 edges between nodes, with each node being equipped
 836 with a 2,089-dimensional feature vector.

³<https://github.com/yandex-research/heterophilous-graphs/tree/main/data>

- 837 • **Actor**⁴ [15] is an actor-centric induced subgraph derived from the broader film-director-actor-writer
838 network, as originally presented by [43]. In this refined network, each node corresponds to an
839 individual actor, and the edges signify the co-occurrence of these actors on the same Wikipedia
840 page. The node features are identified through the presence of certain keywords found within
841 the actors’ Wikipedia entries. For the purpose of classification, the actors are organized into five
842 distinct categories based on the words of the actor’s Wikipedia. Statistically, it contains 7,600
843 nodes and 30,019 edges between nodes, with each node being equipped with a 932-dimensional
844 feature vector.
- 845 • **Flickr** and **Blogcatalog**⁵ [37] are two datasets of social networks, originating from the blog-sharing
846 platform BlogCatalog and the photo-sharing platform Flickr, respectively. Within these datasets,
847 nodes symbolize the individual users of the platforms, while links signify the followship relation-
848 ships that exist between them. In the context of social networks, users frequently create personalized
849 content, such as publishing blog posts or uploading and sharing photos with accompanying tag
850 descriptions. These textual contents are consequently treated as attributes associated with each
851 node. The classification objection is to predict the interest group of each user. Specifically, Flickr
852 contains 7,575 nodes and 479,476 edges between nodes. The graph encompasses a total of 9
853 distinct node classes, with each node being equipped with a 12047-dimensional attribute vector.
854 BlogCatalog contains 5,196 nodes and 343,486 edges between nodes. The graph encompasses a
855 total of 6 distinct node classes, with each node being equipped with 8189-dimensional attributes.
- 856 • **Wikics**⁶ [38] is a dataset curated from Wikipedia, specifically designed for benchmarking the
857 performance of GNNs. It is meticulously constructed around 10 distinct categories that represent
858 various branches of computer science, showcasing a high degree of connectivity. The node features
859 are extracted from the text of the associated Wikipedia articles, leveraging the power of pretrained
860 GloVe word embeddings [44]. These features are computed as the average of the word embeddings,
861 yielding a comprehensive 300-dimensional representation for each node. The dataset encompasses
862 a substantial network of 11,701 nodes interconnected by 431,206 edges.
- 863 • **Pubmed**⁷ [39] is a classical citation network consisting of 19,717 scientific publications with
864 44,338 links between them. The text contents of each publication are treated as their node attributes,
865 and thus each node is assigned a 500-dimensional attribute vector. The target is to predict which of
866 the paper categories each node belongs to, with a total of 3 candidate classes.
- 867 • **Photo**⁸ [29] is one of the Amazon subset network from [29]. Nodes in the graph represent goods
868 and edges represent that two goods are frequently bought together. Given product reviews as
869 bag-of-words node features, each node is assigned a 745-dimensional feature vector. The task is to
870 map goods to their respective product category. It contains 7,650 nodes and 238,162 edges between
871 nodes. The graph encompasses a total of 8 distinct product categories.

872 F.3 Baseline Methods and the Codebase

873 For comprehensive comparisons, we choose 13 representative homophilous and heterophilous GNNs
874 as baseline methods in the benchmark, including (i) Shallow base model: MLP; (ii) Homophilous
875 GNNs: GCN, GAT, GCNII; and (iii) Heterophilous GNNs: H2GCN, MixHop, GBK-GNN, GGCN,
876 GloGNN, HOGGCN, GPR-GNN. Detailed descriptions of some of these methods can be seen in
877 Appendix A.2.

878 To explore the performance of baseline methods on new datasets and facilitate future expansions,
879 we collect the official/reproduced codes from GitHub and integrate them into a unified codebase.
880 Specifically, all methods share the same data loaders and evaluation metrics. One can easily run
881 different methods with only parameters changing within the codebase. The codebase is based on the
882 PyTorch⁹ framework, supporting DGL¹⁰ and PyG¹¹. Detailed usages of the codebase are available in
883 the Readme file of the codebase.

⁴https://github.com/bingzhewei/geom-gcn/tree/master/new_data/film

⁵https://github.com/TrustAGI-Lab/CoLA/tree/main/raw_dataset

⁶<https://github.com/pmernyei/wiki-cs-dataset>

⁷<https://lincs.soe.ucsc.edu/datac>

⁸<https://github.com/shchur/gnn-benchmark>

⁹<https://pytorch.org>

¹⁰<https://www.dgl.ai>

¹¹<https://www.pyg.org>

884 **F.4 Details of Obtaining Benchmark Performance**

885 Following the settings in existing methods, we construct 10 random splits (48%/32%/20% for
886 train/valid/test) for each dataset and report the average performance among 10 runs on them along
887 with the standard deviation.

888 For all baseline methods except MLP, GCN, and GAT, we conduct parameter searches within the
889 search space recommended by the original papers. The searches are based on the NNI framework
890 with an anneal strategy. We use Adam as the optimizer for all methods. Each method has dozens
891 of search trails according to their time costs and the best performances are reported. The currently
892 known optimal parameters of each method are listed in the codebase. We run these experiments
893 on NVIDIA GeForce RTX 3090 GPU with 24G memory. The out-of-memory error during model
894 training is reported as OOM in Table 2.

895 **G More Details about Experiments**

896 In this section, we describe the additional details of the experiments, including experimental settings
897 and results.

898 **G.1 Additional Experimental Settings**

899 Our method has the same experimental settings within the benchmark, including datasets, splits,
900 evaluations, hardware, optimizer and so on as in Appendix F.4.

901 **Parameters Search Space.** We list the search space of parameters in Table 7, where patience is for
902 early stopping, nhidden is the embedding dimension of hidden layers as well as the representation
903 dimension d_r , relu_variant decides ReLU applying before message aggregation or not as in ACM-
904 GCN, structure_info determines whether to use structure information as supplement node features or
905 not.

Table 7: Parameters search space of our method.

Parameters	Range
learning rate	{0.001, 0.005, 0.01, 0.05}
weight_decay	{0, 1e-7, 5e-7, 1e-6, 5e-6, 5e-5, 5e-4}
patience	{200, 400}
dropout	[0, 0.9]
λ	{0, 0.01, 0.1, 1, 10}
layers	{1, 2, 4, 8}
nhidden	{32, 64, 128, 256}
relu_variant	{True, False}
structure_info	{True, False}

906 **Ablation Study.** In the ablation study, there are three variants of our methods: without SM, without
907 DL, without SM and DL. For "without SM", we delete the supplementary messages during message
908 passing, using only messages from node ego and raw neighborhood for combination. For "without
909 DL", we simply set $\lambda = 0$ to delete the discrimination loss. For "without SM and DL", we just
910 combine the above two settings.

911 **G.2 Additional Experimental Results**

912 In this subsection, we show some additional experimental results and analysis.

913 **G.2.1 Additional Results of CM Estimations**

914 The additional visualizations of CM estimations are shown in Figure 5. As we can see, our method
915 can estimate quite accurate CMs among various homophily and class numbers, which provides a
916 good foundation for the construction of supplementary messages.

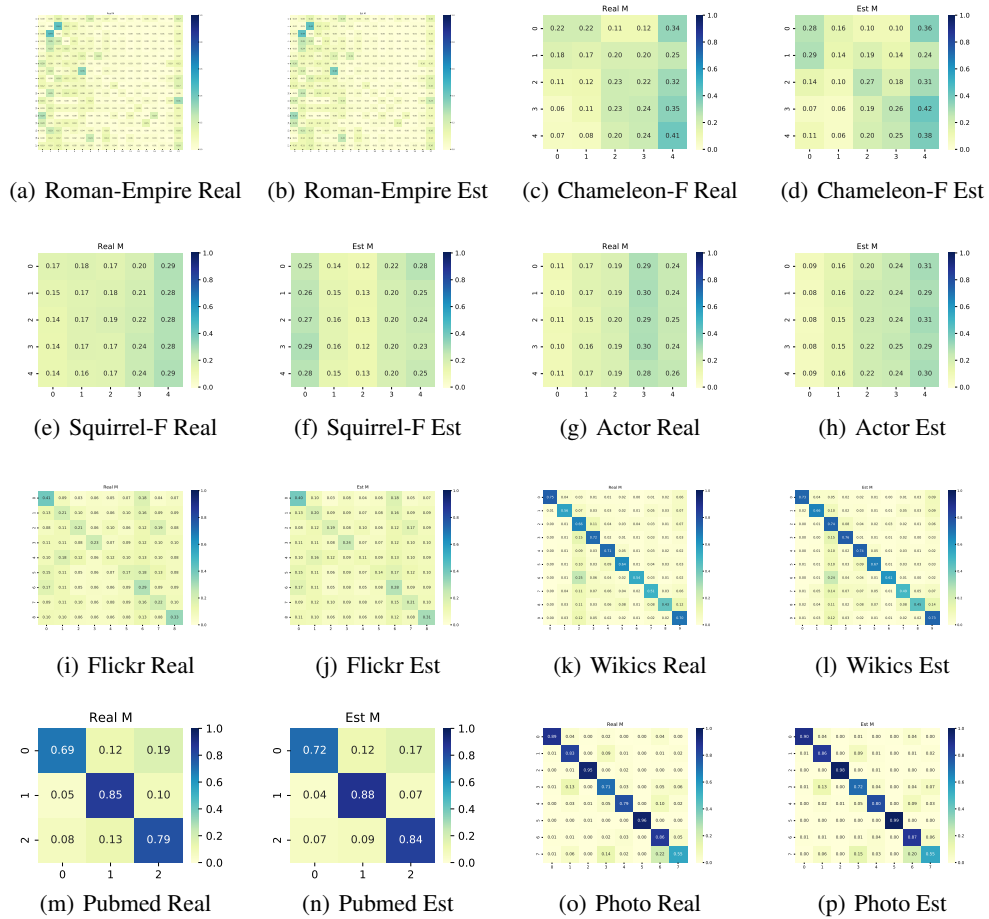


Figure 5: The visualization of real and estimated CMs on other datasets.

917 G.2.2 Additional Performance on Nodes with Various Levels of Degrees.

918 We show the additional performance on nodes with various degrees in Table 8. The results show that
 919 CMGNN can achieve relatively good performance on low-degree nodes, especially on heterophilous
 920 graphs. For the opposite results on homophilous graphs, we guess it may be due to the low-degree
 921 nodes in homophilous graphs having a more discriminative semantic neighborhood, such as a one-hot
 922 form. On the contrary, there are relatively more high-degree nodes with confused neighborhoods due
 923 to the randomness, which leads to the shown results on homophilous graphs.

Table 8: Node classification accuracy comparison (%) among nodes with different degrees.

Dataset	Roman-Empire					Chameleon-F					Actor					
	Deg. Prop.(%)	0~20	20~40	40~60	60~80	80~100	0~20	20~40	40~60	60~80	80~100	0~20	20~40	40~60	60~80	80~100
Ours		88.60	87.00	85.59	86.25	74.33	<u>40.73</u>	45.28	56.02	46.64	39.93	35.56	37.14	<u>38.40</u>	<u>36.03</u>	36.84
ACM-GCN		79.00	77.87	73.52	72.09	53.77	39.51	41.21	<u>52.25</u>	<u>45.80</u>	47.09	34.48	36.58	36.27	34.63	37.46
OrderedGNN		88.60	87.00	<u>85.56</u>	84.68	69.69	43.21	<u>44.51</u>	49.16	38.27	32.23	<u>35.94</u>	38.06	37.87	35.77	<u>37.15</u>
GCNII		86.79	85.14	85.20	<u>84.75</u>	<u>71.09</u>	34.84	42.56	47.50	40.45	<u>41.84</u>	36.89	<u>37.20</u>	38.53	38.02	36.99

Dataset	Squirrel					Pubmed					Photo					
	Deg. Prop.(%)	0~20	20~40	40~60	60~80	80~100	0~20	20~40	40~60	60~80	80~100	0~20	20~40	40~60	60~80	80~100
Ours		45.37	47.10	45.25	34.86	37.10	89.32	89.33	89.31	92.62	89.39	88.88	<u>95.76</u>	96.96	98.27	97.55
ACM-GCN		41.12	44.30	<u>44.22</u>	32.97	42.10	89.60	89.54	89.58	92.02	<u>89.23</u>	<u>89.88</u>	95.20	96.95	98.00	<u>97.56</u>
OrderedGNN		<u>43.78</u>	45.53	43.09	27.90	28.48	89.67	89.37	89.45	<u>92.54</u>	89.02	90.13	95.77	97.14	<u>98.24</u>	97.58
GCNII		43.08	<u>45.55</u>	43.65	<u>33.07</u>	<u>38.05</u>	89.77	<u>89.50</u>	89.24	92.45	88.86	88.89	95.36	<u>97.12</u>	97.83	96.64

924 **G.2.3 Efficiency Study**

925 **Complexity Analysis.** The number of learnable parameters in layer l of CMGNN is $3d_r(d_r + 1) + 9$,
 926 compared to $d_r d_r$ in GCN and $3d_r(d_r + 1) + 9$ in ACM-GCN. The time complexity of layer l is
 927 composed of 3 parts (i) AGGREGATE function: $O(Nd_r^2)$, $O(Nd_r^2 + Md_r)$ and $O(Nd_r^2 + NKd_r)$
 928 for identity neighborhood, raw neighborhood and the supplementary neighborhood respectively, where
 929 $M = |\mathcal{E}|$ denotes the number of edges; (ii) COMBINE function: $O(3N(3d_r + 1) + 12N)$ for adaptive
 930 weights calculating and $O(3N)$ for combination; (iii) FUSE function: $O(1)$ for concatenations.
 931 To this end, the time complexity of CMGNN is $O(Nd_r(3d_r + K + 9) + Md_r + 18N + 1)$, or
 932 $O(Nd_r^2 + Md_r)$ for brevity.

933 **Experimental Running Time.** we report the actual average running time (ms per epoch) of baseline
 934 methods and CMGNN in Table 9 for comparison. The results demonstrate that CMGNN can balance
 935 both performance effectiveness and running efficiency.

Table 9: Efficiency study results of average model running time (ms/epoch). OOM denotes out-of-memory error during the model training.

Method	Roman-Empire	Amazon-Ratings	Chameleon-F	Squirrel-F	Actor	Flickr	BlogCatalog	Wikies	Pubmed	Photo
MLP	7.8	7.0	6.1	6.5	6.3	9.1	6.7	6.4	6.1	5.8
GCN	33.8	33.4	7.9	20.6	34.4	37.2	30.4	25.5	35.6	28.1
GAT	15.9	67.3	10.3	14.0	30.8	66.2	17.6	26.8	33.4	36.0
GCNII	29.4	28.4	37.3	19.6	37.7	84.2	97.6	20.7	258.0	46.9
H2GCN	20.0	31.2	17.2	32.4	55.6	415.7	165.5	332.8	39.0	87.6
MixHop	434.6	486.3	21.9	31.0	30.6	90.4	81.6	277.4	89.5	172.2
GBK-GNN	119.8	191.8	31.0	238.1	157.9	OOM	OOM	198.6	137.0	193.3
GGCN	OOM	OOM	55.7	42.1	199.8	111.2	108.7	226.6	2290.8	105.2
GloGNN	25.4	19.3	121.8	23.3	1292	562.9	30.9	1658.1	43.2	677.4
HOGGCN	OOM	OOM	25.2	54.3	1002.9	707.3	367.4	1406	OOM	655.3
GPR-GNN	15.9	12.5	22.3	23.2	16.7	15.9	14.7	49.8	13.2	13.1
ACM-GCN	56.7	56.7	26.1	29.7	22.5	60.7	31.7	42.4	37.1	40.1
OrderedGNN	86.0	110.8	49.5	60.1	67.8	107.0	88.3	116.9	88.1	78.2
CMGNN	51.5	93.5	62.5	64.7	19.0	52.5	69.8	44.0	102.9	20.4

936 **NeurIPS Paper Checklist**

937 **1. Claims**

938 Question: Do the main claims made in the abstract and introduction accurately reflect the
939 paper's contributions and scope?

940 Answer: [\[Yes\]](#)

941 Justification: The contributions and scope of this paper are included in the abstract and
942 introduction.

943 Guidelines:

- 944 • The answer NA means that the abstract and introduction do not include the claims
945 made in the paper.
- 946 • The abstract and/or introduction should clearly state the claims made, including the
947 contributions made in the paper and important assumptions and limitations. A No or
948 NA answer to this question will not be perceived well by the reviewers.
- 949 • The claims made should match theoretical and experimental results, and reflect how
950 much the results can be expected to generalize to other settings.
- 951 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
952 are not attained by the paper.

953 **2. Limitations**

954 Question: Does the paper discuss the limitations of the work performed by the authors?

955 Answer: [\[Yes\]](#)

956 Justification: The limitations of this work are listed in the Sec 7.

957 Guidelines:

- 958 • The answer NA means that the paper has no limitation while the answer No means that
959 the paper has limitations, but those are not discussed in the paper.
- 960 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 961 • The paper should point out any strong assumptions and how robust the results are to
962 violations of these assumptions (e.g., independence assumptions, noiseless settings,
963 model well-specification, asymptotic approximations only holding locally). The authors
964 should reflect on how these assumptions might be violated in practice and what the
965 implications would be.
- 966 • The authors should reflect on the scope of the claims made, e.g., if the approach was
967 only tested on a few datasets or with a few runs. In general, empirical results often
968 depend on implicit assumptions, which should be articulated.
- 969 • The authors should reflect on the factors that influence the performance of the approach.
970 For example, a facial recognition algorithm may perform poorly when image resolution
971 is low or images are taken in low lighting. Or a speech-to-text system might not be
972 used reliably to provide closed captions for online lectures because it fails to handle
973 technical jargon.
- 974 • The authors should discuss the computational efficiency of the proposed algorithms
975 and how they scale with dataset size.
- 976 • If applicable, the authors should discuss possible limitations of their approach to
977 address problems of privacy and fairness.
- 978 • While the authors might fear that complete honesty about limitations might be used by
979 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
980 limitations that aren't acknowledged in the paper. The authors should use their best
981 judgment and recognize that individual actions in favor of transparency play an impor-
982 tant role in developing norms that preserve the integrity of the community. Reviewers
983 will be specifically instructed to not penalize honesty concerning limitations.

984 **3. Theory Assumptions and Proofs**

985 Question: For each theoretical result, does the paper provide the full set of assumptions and
986 a complete (and correct) proof?

987 Answer: [\[Yes\]](#)

988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025
1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041

Justification: We provide the formalization analysis in Appendix.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: The details of the method and experimental settings are provided in Appendix.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079
1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092

Answer: [Yes]

Justification: The data and code are available in the supplementary material.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: The detailed experimental settings are provided in Appendix.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: We report the average accuracy and the standard deviation as the performance in experiments.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).

- 1093 • It should be clear whether the error bar is the standard deviation or the standard error
1094 of the mean.
- 1095 • It is OK to report 1-sigma error bars, but one should state it. The authors should
1096 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
1097 of Normality of errors is not verified.
- 1098 • For asymmetric distributions, the authors should be careful not to show in tables or
1099 figures symmetric error bars that would yield results that are out of range (e.g. negative
1100 error rates).
- 1101 • If error bars are reported in tables or plots, The authors should explain in the text how
1102 they were calculated and reference the corresponding figures or tables in the text.

1103 8. Experiments Compute Resources

1104 Question: For each experiment, does the paper provide sufficient information on the com-
1105 puter resources (type of compute workers, memory, time of execution) needed to reproduce
1106 the experiments?

1107 Answer: [Yes]

1108 Justification: We list the hardware and software resources along with the space and space
1109 complexity in Appendix.

1110 Guidelines:

- 1111 • The answer NA means that the paper does not include experiments.
- 1112 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
1113 or cloud provider, including relevant memory and storage.
- 1114 • The paper should provide the amount of compute required for each of the individual
1115 experimental runs as well as estimate the total compute.
- 1116 • The paper should disclose whether the full research project required more compute
1117 than the experiments reported in the paper (e.g., preliminary or failed experiments that
1118 didn't make it into the paper).

1119 9. Code Of Ethics

1120 Question: Does the research conducted in the paper conform, in every respect, with the
1121 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

1122 Answer: [Yes]

1123 Justification: The research conducted in this paper conforms, in every respect, with the
1124 NeurIPS Code of Ethics.

1125 Guidelines:

- 1126 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 1127 • If the authors answer No, they should explain the special circumstances that require a
1128 deviation from the Code of Ethics.
- 1129 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
1130 eration due to laws or regulations in their jurisdiction).

1131 10. Broader Impacts

1132 Question: Does the paper discuss both potential positive societal impacts and negative
1133 societal impacts of the work performed?

1134 Answer: [Yes]

1135 Justification: The potential positive societal impacts are provided in the Introduction while
1136 the potential negative societal impacts are meaningless since this work is a foundational
1137 research.

1138 Guidelines:

- 1139 • The answer NA means that there is no societal impact of the work performed.
- 1140 • If the authors answer NA or No, they should explain why their work has no societal
1141 impact or why the paper does not address societal impact.

- 1142 • Examples of negative societal impacts include potential malicious or unintended uses
1143 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
1144 (e.g., deployment of technologies that could make decisions that unfairly impact specific
1145 groups), privacy considerations, and security considerations.
- 1146 • The conference expects that many papers will be foundational research and not tied
1147 to particular applications, let alone deployments. However, if there is a direct path to
1148 any negative applications, the authors should point it out. For example, it is legitimate
1149 to point out that an improvement in the quality of generative models could be used to
1150 generate deepfakes for disinformation. On the other hand, it is not needed to point out
1151 that a generic algorithm for optimizing neural networks could enable people to train
1152 models that generate Deepfakes faster.
- 1153 • The authors should consider possible harms that could arise when the technology is
1154 being used as intended and functioning correctly, harms that could arise when the
1155 technology is being used as intended but gives incorrect results, and harms following
1156 from (intentional or unintentional) misuse of the technology.
- 1157 • If there are negative societal impacts, the authors could also discuss possible mitigation
1158 strategies (e.g., gated release of models, providing defenses in addition to attacks,
1159 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
1160 feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

1162 Question: Does the paper describe safeguards that have been put in place for responsible
1163 release of data or models that have a high risk for misuse (e.g., pretrained language models,
1164 image generators, or scraped datasets)?

1165 Answer: [NA]

1166 Justification: This paper poses no such risk. The datasets we used are all publicly available
1167 online.

1168 Guidelines:

- 1169 • The answer NA means that the paper poses no such risks.
- 1170 • Released models that have a high risk for misuse or dual-use should be released with
1171 necessary safeguards to allow for controlled use of the model, for example by requiring
1172 that users adhere to usage guidelines or restrictions to access the model or implementing
1173 safety filters.
- 1174 • Datasets that have been scraped from the Internet could pose safety risks. The authors
1175 should describe how they avoided releasing unsafe images.
- 1176 • We recognize that providing effective safeguards is challenging, and many papers do
1177 not require this, but we encourage authors to take this into account and make a best
1178 faith effort.

12. Licenses for existing assets

1180 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
1181 the paper, properly credited and are the license and terms of use explicitly mentioned and
1182 properly respected?

1183 Answer: [Yes]

1184 Justification: The datasets and code of baseline methods are publicly available online. We
1185 cite the original paper and mark the URL in both papers and codebase.

1186 Guidelines:

- 1187 • The answer NA means that the paper does not use existing assets.
- 1188 • The authors should cite the original paper that produced the code package or dataset.
- 1189 • The authors should state which version of the asset is used and, if possible, include a
1190 URL.
- 1191 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 1192 • For scraped data from a particular source (e.g., website), the copyright and terms of
1193 service of that source should be provided.

- 1194
- 1195
- 1196
- 1197
- 1198
- 1199
- 1200
- 1201
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
 - For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
 - If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New Assets

1202

1203 Question: Are new assets introduced in the paper well documented and is the documentation

1204 provided alongside the assets?

1205 Answer: [Yes]

1206 Justification: We provide a public codebase along with an illustrative README file.

1207 Guidelines:

- 1208
- 1209
- 1210
- 1211
- 1212
- 1213
- 1214
- 1215
- The answer NA means that the paper does not release new assets.
 - Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
 - The paper should discuss whether and how consent was obtained from people whose asset is used.
 - At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and Research with Human Subjects

1216

1217 Question: For crowdsourcing experiments and research with human subjects, does the paper

1218 include the full text of instructions given to participants and screenshots, if applicable, as

1219 well as details about compensation (if any)?

1220 Answer: [NA]

1221 Justification: This paper does not involve crowdsourcing nor research with human subjects.

1222 Guidelines:

- 1223
- 1224
- 1225
- 1226
- 1227
- 1228
- 1229
- 1230
- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
 - Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
 - According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects

1231

1232

1233 Question: Does the paper describe potential risks incurred by study participants, whether

1234 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)

1235 approvals (or an equivalent approval/review based on the requirements of your country or

1236 institution) were obtained?

1237 Answer: [NA]

1238 Justification: This paper does not involve crowdsourcing nor research with human subjects.

1239 Guidelines:

- 1240
- 1241
- 1242
- 1243
- 1244
- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
 - Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.

1245
1246
1247
1248
1249

- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.