

ECGN: A CLUSTER-AWARE APPROACH TO GRAPH NEURAL NETWORKS FOR IMBALANCED CLASSIFICATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Classifying nodes in a graph is a common problem. The ideal classifier must adapt to any imbalances in the class distribution. It must also use information in the clustering structure of real-world graphs. Existing Graph Neural Networks (GNNs) have not addressed both problems together. We propose the Enhanced Cluster-aware Graph Network (*ECGN*), a novel method that addresses these issues by integrating cluster-specific training with synthetic node generation. Unlike traditional GNNs that apply the same node update process for all nodes, *ECGN* learns different aggregations for different clusters. We also use the clusters to generate new minority-class nodes in a way that helps clarify the inter-class decision boundary. By combining cluster-aware embeddings with a global integration step, *ECGN* enhances the quality of the resulting node embeddings. Our method works with any underlying GNN and any cluster generation technique. Experimental results show that *ECGN* consistently outperforms its closest competitors by up to 11% on some widely-studied benchmark datasets. The GitHub implementation for implementation and replication is publicly available on CodeLink.

1 INTRODUCTION

Graph Neural Networks (GNNs) have shown remarkable success in various tasks involving graph-structured data, including node classification (Kipf & Welling, 2016), link prediction (Zhang & Chen, 2018), and recommender systems (Ying et al., 2018). Indeed, GNNs have achieved state-of-the-art performance in many of these tasks. However, existing methods often expect all node classes and labels to be equally frequent. But in many real-world scenarios, node classes are imbalanced. For instance, most users on social network platforms are legitimate, but a small percentage are bots. This imbalance hurts the accuracy of bot detection (Mohammadrezaei et al., 2018). A similar challenge emerges when classifying websites by topics (Wang et al., 2020). A few topics are extremely popular, while most are rare. The popular topics (the majority classes) tend to dominate the loss function. Hence, the GNN focuses on these classes during training. This undermines the accuracy of the GNN for nodes of the minority classes. Hence, there is a need for GNN models capable of handling class-imbalanced node classification.

The class imbalance problem has been extensively studied in traditional machine learning. The solutions typically fall into three categories. *Data-level methods* balance the class distribution by over-sampling or under-sampling (Chawla et al., 2002; Kubat & Matwin, 1997). *Algorithm-level approaches* adjust the training process by using different misclassification penalties or prior probabilities for different classes (Ling & Sheng, 2008; Cui et al., 2019a). *Hybrid methods* combine both strategies to mitigate class imbalance (Batista et al., 2004a). However, all these methods assume independent and identically distributed data. By its very nature, the graph structure introduces dependencies between the nodes. Hence, such methods can yield suboptimal results when applied directly to graph datasets.

Compounding the class imbalance problem is the issue of *uniform node updates* in GNNs. Recall that GNNs update each node’s embedding using information from the node’s neighbors. The information exchange is mediated by trainable weight matrices. The same weights are uniformly applied to all nodes. However, such uniform node updates can cause two problems.

054 While standard GNNs use uniform update rules that aggregate information from neighboring nodes,
 055 they may not fully capture the rich local patterns and community-specific behaviors present in the
 056 graph. By incorporating cluster-aware updates, we aim to enhance the model’s ability to learn
 057 these localized structures. Graphs have intricate substructures, such as clustered communities or
 058 hubs (Girvan & Newman, 2002). Nodes within the same cluster exhibit higher similarity and
 059 stronger dependencies than nodes from different clusters. Current GNNs can miss these nuanced
 060 local patterns and community-specific behaviors by treating all nodes identically. Second, since the
 061 weights are optimized over the entire graph, they can be biased towards the majority class. Hence,
 062 the node update process is sub-optimal for nodes from the minority class. These problems can lead
 063 to poor embedding quality and underperformance in classification.

064 While imbalanced classification and node clustering are two orthogonal problems, we argue that the
 065 interplay between them can be leveraged to address class imbalance more effectively. Specifically,
 066 clusters capture rich local structures and dependencies within the graph. By incorporating cluster-
 067 aware updates, we can more accurately learn the nuanced relationships between nodes within a
 068 cluster, mitigating the dominance of majority classes during training. Clustering provides a natural
 069 framework for focusing on minority-class nodes in their local context, enabling us to preserve their
 070 distinct patterns and improve their representation quality. Thus, combining these two aspects allows
 071 us to address both class imbalance and the limitations of uniform node updates simultaneously.

072 **Our Contributions:** Although some existing studies have addressed either label imbalance (Zhao
 073 et al., 2021a; Zhou & Gong, 2023a) or cluster-aware updates (Chiang et al., 2019), there is little work
 074 on tackling both issues simultaneously. We propose the *Enhanced Cluster-aware Graph Network*
 075 (*ECGN*) to bridge this gap.

076 *ECGN* operates through a three-phase process. In the **pre-training** phase, we train cluster-specific
 077 GNNs in parallel. These GNNs extract information from local structures in the graph while ensuring
 078 that all embeddings map to the same latent space. The **node generation** phase is a novel way to
 079 generate synthetic nodes for the minority class. In particular, the synthetic node representations
 080 incorporate cluster-specific features. Finally, the **global-integration** phase integrates the outputs of
 081 the previous stages into a cohesive set of node embeddings. These capture the global information
 082 across the entire graph. Our framework can be used with any existing GNN and is applicable whether
 083 or not the clusters are known a priori.

084 *ECGN* makes four significant contributions:

- 086 • **cluster-aware node updates**, that capture local cluster-specific information;
- 087 • **addressing label imbalance** via innovative synthetic cluster-aware node generation;
- 088 • **seamless local to global integration**, allowing the embeddings to learn from different
 089 scales; and
- 090 • **broad applicability**, by enabling any underlying GNN model to be used.

092 We verify the accuracy of *ECGN* on five benchmark datasets and show that we consistently outper-
 093 form our closest competitors, with a **lift of up to 11% in F1 score** on the widely studied Citeseer
 094 dataset. These results confirm the applicability of *ECGN* for a wide range of real-world applications.

096 2 RELATED WORKS

098 We discuss the related work on learning under class imbalance, and Graph Neural Networks.

100 2.1 CLASS IMBALANCE LEARNING

102 Class imbalance in representation learning is a well-established topic in the field of machine learn-
 103 ing, having been extensively studied over the years (He & Garcia, 2009). The primary objective is to
 104 develop an unbiased classifier for a labeled dataset where the distribution is skewed, with majority
 105 classes having a substantially larger number of samples than minority classes. Notable contributions
 106 to this area include re-weighting and re-sampling techniques. Re-weighting methods modify the loss
 107 function by assigning greater importance to minority classes (Lin et al., 2017; Cui et al., 2019b), or
 by enlarging the margins for these classes (Cao et al., 2019; Liu et al., 2019; Menon et al., 2021).

108 On the other hand, re-sampling methods aim to balance the dataset by pre-processing the training
109 samples, employing strategies like over-sampling the minority classes (Chawla et al., 2002), under-
110 sampling the majority classes (Kubat & Matwin, 1997), or a combination of both (Batista et al.,
111 2004b).

112 With advancements in neural networks, re-sampling strategies have evolved to not only include tra-
113 ditional sampling techniques but also to incorporate generative approaches. For instance, modern
114 approaches augment minority class samples through generative methods (Liu et al., 2020), where
115 techniques such as SMOTE (Chawla et al., 2002) generate new samples by interpolating between
116 existing minority samples and their nearest neighbors. Additionally, other methods synthesize mi-
117 nority class samples by transferring knowledge from majority classes (Kim et al., 2020; Wang et al.,
118 2021b). However, most of these existing methods are tailored to independent and identically dis-
119 tributed samples and are not directly applicable to graph-structured data, where the relational context
120 between samples must be considered.

122 2.2 GRAPH NEURAL NETWORKS

123
124 Graph Neural Networks, first introduced in 2005 (Gori et al., 2005), have gained tremendous mo-
125 mentum in recent years with the advancements in deep learning, proving to be highly effective in
126 processing non-Euclidean structured data. GNNs typically operate using a message-passing frame-
127 work, where nodes iteratively gather information from their neighbors to learn low-dimensional
128 embeddings that capture the graph’s structural and feature information (Gilmer et al., 2017). These
129 techniques are generally divided into two categories: spectral-based and spatial-based methods.
130 Spectral-based methods exploit graph signal processing and leverage the graph Laplacian matrix to
131 perform node filtering (Defferrard et al., 2016; Kipf & Welling, 2017; Bianchi et al., 2020), while
132 spatial-based methods aggregate information directly from the local neighborhood of each node
133 based on the graph topology (Veličković et al., 2018; Hamilton et al., 2017; You et al., 2019).

134 Addressing the challenge of class imbalance within GNNs has been an active area of research.
135 Approaches such as GraphSMOTE (Zhao et al., 2021b) extend the popular SMOTE technique to
136 the embedding space of GNNs by synthesizing new minority nodes while also generating additional
137 edges, improving performance in imbalanced settings. Another approach, GraphENS (Park et al.,
138 2022), generates synthetic minority node features by mixing existing nodes from other classes. In
139 contrast, ClusterGCN (Chiang et al., 2019) leverages METIS-based graph partitioning to create
140 subclusters and trains the GNN in an SGD-based framework to capture cluster-specific information.
141 However, while ClusterGCN captures local structural information by operating on clusters, it still
142 applies uniform node updates within each subcluster and across the entire graph, meaning the same
143 update rules and aggregation functions are uniformly applied to all nodes without adapting to their
144 unique local structures or roles within the graph. This uniformity prevents ClusterGCN from fully
145 addressing the issue of uniform node updates during training, leading to a loss of fine-grained local-
global patterns crucial for capturing nuanced relationships and dependencies in graph learning.

146 In addressing the class imbalance problem in graph data, our work builds upon and extends existing
147 research that has explored various strategies for improving node classification performance under
148 imbalanced conditions. Recent studies such Park et al. (2021) and Qian et al. (2022) propose novel
149 methods for mitigating imbalance by modifying graph structures or introducing contrastive learn-
150 ing techniques. Wang et al. (2021a) introduced Distance-wise Prototypical GNNs, which focus on
151 learning class prototypes in an imbalanced setting. Similarly, Song et al. (2022) with TAM, and
152 Zeng et al. (2022) with Imgcl, highlight the importance of incorporating topological awareness and
153 contrastive learning in handling imbalanced datasets. Moreover, Zhou & Gong (2023b) leverages
154 data augmentation to improve minority-class representation. Other works include (Liu et al., 2023)
155 which introduced an interesting topological augmentation framework for class imbalance, and (Li
156 et al., 2023) which proposed a framework that synthesizes harder minor samples and incorporates
157 a SemiMixup module to expand minority class decision boundaries without encroaching on neigh-
158 boring class subspaces. While most of these methods focus on enhancing node embeddings, our
159 approach distinguishes itself by integrating a cluster-aware framework, where we not only focus on
160 adjusting the node updates but also generate synthetic minority-class nodes using a Cluster-Aware
161 SMOTE technique. This results in a more robust decision boundary between classes, particularly
when considering the underlying clustering structure in the graph. Our method, therefore, offers
a complementary approach to the existing literature, enhancing both local cluster information and

global graph integration. We further demonstrate that ECGN outperforms these state-of-the-art techniques across several benchmark datasets, offering a more effective solution for imbalanced node classification.

3 PROPOSED ALGORITHM

We are given a graph $G = (V, E)$, where V represents the set of nodes and E denotes the set of edges. Each node $v_i \in V$ is associated with a feature vector $\mathbf{x}_i \in \mathbb{R}^d$, forming the node feature matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$, where $n = |V|$ is the number of nodes and d is the feature dimension. The graph structure is represented by the adjacency matrix $A \in \{0, 1\}^{n \times n}$, where $A_{ij} = 1$ if there is an edge between nodes v_i and v_j , and $A_{ij} = 0$ otherwise. Each node $v_i \in V$ belongs to a class $y_i \in \{Y_1, Y_2, \dots, Y_c\}$, where c is the number of classes. The class distribution can be imbalanced. The classes for a subset of the nodes are known, and our goal is to predict the classes for the remaining nodes.

Next, we discuss ECGN’s architecture and algorithm, and provide details of our novel node generation step.

3.1 ARCHITECTURE OF ECGN

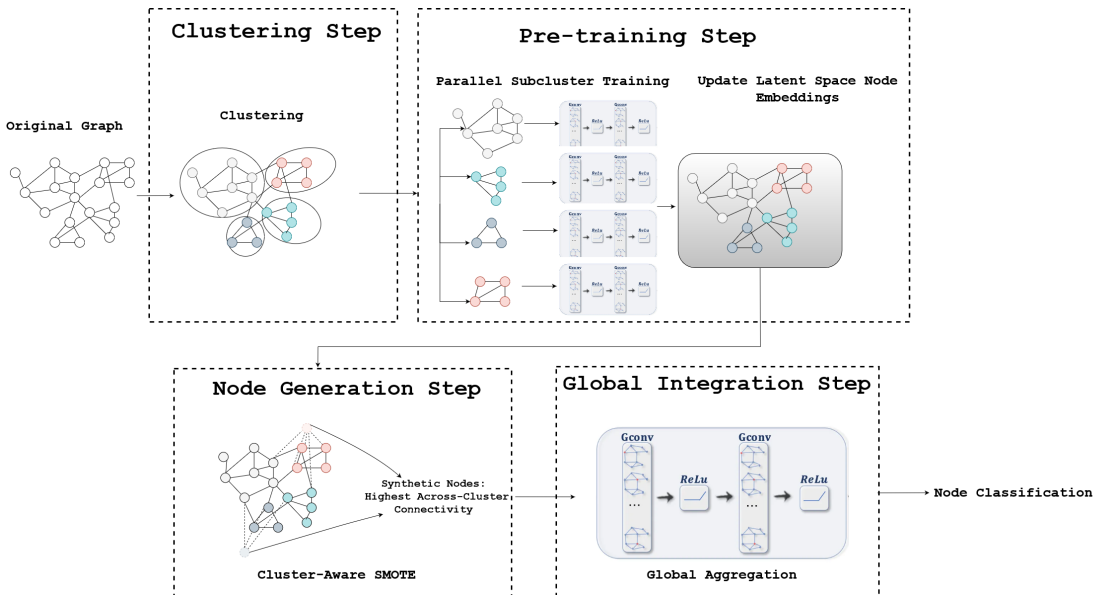


Figure 1: Working framework of ECGN architecture. We perform clustering on the nodes initially, train the sub-clusters independently in parallel, update the node embeddings to the original graph followed by Cluster-Aware SMOTE step and finally global integration.

The architecture of ECGN is presented in Figure 1. The algorithm begins by clustering the graph, unless the clusters are already known. Our method is flexible and can work with various clustering algorithms, though the choice of clustering algorithm does affect the results to some extent. We emphasize that the effectiveness of our method is not tied to any specific clustering approach, ensuring broad applicability across different scenarios. For example, fast algorithms such as Locality Sensitive Hashing (LSH) (Indyk & Motwani, 1998) or METIS (Karypis & Kumar, 1998) can be used for clustering (see Appendix A.4 for details).

Next comes the **pre-training** step. We first create a subgraph for each cluster. This subgraph includes only the nodes and edges within the cluster. Then, for each subgraph, we run a GNN

Algorithm 1 ECGN for Node Classification

```

216 1: Initialize node features:  $\mathbf{X} \in \mathbb{R}^{n \times d}$  and adjacency matrix:  $\mathbf{A} \in \{0, 1\}^{n \times n}$ .
217
218 2: Define the GNN architecture (network configuration) to be used for both pretraining and global
219 integration.
220
221 3: /* Clustering Step */
222 4: Clusters  $\mathcal{C} = \{C_1, C_2, \dots, C_k\}$ , either from prior knowledge, or obtained via Locality Sensitive
223 Hashing (LSH) or METIS partitioning (Appendix A.4).
224
225 5: /* Pre-training Step */
226 6: Initialize GNN parameters  $\theta$ .
227 7: for each cluster  $C_i \in \mathcal{C}$  do
228 8:   Construct subgraph  $G_{C_i} = (V_{C_i}, E_{C_i})$  containing only nodes and edges from  $C_i$ .
229 9:   Train a GNN on  $G_{C_i}$  with the objective of node classification within the cluster.
230 10:  Obtain embeddings  $\mathbf{H}_{C_i} \leftarrow \text{GNN}(G_{C_i} | \theta)$  for the nodes in  $C_i$ .
231 11: end for
232 12: Combine embeddings:  $\mathbf{H} \leftarrow \bigcup_{C_i \in \mathcal{C}} \mathbf{H}_{C_i}$ .
233
234 13: /* Synthetic Nodes Generation Step (Section 3.3) */
235 14: Generate new embeddings for minority class nodes using Cluster-Aware SMOTE .
236 15: Create nodes corresponding to these embeddings, and add edges to them in the graph.
237 16: Call the new graph  $G'$  with embeddings  $\mathbf{H}'$ 
238
239 17: /* Global Integration Step */
240 18:  $\mathbf{H}' \leftarrow$  convolution of  $G'$  using  $\mathbf{H}'$  as node features.
241 19: Perform node classification using the final embeddings  $\mathbf{H}'$ .

```

to generate embeddings for nodes in that cluster. All GNNs are run separately but share the same initialization. The result of this stage is a node embedding $\{\mathbf{h}_i\}$ for every node i in the graph. Unlike embeddings from a global GNN, these embeddings focus on information from the local cluster of each node. Appendix A.2 provides a detailed description.

The next stage tackles the class imbalance problem. For this, we generate new nodes and edge with a new technique called **Cluster-Aware SMOTE**. This differs from standard SMOTE in several ways. First, our method operates on the latent space of the cluster-aware node embeddings \mathbf{h}_i instead of the node features \mathbf{x}_i . Second, unlike SMOTE, our approach focuses on minority-class nodes that lie on cluster boundaries. The intuition is that in many real-world datasets, nodes from a given minority-class collect within one or a few clusters. Hence, the cluster boundaries are a proxy for inter-class decision boundaries. By generating nodes on the cluster borders, we can improve the decision boundary for node classification. Finally, after generating nodes with new embeddings, we link them to existing nodes. This step is needed for the following global aggregation step of *ECGN*. Details of Cluster-Aware SMOTE are presented in Section 3.3.

The last step is **global integration**. Here, we propagation global information throughout the nodes through graph convolution. The result is a set of node embeddings that combine information from local clusters as well as global patterns. Finally, these refined embeddings are used for node classification. Algorithm 1 provides the pseudo-code for *ECGN*.

3.2 PRE-TRAINING STEP

In the pre-training step, we aim to capture representative local features within each cluster by training cluster-specific GNN models independently. This process enhances the embeddings for each node by focusing on local structures relevant to its respective cluster.

We begin with the same initialization across all cluster-specific GraphSAGE models, ensuring consistency in the latent space. Each GraphSAGE model is trained with the same node classification objective but constrained to its assigned cluster subgraph. Specifically, for each cluster C_i , the GraphSAGE model learns node embeddings $\mathbf{h}_i^{(l+1)}$ at layer $l + 1$ as:

$$\mathbf{h}_i^{(l+1)} = \sigma \left(W_{\text{neigh}}^{(l)} \cdot \text{MEAN} \left(\{\mathbf{h}_j^{(l)} : j \in \mathcal{N}(i)\} \right) + W_{\text{self}}^{(l)} \cdot \mathbf{h}_i^{(l)} \right), \quad (1)$$

where:

- $\mathcal{N}(i)$ denotes the neighbors of node i within cluster C_i ,
- $W_{\text{neigh}}^{(l)}$ and $W_{\text{self}}^{(l)}$ are trainable weight matrices for aggregating neighbor and self-node embeddings at layer l , respectively,
- $\text{MEAN}(\cdot)$ computes the element-wise mean of the embeddings from neighbors $\{\mathbf{h}_j^{(l)} : j \in \mathcal{N}(i)\}$,
- σ is a non-linear activation function (e.g., ReLU).

The objective for each cluster-specific GraphSAGE model is formulated as a node classification loss $\mathcal{L}_{\text{cluster}}^{(k)}$ for the k -th cluster:

$$\mathcal{L}_{\text{cluster}}^{(k)} = - \sum_{i \in C_k} \sum_{c=1}^C y_{ic} \log \hat{y}_{ic}, \quad (2)$$

where y_{ic} is the true label of node i for class c , and \hat{y}_{ic} is the predicted probability for class c .

The combination of these cluster-specific embeddings results in enhanced node representations that capture unique, locally-relevant patterns. These embeddings are then passed to subsequent stages for integration with global information.

3.3 CLUSTER-AWARE SMOTE

Previous applications of SMOTE in graph contexts had several limitations. Synthetic nodes have been generated using graph features, but this ignores the link structure of the graph (Zhao et al., 2021b;a). Also, synthetic nodes were derived from minority-class seed nodes that were chosen randomly. If the seeds have poor connectivity, so do the synthetic nodes. Hence, the GNN’s node updates may not efficiently convey information about the minority class.

Our proposed approach, named Cluster-Aware SMOTE, addresses these challenges by leveraging both intra- and inter-cluster connectivity information. Our method prioritizes minority-class nodes that lie on the borders of their clusters. The resulting synthetic nodes lead to a more accurate decision boundary between classes. Also, instead of using the node features, we use the cluster-aware node embeddings from the pre-training step. This ensures that both features and connectivity information is used in creating the synthetic nodes.

The synthetic node generation process involves the following steps:

1. **Identify Highly Connected Nodes:** For each minority node v in class Y_m , we compute its connectivity to nodes in other clusters:

$$\text{connectivity}(v) = \sum_{\substack{u \in V \\ C(u) \neq C(v)}} A_{vu} \quad (3)$$

We select the top k nodes with the highest connectivity scores as the seed nodes for generating synthetic samples.

2. **Nearest Neighbor Selection:** For each seed node v from the minority class, we find its nearest neighbor $\text{nn}(v)$ within the same class in the embedding space:

$$\text{nn}(v) = \arg \min_{u \in Y_m, v \neq u} \|\mathbf{h}_v - \mathbf{h}_u\| \quad (4)$$

where \mathbf{h}_i is the cluster-aware embedding for node i from the pre-training step, and $\|\cdot\|$ denotes the Euclidean distance. For very large graphs, techniques like FAISS (Douze et al., 2024) can be used.

3. **Synthetic Node Generation:** We generate a synthetic node v' by interpolating between the embeddings of v and its nearest neighbor $\text{nn}(v)$:

$$\mathbf{x}_{v'} = (1 - \delta)\mathbf{x}_v + \delta\mathbf{x}_{\text{nn}(v)} \quad (5)$$

where δ is a random variable drawn from a uniform distribution in the range $[0, 1]$.

- 324 4. **Edge Preservation:** We add the synthetic node v' to the graph, and add edges from v' to all
 325 the neighbors of v . Thus, v' inherits the edges of v . We also add a link from v' to v . These
 326 steps ensure that the addition of the synthetic node v' preserves the local graph topology.
 327 The updated adjacency matrix is as follows:

$$328 A_{v'i} = \begin{cases} A_{vi} & \text{if } i \neq v \\ 1 & \text{if } i = v \end{cases} \quad (6)$$

- 329
 330
 331 5. **Oversampling Control:** To control the number of generated nodes, we introduce a pa-
 332 rameter α . For each minority class Y_m , we generate $\alpha \cdot |Y_m|$ synthetic nodes. We restrict
 333 the value of α such that the number of synthetic nodes remains less than 50% of the ma-
 334 jority class size. This ensures that the synthetic nodes do not degrade performance for the
 335 majority class.
 336

337 By focusing on nodes with high inter-cluster connectivity and generating synthetic samples in the
 338 latent space, our approach improves the diversity of synthetic nodes and better captures the underly-
 339 ing graph structure. This not only helps in addressing class imbalance but also enhances the overall
 340 classification performance by providing more representative samples for the minority class.
 341

342 4 EXPERIMENTS

343 We verify the accuracy of *ECGN* on five well-studied benchmark datasets. We first describe the
 344 datasets and the competing baselines. Then, we compare all algorithms on the node classification
 345 task. Finally, we show via ablation studies the need for the various steps of *ECGN*.
 346

347 **Datasets:** We evaluate *ECGN* on several widely-used public datasets for the node classification
 348 task. All the details of datasets and baselines can be found in in Appendix A.1. Specifically, Table 2
 349 shows the statistics and experimental setup for each dataset.
 350

351 **Baselines:** We compared *ECGN* against several state of the art methods. These include GraphSAGE
 352 (with and without cluster information), cluster-aware GNNs such as ClusterGCN (Chiang et al.,
 353 2019), GNNs for imbalanced classification such as GraphSMOTE (Zhao et al., 2021b), **recent state**
 354 **of the art models like GraphENS (Park et al., 2021) and TAM (Song et al., 2022)** and various other
 355 reweighting and oversampling schemes. Appendices A.1 and A.5 provide more details about the
 356 baselines and their hyperparameters. All methods were tested on node classification tasks, and
 357 compared on the basis of their F1 scores. To ensure robust and reliable results, we averaged the
 358 F1-scores over four different random seeds.

359 **Direct Inference from Subclusters:** We evaluated the effect of bypassing global integration by
 360 directly inferring from subclusters without combining their representations into a global model.
 361 This experiment highlights the role of global integration in connecting local subcluster relationships
 362 with the global graph structure. Details of this experiment are provided in Appendix A.6.

363 **Weight Transfer Strategies:** We explored three strategies for transferring pre-trained GNN weights
 364 from subclusters to the global model: Average Weights, Largest Subcluster Weights, and Best Per-
 365 forming Subcluster Weights. These strategies were compared against *ECGN* without weight transfer
 366 to assess their influence on the global model’s performance. For more details, see Appendix A.7.

367 **Sensitivity to Clustering Algorithms:** To evaluate the impact of different clustering methods, in-
 368 cluding METIS, LSH, and Random Clustering, we analyzed how they influence the performance of
 369 our method. This experiment was designed to assess the robustness of our approach to variations in
 370 graph partitioning techniques. Details are provided in Appendix A.9.

371 **Correlation between Clusters and Node Labels:** We analyzed the alignment between clusters and
 372 ground truth node labels (classes) by examining the distribution of class labels within clusters. This
 373 analysis provides insights into whether clustering processes naturally reflect the underlying label
 374 structure. See Appendix A.10 for a detailed explanation.
 375

376 For *ECGN*, we clusters the graphs using METIS. We used 3 clusters for CORA and CITESEER,
 377 7 for Amazon Computers, 40 for Reddit, and 20 for ogbn-arxiv. Section A.8 discusses how the
 number of clusters affect classification accuracy.

Table 1: *Results across different graph benchmark datasets: Mean F1-Scores and Balanced Accuracies are reported along with standard deviations. ECGN with SMOTE outperforms other methods in both F1-Score and Balanced Accuracy. The lift is up to 10 – 12% over the closest competitors (for Citeseer). ECGN with SMOTE is statistically significantly better than all competing methods at the $p < 0.1$ level except for the underlined rows.*

(a) CORA			(b) CITESEER		
Method	F1-Score	Balanced Accuracy	Method	F1-Score	Balanced Accuracy
GraphSAGE	0.655 ± 0.04	0.605 ± 0.04	GraphSAGE Baseline	0.3625 ± 0.04	0.305 ± 0.04
GraphSAGE (+ Cluster Features)	0.680 ± 0.04	0.635 ± 0.04	GraphSAGE (+ Cluster Features)	0.3825 ± 0.04	0.325 ± 0.04
SMOTE	0.690 ± 0.03	0.665 ± 0.02	SMOTE	0.450 ± 0.03	0.430 ± 0.03
Re-Weighting	0.670 ± 0.02	0.630 ± 0.02	Re-Weighting	0.560 ± 0.02	0.530 ± 0.02
EN-Weighting	0.680 ± 0.04	0.635 ± 0.04	EN-Weighting	0.520 ± 0.04	0.500 ± 0.04
Over-Sampling	0.645 ± 0.03	0.590 ± 0.02	Over-Sampling	0.345 ± 0.03	0.300 ± 0.02
CB-Sampling	0.710 ± 0.02	0.680 ± 0.01	CB-Sampling	0.510 ± 0.02	0.490 ± 0.02
GraphSMOTE	0.710 ± 0.03	0.680 ± 0.01	GraphSMOTE	0.590 ± 0.03	0.570 ± 0.02
GraphENS	0.738 ± 0.02	0.712 ± 0.03	GraphENS	0.630 ± 0.02	0.680 ± 0.03
TAM	0.735 ± 0.03	0.720 ± 0.03	TAM	0.625 ± 0.03	0.680 ± 0.03
ClusterGCN	0.727 ± 0.01	0.690 ± 0.03	ClusterGCN	0.580 ± 0.03	0.560 ± 0.03
Cluster-Aware SMOTE only	0.700 ± 0.02	0.670 ± 0.02	Cluster-Aware SMOTE only	0.460 ± 0.03	0.430 ± 0.03
ECGN (Without SMOTE)	0.732 ± 0.03	0.710 ± 0.03	ECGN (Without SMOTE)	0.610 ± 0.03	0.580 ± 0.03
ECGN (With SMOTE)	0.740 ± 0.03	0.720 ± 0.03	ECGN (With SMOTE)	0.650 ± 0.03	0.620 ± 0.03

(c) Reddit			(d) ogbn-arxiv		
Method	F1-Score	Balanced Accuracy	Method	F1-Score	Balanced Accuracy
GraphSAGE Baseline	0.740 ± 0.04	0.700 ± 0.04	GraphSAGE Baseline	0.3675 ± 0.04	0.345 ± 0.04
GraphSAGE (+ Cluster Features)	0.740 ± 0.04	0.700 ± 0.04	GraphSAGE (+ Cluster Features)	0.3825 ± 0.04	0.355 ± 0.04
SMOTE	0.760 ± 0.04	0.730 ± 0.04	SMOTE	0.390 ± 0.04	0.370 ± 0.04
Re-Weighting	0.770 ± 0.05	0.740 ± 0.05	Re-Weighting	0.400 ± 0.05	0.380 ± 0.05
EN-Weighting	0.750 ± 0.04	0.720 ± 0.04	EN-Weighting	0.410 ± 0.05	0.390 ± 0.05
Over-Sampling	0.750 ± 0.04	0.720 ± 0.04	Over-Sampling	0.385 ± 0.04	0.365 ± 0.04
CB-Sampling	0.710 ± 0.05	0.690 ± 0.05	CB-Sampling	0.390 ± 0.04	0.370 ± 0.04
GraphSMOTE	0.770 ± 0.05	0.740 ± 0.05	GraphSMOTE	0.410 ± 0.04	0.390 ± 0.04
GraphENS	0.780 ± 0.05	0.750 ± 0.05	GraphENS	0.440 ± 0.05	0.420 ± 0.05
TAM	0.770 ± 0.05	0.740 ± 0.05	TAM	0.430 ± 0.05	0.410 ± 0.05
ClusterGCN	0.760 ± 0.04	0.730 ± 0.04	ClusterGCN	0.430 ± 0.04	0.410 ± 0.04
ECGN (Without SMOTE)	0.770 ± 0.05	0.740 ± 0.05	ECGN (Without SMOTE)	0.442 ± 0.04	0.430 ± 0.05
ECGN (With SMOTE)	0.790 ± 0.05	0.760 ± 0.05	ECGN (With SMOTE)	0.450 ± 0.05	0.430 ± 0.05

(e) Amazon Computers		
Method	F1-Score	Balanced Accuracy
GraphSAGE Baseline	0.750 ± 0.03	0.720 ± 0.03
GraphSAGE (+ Cluster Features)	0.760 ± 0.03	0.730 ± 0.03
SMOTE	0.760 ± 0.03	0.730 ± 0.03
Re-Weighting	0.730 ± 0.03	0.700 ± 0.03
EN-Weighting	0.740 ± 0.03	0.720 ± 0.03
Over-Sampling	0.710 ± 0.04	0.670 ± 0.04
CB-Sampling	0.700 ± 0.04	0.680 ± 0.04
GraphSMOTE	0.760 ± 0.03	0.730 ± 0.03
GraphENS	0.770 ± 0.04	0.740 ± 0.04
TAM	0.760 ± 0.04	0.730 ± 0.04
ClusterGCN	0.740 ± 0.03	0.720 ± 0.03
ECGN (Without SMOTE)	0.767 ± 0.04	0.740 ± 0.04
ECGN (With SMOTE)	0.770 ± 0.01	0.760 ± 0.01

4.1 RESULTS

Table 1 shows the classification accuracy for all the datasets. We observe the following.

ECGN consistently outperforms other models across all datasets. The closest competitors are **GraphENS, TAM**, ClusterGCN and GraphSMOTE. However, *ECGN*'s F1-scores are higher by an average of nearly 5%. *ECGN* outperforms its closest competitors by up to 11% (e.g., on the Citeseer dataset).

Cluster-aware node updates are necessary. Consider two seemingly simple alternatives to the cluster-aware node updates of *ECGN*. One is to just provide the clusters as features to GraphSAGE.

432 The second is to just use the Cluster-Aware SMOTE without the pre-training step of *ECGN*. How-
433 ever, *ECGN* outperforms the former by 21% on average, and the latter by 14% on average.
434

435 **Cluster-Aware SMOTE improves classification accuracy.** The F1 score of *ECGN* using Cluster-
436 Aware SMOTE is 3% higher than *ECGN* without this step. For Citeseer, the difference increases to
437 6%. Thus, Cluster-Aware SMOTE adds value.

438
439 **Global integration is essential for performance.** Skipping global integration results in drastic F1-
440 score reductions, with scores dropping to 0.26 for Citeseer and 0.32 for Cora. These results highlight
441 that global integration is critical to effectively capture relationships between local subclusters and
442 the global graph structure. Detailed results are in Appendix A.6.

443 **Direct training outperforms weight transfer strategies.** The weight transfer strategies (Aver-
444 age Weights, Largest Subcluster Weights, Best Performing Subcluster Weights) slightly improve
445 performance, with the highest F1-score reaching 0.67 on Cora. However, these approaches still un-
446 derperform compared to direct global model training in *ECGN*, which achieves consistently higher
447 accuracy. Refer to Appendix A.7 for the full comparison.

448 ***ECGN* is robust to different clustering strategies.** Our method demonstrates robust performance
449 across various clustering algorithms, including METIS, LSH, and Random Clustering. This flexi-
450 bility underscores the adaptability of *ECGN* to different graph partitioning techniques without sig-
451 nificant degradation in performance. Details are provided in Appendix A.9.

452 **Clusters align meaningfully but not perfectly with class labels.** The analysis shows that clusters
453 often align with specific class labels. For example, in Cora, Cluster 1 contains over 70% of nodes
454 from Class 2, while in Citeseer, Cluster 1 has a nearly equal mix of Classes 0 and 1. Similarly, in
455 the Amazon Computers dataset, Cluster 7 has over 80% of nodes from Class 4. However, not all
456 clusters exhibit such strong alignment. These findings highlight that clustering processes capture
457 meaningful patterns but do not always perfectly reflect class labels. See Appendix A.10 for a deeper
458 analysis.

459
460 These experiments further validate the structure and robustness of *ECGN*, demonstrating its ability
461 to generalize effectively across both local subcluster models and global models.
462

463 5 CONCLUSION AND LIMITATIONS

464
465 In this paper, we introduced the Enhanced Cluster-aware Graph Network (*ECGN*), a novel frame-
466 work designed to address the challenges of class imbalance and subcluster-specific training in graph
467 neural networks. By integrating cluster-specific updates, synthetic node generation, and a global
468 integration step, *ECGN* demonstrates significant improvements in classification performance on im-
469 balanced datasets. Our experimental results show that *ECGN* not only enhances the representation
470 of minority classes but also maintains the structural integrity of the original graph, leading to more
471 accurate and robust predictions. We also stated that modifying *ECGN* either by bypassing global
472 tuning or integrating weight transfer learning hurts the performance. A detailed analysis of why
473 *ECGN* works is provided in Appendix A.3 for more clarity.

474 However, there are limitations to our approach. The reliance on subcluster partitioning may intro-
475 duce sensitivity to the quality of the clustering algorithm, and potentially impact the overall per-
476 formance if the clusters are not well-formed. Additionally, the synthetic node generation process,
477 while beneficial for handling imbalance, may introduce noise if not carefully managed, especially in
478 graphs with highly complex structures. Future work will focus on refining these aspects for broader
479 graph tasks like graph/link predictions, and evaluating *ECGN* on more diverse graph datasets to
480 further validate its effectiveness.

481 REFERENCES

482
483 Gustavo E. A. P. A. Batista, Ronaldo C. Prati, and Maria Carolina Monard. A study of the behavior
484 of several methods for balancing machine learning training data. *ACM SIGKDD Explorations*
485 *Newsletter*, 6(1):20–29, June 2004a. ISSN 1931-0153. doi: 10.1145/1007730.1007735. URL
<http://dx.doi.org/10.1145/1007730.1007735>.

- 486 Gustavo EAPA Batista, Ronaldo C Prati, and Maria Carolina Monard. A study of the behavior
487 of several methods for balancing machine learning training data. *ACM SIGKDD Explorations*
488 *Newsletter*, 6(1):20–29, 2004b.
- 489 Filippo Maria Bianchi, Daniele Grattarola, and Lorenzo Livi. Spectral clustering with graph neural
490 networks for graph pooling. *International Conference on Machine Learning*, pp. 11000–11010,
491 2020.
- 492 J. W. Butler. Machine sampling from given probability distributions. In *Symposium on Monte Carlo*
493 *Methods*, pp. 249–264. Wiley New York, 1956.
- 494 Kaidi Cao, Colin Wei, Adrien Gaidon, Nuno Arechiga, and Tengyu Ma. Learning imbalanced
495 datasets with label-distribution-aware margin loss. In *Advances in Neural Information Processing*
496 *Systems*, volume 32, 2019.
- 497 Nitesh V Chawla, Kevin W Bowyer, Lawrence O Hall, and W Philip Kegelmeyer. Smote: synthetic
498 minority over-sampling technique. *Journal of artificial intelligence research*, 16:321–357, 2002.
- 499 Wei-Lin Chiang, Xuanqing Liu, Si Si, Yang Li, Samy Bengio, and Cho-Jui Hsieh. Cluster-gcn:
500 An efficient algorithm for training deep and large graph convolutional networks. In *Proceedings*
501 *of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*,
502 KDD '19. ACM, July 2019. doi: 10.1145/3292500.3330925. URL [http://dx.doi.org/](http://dx.doi.org/10.1145/3292500.3330925)
503 [10.1145/3292500.3330925](http://dx.doi.org/10.1145/3292500.3330925).
- 504 Yin Cui, Menglin Jia, Tsung-Yi Lin, Yang Song, and Serge Belongie. Class-balanced loss based on
505 effective number of samples, 2019a. URL <https://arxiv.org/abs/1901.05555>.
- 506 Yin Cui, Menglin Jia, Tsung-Yi Lin, Yang Song, and Serge Belongie. Class-balanced loss based on
507 effective number of samples. In *Proceedings of the IEEE/CVF Conference on Computer Vision*
508 *and Pattern Recognition*, pp. 9268–9277, 2019b.
- 509 Michaël Defferrard, Xavier Bresson, and Pierre Vandergheynst. Convolutional neural networks
510 on graphs with fast localized spectral filtering. In *Advances in Neural Information Processing*
511 *Systems*, pp. 3844–3852, 2016.
- 512 Matthijs Douze, Alexandr Guzhva, Chengqi Deng, Jeff Johnson, Gergely Szilvasy, Pierre-
513 Emmanuel Mazaré, Maria Lomeli, Lucas Hosseini, and Hervé Jégou. The faiss library, 2024.
- 514 Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. Neural
515 message passing for quantum chemistry. *International Conference on Machine Learning*, pp.
516 1263–1272, 2017.
- 517 M. Girvan and M. E. J. Newman. Community structure in social and biological networks. *Proceed-*
518 *ings of the National Academy of Sciences*, 99(12):7821–7826, June 2002. ISSN 1091-6490. doi:
519 10.1073/pnas.122653799. URL <http://dx.doi.org/10.1073/pnas.122653799>.
- 520 Marco Gori, Gabriele Monfardini, and Franco Scarselli. A new model for learning in graph domains.
521 In *Proceedings of the 2005 IEEE International Joint Conference on Neural Networks, 2005.*,
522 volume 2, pp. 729–734. IEEE, 2005.
- 523 William L Hamilton, Rex Ying, and Jure Leskovec. Inductive representation learning on large
524 graphs. In *Advances in Neural Information Processing Systems*, pp. 1024–1034, 2017.
- 525 Haibo He and Eduardo A Garcia. Learning from imbalanced data. *IEEE Transactions on knowledge*
526 *and data engineering*, 21(9):1263–1284, 2009.
- 527 Piotr Indyk and Rameez Motwani. Approximate nearest neighbors: Towards removing the curse of
528 dimensionality. In *Proceedings of the Thirtieth Annual ACM Symposium on Theory of Computing*
529 *(STOC)*, pp. 604–613. ACM, 1998. doi: 10.1145/276698.276876.
- 530 Nathalie Japkowicz and Shaju Stephen. The class imbalance problem: A systematic study1. *In-*
531 *telligent Data Analysis*, 6(5):429–449, November 2002. ISSN 1088-467X. doi: 10.3233/
532 ida-2002-6504. URL <http://dx.doi.org/10.3233/IDA-2002-6504>.

- 540 George Karypis and Vipin Kumar. A fast and high quality multilevel scheme for partition-
541 ing irregular graphs. *SIAM Journal on Scientific Computing*, 20(1):359–392, 1998. doi:
542 10.1137/S1064827595287997.
- 543
- 544 Jeongwoo Kim, Kangwook Kim, Suha Kim, and Sungroh Kim. M2m: Imbalanced classification
545 via major-to-minor translation. In *Proceedings of the IEEE/CVF Conference on Computer Vision*
546 *and Pattern Recognition*, pp. 13896–13905, 2020.
- 547 Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional net-
548 works, 2016. URL <https://arxiv.org/abs/1609.02907>.
- 549
- 550 Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional net-
551 works. In *International Conference on Learning Representations*, 2017.
- 552
- 553 Miroslav Kubat and Stan Matwin. Addressing the curse of imbalanced training sets: one-sided se-
554 lection. In *Proceedings of the Fourteenth International Conference on Machine Learning (ICML)*,
555 pp. 179–186, 1997.
- 556 Wen-Zhi Li, Chang-Dong Wang, Hui Xiong, and Jian-Huang Lai. Graphsha: Synthesizing harder
557 samples for class-imbalanced node classification, 2023. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2306.09612)
558 [2306.09612](https://arxiv.org/abs/2306.09612).
- 559
- 560 Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense
561 object detection. In *Proceedings of the IEEE international conference on computer vision*, pp.
562 2980–2988, 2017.
- 563 Charles X. Ling and Victor S. Sheng. Cost-sensitive learning and the class imbalance problem.
564 2008. URL <https://api.semanticscholar.org/CorpusID:7601043>.
- 565
- 566 Boyan Liu, Xiang Zhang, Minki Song, and Jufu Feng. Deep generative oversampling for imbalanced
567 data classification. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34,
568 pp. 4606–4613, 2020.
- 569
- 570 Weyang Liu, Yandong Wen, Zhiding Yu, Meng Li, Bhiksha Raj, and Le Song. Large-margin soft-
571 max loss for convolutional neural networks. In *Proceedings of the 36th International Conference*
572 *on Machine Learning*, pp. 5075–5084, 2019.
- 573 Zhining Liu, Ruizhong Qiu, Zhichen Zeng, Hyunsik Yoo, David Zhou, Zhe Xu, Yada Zhu, Kommy
574 Weldemariam, Jingrui He, and Hanghang Tong. Class-imbalanced graph learning without class
575 rebalancing, 2023. URL <https://arxiv.org/abs/2308.14181>.
- 576
- 577 Aditya Krishna Menon, Sadeep Jayasumana, Abhinav Singh Rawat, Himanshu Jain, Andreas Veit,
578 and Sanjiv Kumar. Long-tail learning via logit adjustment. In *International Conference on Learn-*
579 *ing Representations*, 2021.
- 580 Mohammadreza Mohammadrezaei, Mohammad Ebrahim Shiri, and Amir Masoud Rahmani. Ident-
581 ifying fake accounts on social networks based on graph analysis and classification algorithms.
582 *Security and Communication Networks*, 2018:1–8, August 2018. ISSN 1939-0122. doi:
583 10.1155/2018/5923156. URL <http://dx.doi.org/10.1155/2018/5923156>.
- 584
- 585 Joonhyung Park, Jaeyun Song, and Eunho Yang. Graphens: Neighbor-aware ego network synthesis
586 for class-imbalanced node classification. In *International Conference on Learning Representa-*
587 *tions (ICLR)*, 2021.
- 588
- 589 Minseong Park, Soeun Park, Jiyoung Park, and U Kang. Graphens: Exploiting edge neighborhood
590 similarity for graph-based semi-supervised learning. In *Proceedings of the 2022 SIAM Interna-*
591 *tional Conference on Data Mining (SDM)*, pp. 576–584. SIAM, 2022.
- 592
- 593 Yiyue Qian, Chunhui Zhang, Yiming Zhang, Qianlong Wen, Yanfang Ye, and Chuxu Zhang. Co-
modality graph contrastive learning for imbalanced node classification. *Advances in Neural In-*
formation Processing Systems (NeurIPS), 35:15862–15874, 2022.

- 594 Jaeyun Song, Joonhyung Park, and Eunho Yang. Tam: topology-aware margin loss for class-
595 imbalanced node classification. In *International Conference on Machine Learning (ICML)*, pp.
596 20369–20383, 2022.
- 597 Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua
598 Bengio. Graph attention networks. In *International Conference on Learning Representations*,
599 2018.
- 601 Yu Wang, Charu C. Aggarwal, and Tyler Derr. Distance-wise prototypical graph neural network in
602 node imbalance classification. 2021a.
- 603 Zheng Wang, Xiaojun Ye, Chaokun Wang, Jian Cui, and Philip S. Yu. Network embedding with
604 completely-imbalanced labels. 2020. doi: 10.48550/ARXIV.2007.03545. URL [https://](https://arxiv.org/abs/2007.03545)
605 arxiv.org/abs/2007.03545.
- 607 Zifeng Wang, Zhenyu Chai, Guang Zhou, and Yue Zhang. Tackling class imbalance with
608 distribution-based implicit generative model. In *Proceedings of the 27th ACM SIGKDD Con-*
609 *ference on Knowledge Discovery & Data Mining*, pp. 47–57, 2021b.
- 611 Rex Ying, Ruining He, Kaifeng Chen, Pong Eksombatchai, William L. Hamilton, and Jure
612 Leskovec. Graph convolutional neural networks for web-scale recommender systems. In *Pro-*
613 *ceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data*
614 *Mining, KDD '18*. ACM, July 2018. doi: 10.1145/3219819.3219890. URL [http://dx.doi.](http://dx.doi.org/10.1145/3219819.3219890)
615 [org/10.1145/3219819.3219890](http://dx.doi.org/10.1145/3219819.3219890).
- 616 Jiaxuan You, Rex Ying, and Jure Leskovec. Position-aware graph neural networks. In *International*
617 *Conference on Machine Learning*, pp. 7134–7143. PMLR, 2019.
- 618 Liang Zeng, Lanqing Li, Zi-Chao Gao, Peilin Zhao, and Jian Li. Imgcl: Revisiting graph contrastive
619 learning on imbalanced node classification. In *AAAI Conference on Artificial Intelligence*, 2022.
- 621 Muhan Zhang and Yixin Chen. Link prediction based on graph neural networks, 2018. URL
622 <https://arxiv.org/abs/1802.09691>.
- 623 Tianxiang Zhao, Xiang Zhang, and Suhang Wang. Graphsmote: Imbalanced node classification on
624 graphs with graph neural networks. In *Proceedings of the 14th ACM International Conference*
625 *on Web Search and Data Mining, WSDM '21*, pp. 833–841. ACM, March 2021a. doi: 10.1145/
626 3437963.3441720. URL <http://dx.doi.org/10.1145/3437963.3441720>.
- 627 Tong Zhao, Yu Zhang, and Marinka Zitnik. Graphsmote: Imbalanced node classification on graphs
628 with graph neural networks. *arXiv preprint arXiv:2103.02770*, 2021b.
- 629 Mengting Zhou and Zhiguo Gong. Graphsr: A data augmentation algorithm for imbalanced node
630 classification, 2023a. URL <https://arxiv.org/abs/2302.12814>.
- 631 Mengting Zhou and Zhiguo Gong. Graphsr: A data augmentation algorithm for imbalanced node
632 classification. In *AAAI Conference on Artificial Intelligence*, 2023b.

637 A APPENDIX

639 A.1 DETAILS OF BASELINES AND EXPERIMENTAL DATASETS

641 In this section, we provide the details of the baselines implemented and experimental settings.

643 Table 2: Experimental settings to simulate imbalanced scenario for each datasets

Dataset	# Classes	# Imbalanced Classes	Majority Class Samples	Minority Class Samples	Total Nodes	Total Edges	Validation Nodes	Testing Nodes
Cora	7	3	200	20	2,708	5,429	2,050	1,426
CiteSeer	6	3	200	20	3,327	4,732	2,324	1,939
Amazon Computers	10	5	800	50	13,352	245,058	9,299	7,053
Reddit	40	10	1500	100	232,965	114,615,892	163,776	118,953
ogbn-arxiv	40	10	1500	100	169,343	1,166,851	118,540	86,348

- 648 • **Cluster-GCN:** A most popular GCN algorithm that is suitable for SGD-based training
649 by exploiting the graph clustering structure based on METIS partitioning. (Chiang et al.,
650 2019).
- 651 • **GraphSMOTE:** An oversampling method specifically designed for graphs that generates
652 synthetic minority nodes by interpolating between existing nodes within the minority class
653 (Zhao et al., 2021b).
- 654 • **Re-Weighting:** A classic cost-sensitive approach that adjusts the loss function with weights
655 inversely proportional to the number of samples in each class (Japkowicz & Stephen, 2002).
- 656 • **EN-Weighting:** A variant of the re-weighting method, which assigns weights based on the
657 Effective Number of samples in each class (Cui et al., 2019b).
- 658 • **Over-Sampling:** A traditional re-sampling method where minority nodes are repeatedly
659 sampled until each minority class has the same number of samples as the majority classes.
- 660 • **CB-Sampling:** A re-sampling method inspired by (Butler, 1956), which first selects a class
661 and then randomly samples a node from that class.
- 662 • **RU-Selection:** A baseline model that supplements the minority class by randomly select-
663 ing unlabeled nodes with pseudo-labels corresponding to the minority class until the class
664 distribution is balanced.
- 665 • **SU-Selection:** An extension of RU-Selection that selects unlabeled nodes based on their
666 similarity to the minority class, rather than random selection.

667 Here, we provide the details and explain the settings for the imbalanced scenario.
668

- 670 • **Cora Dataset:** Contains 2708 scientific publications categorized into 7 classes with 5429
671 links. We simulated a highly imbalanced scenario by sampling only 30% of the total sam-
672 ples available for the last 3 classes. Full-batch GD training was done, and number of
673 METIS partition clusters were fixed to be 3. Synthetic nodes were added such that the
674 minority class samples increases to 100 from 20 for each of the imbalanced class.
- 675 • **Citeseer Dataset:** Contains 3327 scientific publications classified into 6 categories with
676 4732 links. Full-batch GD training was done, and number of METIS partition clusters were
677 fixed to be 3. Synthetic nodes were added such that the minority class samples increases to
678 100 from 20 for each of the imbalanced class.
- 679 • **Reddit Dataset:** Consists of posts made by users on the Reddit online discussion forum,
680 categorized into 50 classes with over 230K nodes and 11M edges. The training was done
681 with stochastic neighborhood sampling with batch size of 1024. The number of METIS
682 partition clusters were fixed to be 40. Synthetic nodes were added such that the minority
683 class samples increases to 400 from 50 for each of the imbalanced class.
- 684 • **Amazon Computers Dataset:** Contains 13,752 nodes categorized into 10 classes with
685 245,861 edges. Full-batch GD training was done, and the number of METIS partition
686 clusters was fixed to 7. Synthetic nodes were added such that the minority class samples
687 increases to 600 from 100 for each of the imbalanced class.
- 688 • **ogbn-arxiv Dataset:** Comprises 169,343 scientific publications from arXiv, categorized
689 into 40 classes with 1,166,243 edges. To simulate an imbalanced scenario, we sampled
690 only 100 nodes for the last 10 classes. The training was done with stochastic neighborhood
691 sampling with batch size of 1024. The number of METIS partition clusters were fixed to
692 be 20. Synthetic nodes were added such that the minority class samples increases to 600
693 from 100 for each of the imbalanced classes.

694 A.2 ORIGINAL GRAPHSAGE VS. SUBCLUSTERED GRAPHSAGE

696 **Original GraphSAGE** GraphSAGE (Hamilton et al., 2017) generates node embeddings by ag-
697 gregating features from a node’s local neighborhood. Given a graph $G = (V, E)$ with $N = |V|$
698 nodes and initial node features $\mathbf{X} \in \mathbb{R}^{N \times F}$, the embedding of node v at layer k is updated as:

700 1. Neighborhood Aggregation:

$$701 \mathbf{h}_{\mathcal{N}(v)}^{(k)} = \text{AGGREGATE}^{(k)} \left(\left\{ \mathbf{h}_u^{(k-1)} \mid u \in \mathcal{N}(v) \right\} \right), \quad (7)$$

where $\mathcal{N}(v)$ is the set of neighbors of node v , and $\mathbf{h}_u^{(k-1)}$ is the embedding from the previous layer.

2. **Node Embedding Update:**

$$\mathbf{h}_v^{(k)} = \sigma \left(\mathbf{W}^{(k)} \cdot \text{CONCAT} \left(\mathbf{h}_v^{(k-1)}, \mathbf{h}_{\mathcal{N}(v)}^{(k)} \right) \right), \tag{8}$$

where $\mathbf{W}^{(k)}$ is the weight matrix, σ is an activation function, and $\mathbf{h}_v^{(0)} = \mathbf{x}_v$.

This process is repeated for K layers to capture K -hop neighborhood information. The final embeddings $\mathbf{h}_v^{(K)}$ are used for tasks like node classification.

Subcluster-Based GraphSAGE In *ECGN*, we enhance GraphSAGE by incorporating cluster-specific information:

1. **Graph Partitioning:** Divide G into M disjoint subclusters $\{G_1, G_2, \dots, G_M\}$, with corresponding feature matrices \mathbf{X}_i .
2. **Localized Learning:** For each subcluster G_i , perform GraphSAGE focusing only on Cluster-Aware edges:

$$\mathbf{h}_v^{(k)} = \sigma \left(\mathbf{W}^{(k)} \cdot \text{CONCAT} \left(\mathbf{h}_v^{(k-1)}, \text{AGGREGATE}_i^{(k)} \left(\left\{ \mathbf{h}_u^{(k-1)} \mid u \in \mathcal{N}_i(v) \right\} \right) \right) \right), \tag{9}$$

where $\mathcal{N}_i(v)$ denotes Cluster-Aware neighbors.

3. **Embedding Compilation:** Combine embeddings from all subclusters:

$$\mathbf{H}^{(K)} = \begin{pmatrix} \mathbf{H}_1^{(K)} \\ \mathbf{H}_2^{(K)} \\ \vdots \\ \mathbf{H}_M^{(K)} \end{pmatrix}. \tag{10}$$

4. **Global Integration:** Perform an additional GraphSAGE layer over G to integrate global information:

$$\mathbf{h}_v^{(\text{final})} = \sigma \left(\mathbf{W}^{(K+1)} \cdot \text{CONCAT} \left(\mathbf{h}_v^{(K)}, \text{AGGREGATE}^{(K+1)} \left(\left\{ \mathbf{h}_u^{(K)} \mid u \in \mathcal{N}(v) \right\} \right) \right) \right). \tag{11}$$

Key Advantages

- **Enhanced Local Patterns:** Captures fine-grained structures within clusters.
- **Computational Efficiency:** Allows parallel processing of subclusters.
- **Global Coherence:** Global aggregation integrates inter-cluster relationships.
- **Improved Handling of Imbalance:** Clustering aids in addressing class imbalance by focusing on underrepresented nodes within clusters.

By combining localized learning with global integration, the subcluster-based approach in *ECGN* effectively captures both local and global graph structures, leading to improved performance in node classification tasks.

A.3 VISUALIZING THE CLUSTERED COMMUNITIES AND ANALYZING THE SUB-CLUSTERED APPROACH

In this section, we visualize the clustered communities within the Cora and Citeseer datasets and we try to provide a theoretical explanation of why our sub-clustered approach is effective. We selected these datasets due to their manageable size and well-documented structure, which makes them ideal candidates for visual analysis. We divided the datasets into three clusters using the METIS algorithm and present the visualizations below.

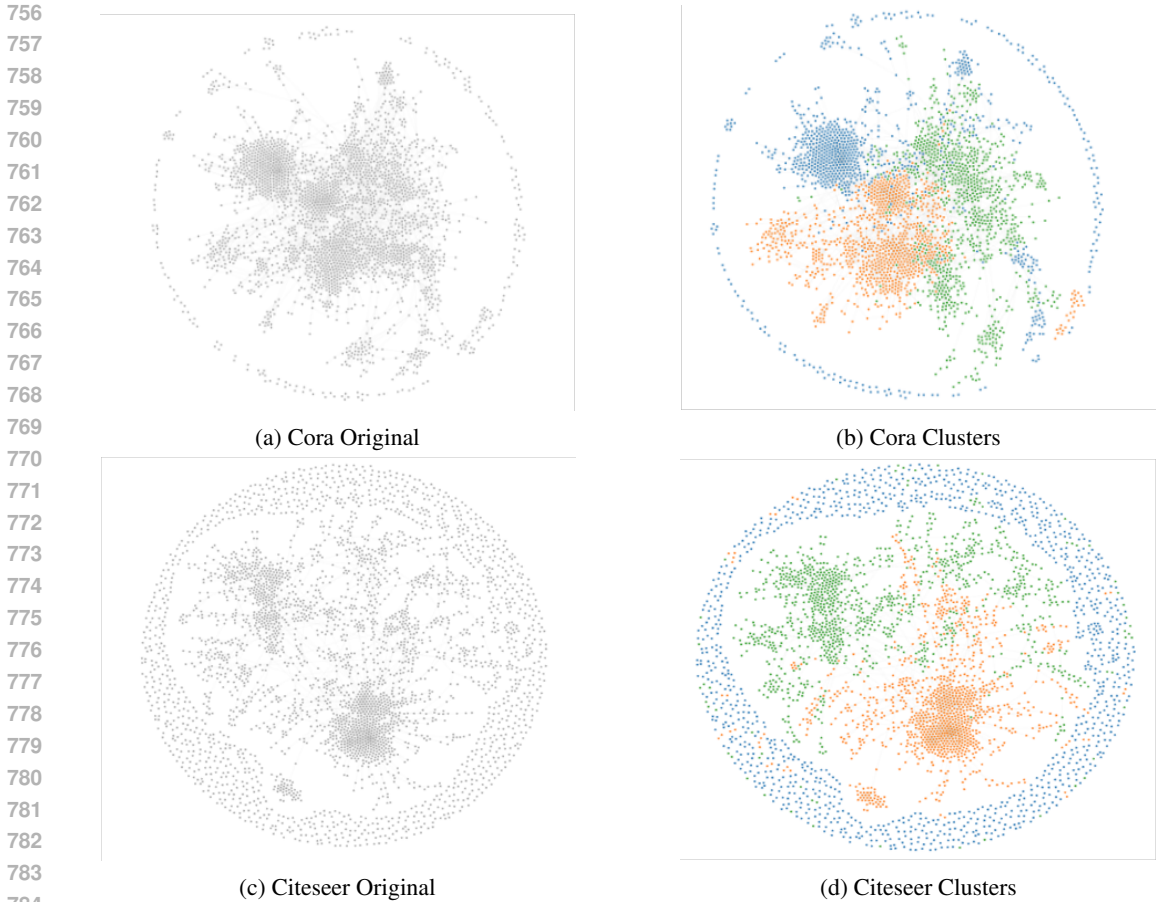


Figure 2: Visualizations of the original and clustered versions of the Cora and Citeseer datasets. The left column shows the original datasets, while the right column shows the datasets divided into three clusters using METIS clustering.

The figures in Figure 2 depict both the original and clustered versions of the Cora and Citeseer datasets. The original graphs (Figure 2a and Figure 2c) exhibit dense connectivity, which often leads to an entangled representation where the underlying community structure is not immediately apparent. By applying the METIS algorithm, we break down these dense graphs into distinct clusters (Figure 2b and Figure 2d), revealing the internal structure of the communities.

In this section, we try to provide a theoretical explanation of why our sub-clustered approach is effective.

Why does the Sub-Clustered Approach Work? The effectiveness of the sub-clustered approach can be theoretically explained through the following principles:

1. Capturing Localized Patterns When the graph is clustered into k subgraphs using METIS, we obtain subgraphs G_1, G_2, \dots, G_k such that:

$$G_i = (V_i, E_i), \quad \text{where } V_i \subseteq V \quad \text{and} \quad E_i \subseteq E$$

The feature matrix for each subgraph is $\mathbf{X}_i \in \mathbb{R}^{n_i \times d}$, where $n_i = |V_i|$ is the number of nodes in subgraph G_i .

By training on these subgraphs independently, the model learns localized patterns within each G_i , which are typically more homogeneous and easier to capture than the global patterns in G . The local

810 loss function for each subgraph can be expressed as:

$$811 \mathcal{L}_i = \frac{1}{n_i} \sum_{v_j \in V_i} \mathcal{L}(f(\mathbf{x}_j), y_j)$$

812 where $f(\mathbf{x}_j)$ is the model’s prediction for node v_j , and y_j is the true label. Training on localized
813 loss functions \mathcal{L}_i allows the model to optimize performance within each cluster before aggregating
814 the knowledge during global integration.

815 **2. Reducing Computational Complexity** The computational complexity of training a GNN on a
816 large graph G is often dominated by the cost of message passing and aggregation across the entire
817 graph. However, by decomposing G into smaller subgraphs G_1, G_2, \dots, G_k , and being able to train
818 them parallel independently, the computational cost is significantly reduced.

819 The overall complexity can be approximated as:

$$820 \text{Total Complexity} \approx \sum_{i=1}^k \mathcal{O}(|E_i|)$$

821 where $|E_i|$ is the number of edges in subgraph G_i . Since each $|E_i|$ is smaller than $|E|$ (the total
822 number of edges in the original graph), the sub-clustered approach leads to more efficient training.

823 **3. Addressing Imbalanced Data** In imbalanced graphs, certain classes of nodes may be under-
824 represented, making it difficult for the model to learn their characteristics. By isolating these nodes
825 within sub-clusters, the model can pay more focused attention to the minority classes.

826 Let C be the set of classes in the graph, with $|C_{\min}|$ and $|C_{\text{maj}}|$ representing the number of nodes in
827 the minority and majority classes, respectively. After clustering, the number of minority nodes in a
828 subgraph G_i can be denoted as $|C_{\min,i}|$. The training process can now focus on balancing the loss
829 contributions:

$$830 \mathcal{L}_i^{\text{balance}} = \frac{1}{|C_{\min,i}|} \sum_{v_j \in C_{\min,i}} \mathcal{L}(f(\mathbf{x}_j), y_j) + \frac{1}{|C_{\text{maj},i}|} \sum_{v_j \in C_{\text{maj},i}} \mathcal{L}(f(\mathbf{x}_j), y_j)$$

831 This ensures that the minority class nodes have a more significant influence on the model’s learning
832 process within each cluster.

833 **4. Global Structure Integration** After the initial training on sub-clusters, the model undergoes
834 global aggregation on the global graph G . This step integrates the knowledge learned from each
835 sub-cluster and ensures that node representations are coherent across the entire graph. The global
836 integration process can be represented as:

$$837 \mathcal{L}_{\text{global}} = \frac{1}{n} \sum_{v_j \in V} \mathcal{L}(f(\mathbf{x}_j), y_j)$$

838 This global loss function aligns the local representations and improves the overall performance of
839 the model.

840 By training on these clusters and then performing the global integration, the model leverages both
841 localized knowledge and global context, resulting in more accurate and generalizable node repre-
842 sentations. The sub-clustered approach mitigates the risk of overfitting to dominant structures and
843 promotes a more balanced and comprehensive understanding of the graph.

844 A.4 CLUSTERING AND COMMUNITY DETECTION TECHNIQUES

845 A.4.1 LOCALITY-SENSITIVE HASHING CLUSTERING (FEATURE BASED CLUSTERING)

846 We present an standard version of LSH clustering algorithm to efficiently handle large-scale datasets
847 with millions of nodes. The algorithm uses sparse random projections to group similar feature
848 vectors into clusters.

849 Algorithm Description

Given a feature matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$, where n is the number of samples and d is the number of features, the goal is to cluster the samples based on their similarity. The optimized algorithm proceeds as follows:

1. **Hash Table Creation:** We create T hash tables, each using P random projections. Each hash table is represented by a sparse random projection matrix $\mathbf{R} \in \mathbb{R}^{d \times P}$.

$$\mathbf{R}_t \sim \text{SparseRandomProjection}(d, P) \quad \text{for } t = 1, 2, \dots, T$$

2. **Hashing Feature Vectors:** Each feature vector $\mathbf{x}_i \in \mathbb{R}^d$ is projected into a lower-dimensional space using the hash tables. The projection is followed by taking the sign of the resulting vector to create hash keys.

$$\mathbf{h}_i^t = \text{sign}(\mathbf{x}_i \mathbf{R}_t) \quad \text{for } i = 1, 2, \dots, n \quad \text{and } t = 1, 2, \dots, T$$

The sign function is applied element-wise, resulting in a hash key $\mathbf{h}_i^t \in \{-1, 1\}^P$.

3. **Bucket Assignment:** Each hash key is used to group feature vectors into buckets. A bucket $\mathcal{B}_t(\mathbf{h})$ contains all vectors that share the same hash key \mathbf{h} for the t -th hash table.

$$\mathcal{B}_t(\mathbf{h}) = \{i \mid \mathbf{h}_i^t = \mathbf{h}\}$$

4. **Merging Buckets:** We merge buckets from all hash tables into preliminary clusters. Each node is assigned to a cluster based on its initial bucket assignments, ensuring unique assignments.

$$\text{clusters}[i] = \text{cluster_id} \quad \text{if } i \in \mathcal{B}_t(\mathbf{h}) \quad \forall t$$

5. **Cluster Refinement:** Each preliminary cluster is refined by computing the centroid of its feature vectors and using cosine similarity with a threshold to ensure nodes belong to the most similar cluster.

$$\text{similarity}(\mathbf{x}_i, \mathcal{C}) = \frac{\mathbf{x}_i \cdot \mathbf{c}_{\mathcal{C}}}{\|\mathbf{x}_i\| \|\mathbf{c}_{\mathcal{C}}\|} \quad \text{where } \mathbf{c}_{\mathcal{C}} = \frac{1}{|\mathcal{C}|} \sum_{\mathbf{x}_j \in \mathcal{C}} \mathbf{x}_j$$

Each node i is assigned to cluster \mathcal{C}^* if $\text{similarity}(\mathbf{x}_i, \mathcal{C}^*) > 0.5$.

6. **Final Cluster Formation:** The final clusters are formed by ensuring each node belongs to one and only one cluster.

$$\mathcal{C}_k \rightarrow \mathbf{C}_k \quad \text{where } \mathbf{C}_k \in \mathbb{R}^{|\mathcal{C}_k| \times d}$$

The algorithm ensures efficient clustering of high-dimensional data by leveraging the properties of locality-sensitive hashing and sparse random projections. The resulting clusters can be used in subsequent tasks such as classification, anomaly detection, and data summarization.

A.4.2 METIS PARTITIONING (STRUCTURE BASED CLUSTERING)

METIS partitioning is a graph partitioning technique designed to divide a graph into smaller, roughly equal-sized subgraphs while minimizing the edge cuts between them. The primary objective is to balance the load across subgraphs and reduce the communication volume in parallel computing environments.

Given a graph $G = (V, E)$ with vertices V and edges E , the goal is to partition G into k subgraphs G_1, G_2, \dots, G_k such that:

1. The size of each subgraph is approximately equal, i.e., $|V_i| \approx \frac{|V|}{k}$ for $i = 1, 2, \dots, k$.
2. The number of edges cut, denoted as $\text{cut}(G)$, is minimized. This is mathematically represented as minimizing the sum of weights of edges that have endpoints in different subgraphs:

$$\text{cut}(G) = \sum_{\substack{(u,v) \in E \\ u \in G_i, v \in G_j \\ i \neq j}} w(u, v)$$

where $w(u, v)$ is the weight of the edge between nodes u and v .

METIS employs a multilevel approach, which involves three main phases:

1. **Coarsening Phase:** The graph is iteratively coarsened by collapsing vertices and edges to form a series of progressively smaller graphs.
2. **Partitioning Phase:** A partitioning algorithm, often a variant of the *Kernighan-Lin* or *Fiduccia-Mattheyses* heuristic, is applied to the smallest graph to obtain an initial partition.
3. **Uncoarsening Phase:** The initial partition is projected back through the series of intermediate graphs, refining the partition at each level to improve the quality of the final partition.

This multilevel approach ensures that the partitioning process is both efficient and effective in producing high-quality partitions with balanced subgraph sizes and minimal edge cuts.

A.5 EXPERIMENTAL SETTINGS FOR BASELINE EXPERIMENTS

In our experiments, we use METIS partitioning to create subclusters for all the datasets. The experiments were configured with consistent training hyperparameters across datasets, including an initial learning rate of 0.01 for most datasets (with Reddit using 0.001), and the Adam optimizer. Each experiment ran for up to 1500 epochs, with early stopping after 40 steps if the validation performance did not improve. A batch size of 128 was used for Cora and Citeseer, while larger datasets such as Reddit, AmazonComputer, and ogbn_arxiv used batch sizes of 1024 or 2048 to accommodate their size. For model architecture, we adopted 2-layer GraphSAGE with a layer dimension of 128 for most datasets, though Reddit employed a smaller 64-dimensional GNN with 1 layer. We used the 'mean' aggregator for message passing and allowed for dynamic learning rates across layers. Full-batch training was used for Cora, Citeseer, and ogbn_arxiv, whereas AmazonComputer and Reddit were trained with neighborhood sampling (as documented in <https://docs.dgl.ai/en/0.8.x/guide/minibatch-node.html#guide-minibatch-node-classification-sampler>) with 4-layer deep neighborhood samples with sizes [4,4,4,4] due to their size. Additionally, datasets were clustered, with the number of clusters ranging from 3 (Cora, Citeseer), 7(AmazonComputer), 20(ogbn-arxiv) to 40 (Reddit).

A.6 DIRECT INFERENCE FROM SUBCLUSTERS: WHY DO WE NEED GLOBAL INTEGRATION?

In this section, we explore the impact of bypassing the global integration step and directly inferring from subclusters. This approach leverages local structures within each subcluster but neglects the global graph structure. To evaluate this, we conducted experiments on the Cora and Citeseer datasets, which offer manageable complexity for detailed analysis.

Table 3: Performance results when directly inferring from subclusters without global integration.

Dataset	Num Clusters	F1-Score Without Global Integration	Best ECGN F1-Score
Citeseer	3	0.26	0.65
Cora	3	0.32	0.74

Table 3 shows that skipping global integration leads to significantly lower F1-scores: 0.26 for Citeseer and 0.32 for Cora. This highlights the importance of integrating the global graph after subcluster training. Without it, the model learns *Cluster-Aware* relations but fails to generalize and learn *inter-cluster* relations, causing lower performance.

In conclusion, the global integration step is crucial as it bridges the gap between local subcluster structures and the overarching global graph. It ensures that the final node embeddings are both locally accurate and globally consistent, leading to better performance, as evidenced by the increased F1-scores after global integration.

A.7 REUSING GNN WEIGHTS FROM PRE-TRAINING IN GLOBAL INTEGRATION

We can think of ECGN as a transfer learning approach. In the pre-training step, we learn separately from each cluster. Then, we transfer the learnt embeddings to the global integration step.

This ensures a balance between local and global structures in the graph, which yields the strong performance of *ECGN*.

Extending this idea, we can ask: what if we transferred the GNN model weights from pre-training alongside the node embeddings? To explore this, we experimented with three different strategies for transferring weights:

1. **Average Weights:** Initialize the weights of the global GNN with the averaged weights of all subcluster GNNs.
2. **Largest Subcluster Weights:** GNN weights from the largest subcluster are transferred to the global model.
3. **Best Performing Subcluster Weights:** Weights from the best-performing subcluster are transferred to the global model.

Table 4: Performance results when transferring weights along with feature representations across different strategies.

Dataset	Weight Transfer Strategy	F1-Score with Weight Transfer	Best <i>ECGN</i> F1-Score
Citeseer	Average	0.51	0.65
	Largest	0.65	0.65
	Best	0.66	0.65
Cora	Average	0.52	0.74
	Largest	0.66	0.74
	Best	0.67	0.74
Amazon Computers	Average	0.54	0.78
	Largest	0.68	0.78
	Best	0.69	0.78

The results in Table 4 show that transferring pre-trained weights alongside the node embeddings yields mixed outcomes. The *Average Weights* strategy consistently performed the worst, likely because averaging diluted the unique structural information from each subcluster. The *Largest Subcluster Weights* and *Best Performing Subcluster Weights* improved performance over averaging but did not outperform *ECGN* without weight transfer. Overall, we find that the best performance comes from ignoring the pre-trained GNN weights in the global integration. This is what *ECGN* does.

We believe the reason for the above results lies in the delicate balancing act between learning from the graph’s local structure and its global context. The introduction of pre-trained weights \mathbf{W} into the global model appears to disrupt this balance. The pre-trained weights are perhaps too tailored to specific clusters. So, they may not generalize well when applied to the entire graph.

In conclusion, weight transfer offers no significant advantage and may reduce performance. The success of *ECGN* lies in combining local embeddings through global integration, capturing both local and global structures without transferring subcluster-specific weights.

A.8 SELECTING THE NUMBER OF CLUSTERS

The number of clusters affects both the effectiveness of the cluster-aware embeddings and *ECGN*’s computational efficiency. The key challenge lies in balancing the capture of fine-grained local patterns with the retention of important global structures.

In our experiments, we observed that the model’s performance is relatively robust to the number of clusters within a reasonable range. As shown in Figure 3, the F1-Score varies with the number of clusters across different datasets. For datasets like Cora and Citeseer, using **3 to 5 clusters** yielded comparable results, while extreme values—either too low or too high—led to decreased performance. This suggests that while the number of clusters is important, the model can tolerate variations without significant loss of effectiveness.

1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
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

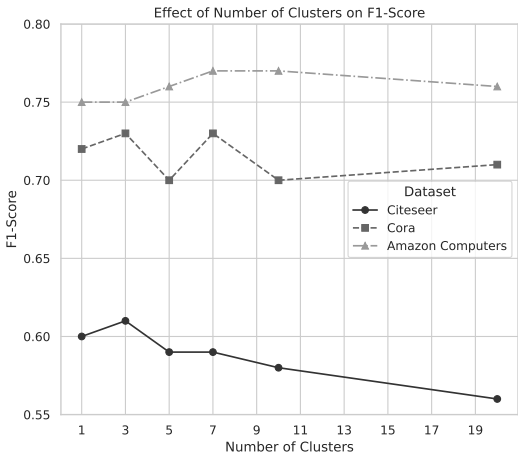


Figure 3: Effect of Number of Clusters on F1-Score for Various Datasets

A.9 SENSITIVITY TO CLUSTERING ALGORITHMS

We evaluated the impact of different clustering algorithms on our method’s performance by conducting experiments using METIS (a structure-based clustering algorithm), Locality-Sensitive Hashing (LSH, a feature-based clustering approach), and Random Clustering. These experiments were performed on the Cora and Citeseer datasets to assess how varying clustering strategies influence the overall effectiveness of our approach.

Table 5: Impact of Different Clustering Algorithms on Performance (F1-Score)

Clustering Algorithm	Cora	Citeseer
METIS	0.740 ± 0.03	0.650 ± 0.03
LSH	0.729 ± 0.02	0.639 ± 0.02
Random	0.690 ± 0.05	0.610 ± 0.03

As anticipated, METIS outperforms both LSH and Random Clustering on both datasets. METIS leverages the underlying graph structure to create clusters that reflect the inherent connectivity of the data, resulting in higher F1-scores. LSH, being a feature-based method, also provides meaningful clusters, though with slightly lower performance compared to METIS. Interestingly, Random Clustering achieves competitive results despite lacking semantic or structural coherence. This highlights the robustness of our method, which effectively utilizes localized training within clusters to mitigate class imbalance issues. By grouping and balancing nodes locally, our approach reduces the dominance of majority classes, thereby enhancing overall performance even when the clustering quality is suboptimal.

The performance of Random Clustering, although lower than METIS and LSH, remains surprisingly competitive. This can be attributed to our method’s ability to perform localized training within each cluster, ensuring that each smaller group of nodes is more balanced. This reduces the dominance of majority classes within each cluster, allowing the model to learn more effectively. Therefore, the method demonstrates significant robustness, maintaining effectiveness even when the clustering lacks inherent semantic or structural coherence.

In scenarios where clustering aligns perfectly with class labels, the performance of our method could potentially see further improvements due to reduced intra-cluster variance and better separation between classes. However, it is important to note that our method does not rely on such perfect alignment and achieves strong performance even with imperfect clustering, underscoring its flexibility and adaptability to various clustering qualities.

Overall, the ablation study and subsequent analysis demonstrate that our method is both robust and flexible, capable of maintaining high performance across different clustering algorithms. This

adaptability is a key strength, allowing our approach to be effectively applied in various contexts and with different types of data clustering strategies.

A.10 CORRELATION BETWEEN CLUSTERS AND NODE LABELS(CLASSES)

To better understand the relationship between clusters and node labels (classes), we conducted an analysis examining how nodes with specific labels are distributed across clusters. Clustering algorithms, such as Locality-Sensitive Hashing (LSH) and METIS, group nodes based on their features or structural properties. However, these clusters are formed independently of the node labels, making it important to investigate the alignment (or lack thereof) between clusters and classes.

This analysis was performed on three datasets: Cora, Citeseer, and Amazon Computers. Figures 4, 5, and 6 illustrate the distribution of node classes within clusters for each dataset. Each bar plot in these figures shows the percentage of nodes from different classes within a specific cluster.

Cora Dataset: In the Cora dataset (Figure 4), certain clusters exhibit a strong alignment with specific classes. For instance, Cluster 1 contains over 70% of nodes from Class 2, suggesting that the clustering algorithm successfully grouped nodes with similar features from this class. However, other clusters, such as Cluster 3, display a more diverse mix of classes, indicating that nodes with overlapping features across classes can be grouped together.

Citeseer Dataset: The Citeseer dataset (Figure 5) reveals a more fragmented relationship between clusters and classes. For example, Cluster 2 is dominated by Class 2, while Cluster 1 shows an almost equal mix of Classes 0 and 1. This variation reflects both the dataset’s inherent label imbalance and the clustering algorithm’s sensitivity to feature similarities. Such patterns emphasize that clusters may not always align perfectly with specific classes.

Amazon Computers Dataset: The Amazon Computers dataset, with its larger number of classes, exhibits both highly concentrated and diverse clusters (Figure 6). For example, Cluster 7 is dominated by Class 4, with over 80% of its nodes belonging to this class. In contrast, Cluster 2 contains nodes from a wider range of classes, showcasing the clustering algorithm’s adaptability in handling datasets with complex structures and larger numbers of classes.

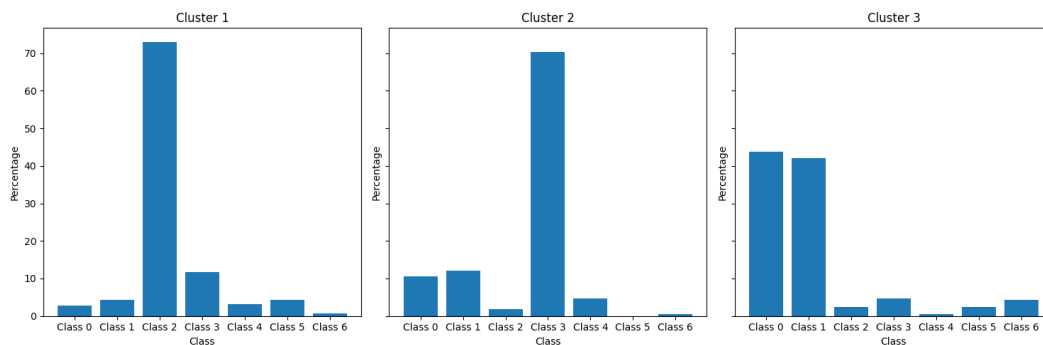


Figure 4: Class distributions across clusters for the Cora dataset.

Key Observations:

- Clusters often show a concentration of nodes from one or two dominant classes, indicating that clustering algorithms effectively group nodes with similar features or structural relationships.
- In some cases, nodes from different classes are grouped together, especially when features overlap or structural relationships between nodes span across classes.
- Datasets with more classes and higher complexity, such as Amazon Computers, display a mix of highly concentrated and diverse clusters, reflecting the adaptability of the clustering approach to different data characteristics.

Implications for Methodology: This experiment demonstrates that while clusters and classes may align in some cases, this alignment is not deterministic. The observed correlations indicate that

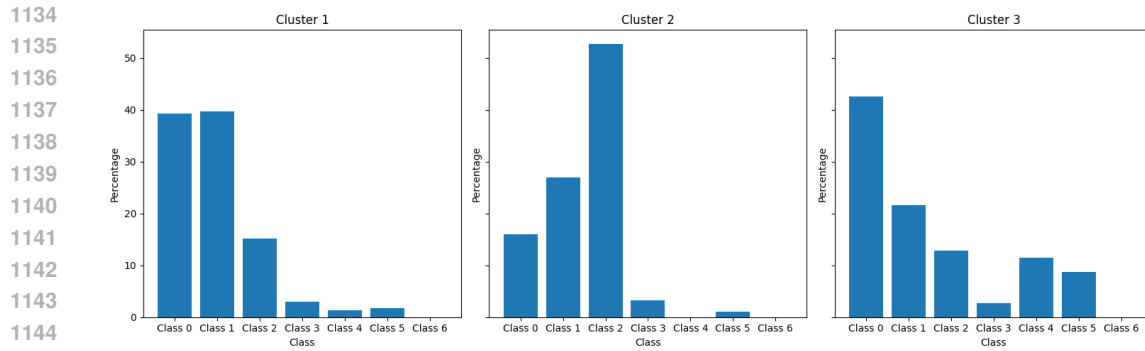


Figure 5: Class distributions across clusters for the Citeseer dataset.

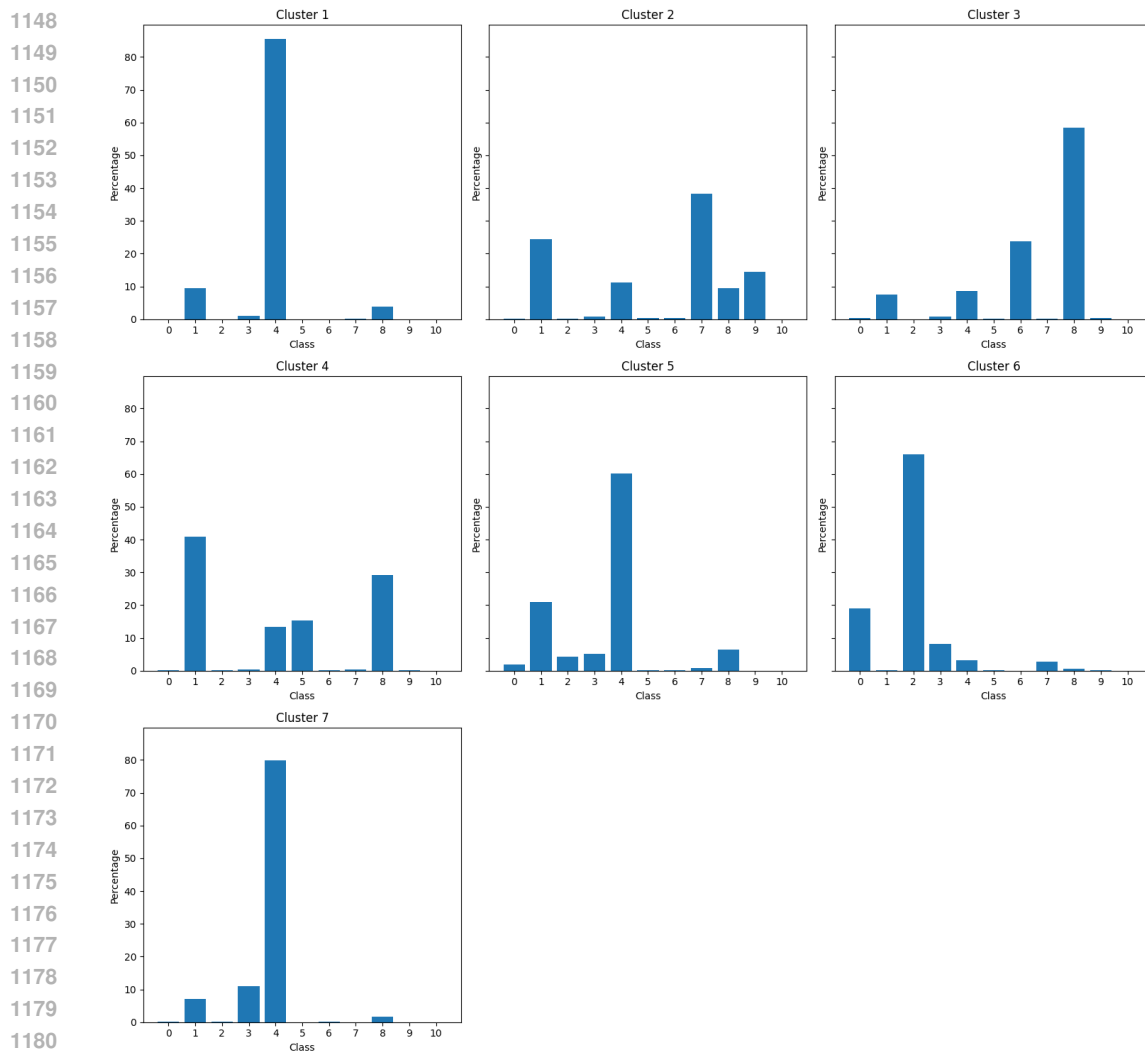


Figure 6: Class distributions across clusters for the Amazon dataset.

clustering captures shared structural or feature-based similarities among nodes, which can serve as a foundation for localized training. However, the presence of mixed-class clusters highlights the need for downstream methods that can handle intra-cluster class diversity effectively.