

National Cheng Kung University

MS Degree in Intelligent Computing

Master's Thesis

融合資料補全、模型剪枝與知識蒸餾聯合學習之輕量穩健圖神經網路

iPaDGNN: Joint Learning with Imputation, Pruning, and Distillation for Robust  
Lightweight Graph Neural Networks

學生：李旻昊

Student：Man-Ho Li

指導老師：李政德 博士

Advisor：Dr. Cheng-Te Li

July 2025

# 國立成功大學

National Cheng Kung University

## 碩士論文合格證明書

Certificate of Approval for Master's Thesis

題目：融合資料補全、模型剪枝與知識蒸餾聯合學習之輕量穩健圖神經網路

Title: iPaDGNN: Joint Learning with Imputation, Pruning, and Distillation for Robust Lightweight Graph Neural Networks

研究生：李旻昊

本論文業經審查及口試合格特此證明

This is to certify that the Master's Thesis of LI MAN HO

論文考試委員：

has passed the oral defense under the decision of the following committee members

莊坤達	李政德	洪弘毅
朱威達	黃仁暉	

指導教授 Advisor(s)：李政德

單位主管 Chair/Director：智慧運算碩士學位學程主任 舒宇宸

(單位主管是否簽章授權由各院、系(所、學位學程)自訂)  
(The certificate must be signed by the committee members and advisor(s). Each department/graduate institute/degree program can determine whether the chair/director also needs to sign.)

2025/7/24

## 摘要

圖神經網路 (Graph Neural Networks, GNNs) 於多種圖結構資料分析任務中表現卓越。然而，其在實務應用上仍面臨三大環環相扣的挑戰：資料不完整、高昂的訓練成本，以及過高的推論開銷。儘管現有方法各自處理單一問題，卻未能提供一個整體的解決方案。

為解決此問題，本研究提出 **iPaD** (Joint Learning with Imputation, Pruning, And Distillation for Robust Lightweight Graph Neural Networks)，一個建立在「**協同優化**」核心原則之上的統一框架，旨在對 GNN 的完整流程進行聯合優化。本框架透過一套連貫的流程來實現此設計哲學：首先，以高效的特徵補全技術確保資料的完整性；接著，採用一創新的代理剪枝策略，打造出一個輕量且強健的 GNN 教師模型；最終，再將其知識蒸餾至一個更高效的 MLP 學生模型中，以達成快速推論的目的。

大量的實驗證明，iPaD 不僅能在高達 99.5% 特徵缺失的資料集上保持高準確性，其**推論速度更比基準 GNN 教師模型提升超過五倍**。本研究證明，此一整體性的設計方法，為建構能夠實際應用於真實世界、兼具魯棒性、高效率與可部署性的 GNN，提供了一條實用且強而有力的途徑。

**關鍵字：**圖神經網路, 資料補全, 模型剪枝, 知識蒸餾, 協同學習

# Abstract

Graph Neural Networks (GNNs) have demonstrated remarkable performance in various graph-structured data analysis tasks. However, their practical deployment is severely hindered by a trio of interconnected challenges: incomplete data, high training costs, and prohibitive inference overhead. While existing methods tackle these issues in isolation, they fail to provide a holistic solution.

To address this gap, we introduce **iPaD** (Joint Learning with **I**mputation, **P**runing, **A**nd **D**istillation for Robust Lightweight Graph Neural Networks), a unified framework built on the core principle of **synergistic optimization** across the entire GNN pipeline. Our framework actualizes this philosophy by first ensuring data integrity with an efficient feature imputation technique. It then employs a novel proxy-based pruning strategy to create a lightweight yet robust GNN teacher, whose knowledge is subsequently distilled into an even more efficient MLP student for fast inference.

Extensive experiments show that iPaD not only maintains high accuracy on datasets with up to 99.5% of features missing but also achieves over a **5x reduction in inference time** compared to its baseline GNN teacher. Our work demonstrates that this holistic approach provides a practical and powerful pathway for building robust, efficient, and deployable GNNs for real-world applications.

**Keyword:** Graph Neural Network, Data Imputation, Model Pruning, Knowledge Distillation, Collaborative Learning

## 誌謝

在本論文研究與撰寫的過程中，我獲得了許多寶貴的指導與支持，謹在此致上誠摯的感謝。

首先，衷心感謝指導教授李政德老師，在整個研究過程中給予我嚴謹的學術訓練、耐心的指導與不吝的建議；無論是在研究方向的確立、方法設計，或是論文撰寫的細節上，都使我受益良多。同時感謝口試委員黃仁暉老師、莊坤達老師、朱威達老師與陳弘軒老師，撥冗審閱本文並提供寶貴意見，使本研究得以更加完善。

感謝實驗室的所有成員，在研究構想討論、程式開發與實驗執行等方面提供了大量的協助與啟發。與大家的交流與合作，不僅豐富了我的學術視野，也讓研究過程充滿樂趣。此外，特別感謝好友曾品昕，在研究受挫與趕稿階段給予我堅定的支持與陪伴，讓我得以調整步伐、持續前行。

最後，感謝我的家人一路以來的支持與鼓勵，給予我穩定的後盾與無條件的關愛，讓我能專注於學術研究並堅持到最後。

謹以此文，向所有在本研究中給予幫助與支持的人士，致上誠摯的謝意。

# Table of Contents / 目錄

摘要	i
Abstract	ii
誌謝	iii
Table of Contents / 目錄	iv
List of Tables / 表格	vii
List of Figures / 圖片	xi
<b>Chapter 1. Introduction</b>	<b>1</b>
1.1. Motivation . . . . .	1
1.2. Challenges . . . . .	2
1.2.1. Incomplete Data . . . . .	2
1.2.2. High Training Cost . . . . .	2
1.2.3. Inefficient Inference . . . . .	2
1.2.4. Lack of Unified Solutions . . . . .	2
1.3. Research Goal . . . . .	3
1.4. Contributions . . . . .	3
1.5. Thesis Organization . . . . .	4
<b>Chapter 2. Related Work</b>	<b>5</b>
2.1. Feature Imputation in Graph Neural Networks . . . . .	5
2.2. Model Pruning in Graph Neural Networks . . . . .	6
2.3. Knowledge Distillation for Graph Models . . . . .	7
2.4. High-Level Comparison of SOTA GNNs . . . . .	8
2.5. Chapter Summary . . . . .	9
<b>Chapter 3. Problem Statement</b>	<b>11</b>
3.1. Problem Background . . . . .	11
3.2. Notations . . . . .	12
3.3. Formal Statement . . . . .	12
3.4. Summary . . . . .	14
<b>Chapter 4. Methodology</b>	<b>15</b>
4.1. Framework Overview . . . . .	15
4.2. APCFI: Approximate Pseudo-Confidence Feature Imputation . . . . .	17
4.2.1. Motivation and Core Idea . . . . .	17
4.2.2. Technical Workflow . . . . .	18
4.2.3. Implementation and Pseudocode . . . . .	22

4.3.	TunedGNN: Architecture Optimization for Robustness . . . . .	25
4.3.1.	Key Design Components . . . . .	25
4.4.	MPP: Mirror Projection Pruning . . . . .	27
4.4.1.	Mirror Projection Pruning (MPP): Theory and Workflow . . . . .	27
4.5.	Teacher Training . . . . .	31
4.6.	MP-KRD: Mirror Projection Knowledge-inspired Reliable Distillation . . . . .	32
4.6.1.	Motivation: The Need for Reliable Distillation . . . . .	33
4.6.2.	The MP-KRD Workflow . . . . .	33
4.6.3.	Deployment . . . . .	37
4.7.	Innovation Summary . . . . .	38
<b>Chapter 5.</b>	<b>Experiment</b>	<b>41</b>
5.1.	Experimental Setup . . . . .	41
5.1.1.	Datasets . . . . .	41
5.1.2.	Compared Methods . . . . .	42
5.1.3.	Training Details . . . . .	42
5.1.4.	Missing Data Settings . . . . .	42
5.1.5.	Dataset Split Strategy . . . . .	42
5.2.	Validating APCFI: A Superior Balance of Accuracy and Efficiency . . . . .	43
5.2.1.	Efficiency Comparison . . . . .	43
5.2.2.	Accuracy under High Missing Rates . . . . .	44
5.3.	Baseline Performance: Establishing Competitiveness with Complete Data . . . . .	46
5.4.	Validating MPP: Quantifying Gains in Training and Inference Efficiency . . . . .	48
5.4.1.	Impact on Computational Costs . . . . .	48
5.4.2.	Accuracy-Efficiency Trade-off and the “Sweet Spot” . . . . .	49
5.5.	Stress Test: Evaluating Robustness under Extreme Data Corruption . . . . .	50
5.5.1.	Performance under Edge Missing . . . . .	50
5.5.2.	Performance under Feature Missing . . . . .	52
5.6.	Dissecting the Framework: An Ablation Study on Core Modules . . . . .	56
5.6.1.	Ablation Results with GCN Backbone . . . . .	56
5.6.2.	Ablation Results with SAGE Backbone . . . . .	57
5.6.3.	Conclusion of Ablation Study . . . . .	58
5.7.	Component-Level Validation: Internal Ablation Studies . . . . .	59
5.7.1.	Validating the Robustness of the TunedGNN Backbone . . . . .	59
5.7.2.	Validating the APCFI Design: A Component-wise Ablation Study . . . . .	59
5.7.3.	Validating the MP-KRD Design . . . . .	61
<b>Chapter 6.</b>	<b>Discussion and Conclusion</b>	<b>65</b>
6.1.	Summary of Findings and Contributions . . . . .	65
6.2.	Limitations . . . . .	66
6.3.	Future Work . . . . .	66
6.4.	Concluding Remarks . . . . .	67
	<b>References</b>	<b>68</b>
	<b>Appendix A. Hyperparameters</b>	<b>73</b>
A.1.	Analysis of APCFI Hyperparameters . . . . .	73

A.2. Analysis of MP-KRD Hyperparameters . . . . .	74
<b>Appendix B. Generalization Ability</b>	<b>75</b>
B.1. Generalization of iPaD to Various GNN Architectures . . . . .	75
<b>Appendix C. Effect of MPP Pruning on Teacher Models</b>	<b>77</b>
C.1. Rationale for Applying MPP Pruning to Teacher Models . . . . .	77
<b>Appendix D. Joint Training Imputer</b>	<b>78</b>
D.1. Joint Training of FP-based Imputer with MPP and KRD . . . . .	78
<b>Appendix E. Complex Missing Conditions</b>	<b>79</b>
E.1. Performance under Combined Feature and Edge Missing Conditions . . . .	79



## List of Tables / 表格

2.1	<b>Comparison of representative feature imputation methods.</b>	6
2.2	<b>Comparison of representative GNN pruning methods.</b>	7
2.3	<b>Comparison of representative knowledge distillation methods for graph models.</b>	8
2.4	<b>High-level comparison of recent SOTA GNN models on efficiency, and robustness.</b> The asterisk (*) indicates that robustness is not the primary focus of the original paper but is evaluated in our study.	9
2.5	<b>Holistic comparison of representative methods across the core challenges of Completeness, Training Efficiency, and Scalability.</b> Our proposed iPaD is the only framework designed to excel in all three aspects.	10
3.1	<b>Summary of main notations used throughout the paper.</b> This table lists the symbols, definitions, and descriptions for graph structure, feature matrices, and related variables.	12
4.1	<b>Imputation Module (APCFI) Innovations.</b> Detailed breakdown of the APCFI components, specifying the source of each technique and categorizing the level of innovation.	39
4.2	<b>Pruning Module (MPP) Innovations.</b> Comprehensive summary of pruning techniques, including parameter transfer, channel pruning, and the proposed proxy pruning strategy using MLP to accelerate GNN training and enhance KD.	39
4.3	<b>Distillation Module (MP-KRD) Innovations.</b> Overview of MP-KRD components, covering KDMLP, curriculum learning, the proposed homogeneous student structure, and the improved CWDLoss for GNN adaptation.	40
5.1	<b>Statistics of Benchmark Datasets Used in Node Classification Tasks.</b> The table summarizes the number of nodes, edges, features, classes, and the evaluation metric for each dataset.	41
5.2	<b>Runtime comparison of different feature imputation methods across datasets.</b> The table reports the runtime (in seconds) for FP, PCFI, and APCFI methods on several benchmark datasets. APCFI achieves comparable efficiency to FP and significantly outperforms PCFI in computational cost.	44
5.3	<b>Performance comparison of various feature imputation methods under structural and uniform missing scenarios.</b> The table reports the accuracy (%) on five datasets for each method. Methods marked with * are directly introduced from the PCFI paper. APCFI achieves the best or near-best performance and average rank under both missing types. The best, second-best, and third-best results are highlighted in <b>green</b> , <b>cyan</b> , and <b>red</b> .	45
5.4	<b>SOTA comparison on complete classic datasets.</b> The best, second-best, and third-best results are highlighted in <b>green</b> , <b>cyan</b> , and <b>red</b> , respectively. Results marked with * are from TunedGNN’s paper.	46

5.5	<b>SOTA comparison on complete medium datasets.</b> The best, second-best, and third-best results are highlighted in <b>green</b> , <b>cyan</b> , and <b>red</b> , respectively. Results marked with * are from TunedGNN’s paper. . . . .	47
5.6	<b>Impact of MPP pruning on FLOPs and inference time of the GCN model.</b> Increasing the prune rate substantially reduces both computational cost (GFLOPs) and inference time, demonstrating the efficiency gains enabled by the MPP pruning strategy. . . . .	49
5.7	<b>Test accuracy of TunedGCN [1] on various datasets under different edge missing ratios.</b> The table reports accuracy under full edge, 30%, 60%, and 90% edge missing conditions. Average accuracy and relative decreases are shown at the bottom. . . . .	51
5.8	<b>Test accuracy of iPaD-GCN on various datasets under different edge missing ratios.</b> The table reports accuracy under full edge, 30%, 60%, and 90% edge missing conditions. Average accuracy and relative decreases are provided at the bottom. . . . .	52
5.9	<b>Comparison under 99.5% uniform feature missing: vanilla GNN, MPP-Tuned, and iPaD with common imputation methods on classic datasets.</b> Model names indicate the imputation strategy. The best results are highlighted in <b>green</b> , the second-best in <b>cyan</b> , and the third-best in <b>red</b> . . . . .	53
5.10	<b>Comparison under 99.5% uniform feature missing: vanilla GNN, SG-Former, Polynormer, TunedGNN, MPP-Tuned, and iPaD with common imputation methods on medium datasets.</b> Model names indicate the imputation strategy. The best results are highlighted in <b>green</b> , the second-best in <b>cyan</b> , and the third-best in <b>red</b> . . . . .	53
5.11	<b>Comparison under 99.5% structural feature missing: vanilla GNN, MPP-Tuned, and iPaD with common imputation methods on classic datasets.</b> Model names indicate the imputation strategy. The best results are highlighted in <b>green</b> , the second-best in <b>cyan</b> , and the third-best in <b>red</b> . . . . .	54
5.12	<b>Comparison under 99.5% structural feature missing: vanilla GNN, SG-Former, Polynormer, TunedGNN, MPP-Tuned, and iPaD with common imputation methods on medium datasets.</b> Model names indicate the imputation strategy. The best results are highlighted in <b>green</b> , the second-best in <b>cyan</b> , and the third-best in <b>red</b> . . . . .	55
5.13	<b>Ablation study of tuned, MPP, and MP-KRD modules on the Photo and CS datasets (GCN Model).</b> The table reports accuracy, training speed (epoch/s), and inference time (ms) for different combinations of modules. When MP-KRD is used, (*) indicates the speed measured during knowledge distillation. . . . .	57
5.14	<b>Ablation study of tuned, MPP, and MP-KRD modules on the Photo and CS datasets (GraphSAGE Model).</b> The table reports accuracy, training speed (epoch/s), and inference time (ms) for different combinations of modules. When MP-KRD is used, (*) indicates the speed measured during knowledge distillation. . . . .	58

5.15	Ablation study of APCFI components on runtime performance. The table reports the runtime (in seconds) for the two main stages of imputation (SPD calculation and Diffusion). Each row corresponds to enabling a key innovation, demonstrating how each component contributes to the final efficiency.	60
5.16	<b>Ablation study of soft-PC in APCFI.</b> The table reports test accuracy on Cora, CiteSeer, PubMed, Photo, and Computers datasets for the APCFI method with and without soft-PC. Avg Accuracy indicates the average accuracy across all datasets.	61
5.17	<b>Ablation study of KRD, MP-MLP, and CWD-loss components in the MP-KRD module on classic datasets using the GCN model.</b> The table reports test accuracy on Cora, CiteSeer, and PubMed for different component combinations. Avg Accuracy indicates the average accuracy and relative improvement. The best, and second-best results are highlighted in <b>green</b> , and <b>red</b> .	62
5.18	<b>Ablation study of KRD, MP-MLP, and CWD-loss components in the MP-KRD module on classic datasets using the SAGE model.</b> The table reports test accuracy on Cora, CiteSeer, and PubMed for different component combinations. Avg Accuracy indicates the average accuracy and relative improvement. The best, and second-best results are highlighted in <b>green</b> , and <b>red</b> .	62
5.19	<b>Ablation study of KRD, MP-MLP, and CWD-loss components in the MP-KRD module on medium-scale datasets using the GCN model.</b> The table reports test accuracy on Computers, Photo, WikiCS, CS, and Physics datasets for different component combinations. Avg Accuracy indicates the average accuracy and relative change. The best, and second-best results are highlighted in <b>green</b> , and <b>red</b> .	63
5.20	<b>Ablation study of KRD, MP-MLP, and CWD-loss components in the MP-KRD module on medium-scale datasets using the SAGE model.</b> The table reports test accuracy on Computers, Photo, WikiCS, CS, and Physics datasets for different component combinations. Avg Accuracy indicates the average accuracy and relative change. The best, and second-best results are highlighted in <b>green</b> , and <b>red</b> .	64
B.1	<b>Efficiency improvement of the iPaD framework over various SOTA models across four datasets.</b> For each dataset (Computers, Photo, WikiCS, CS), accuracy (%) is reported alongside inference time (ms) in parentheses.	76
C.1	<b>Training efficiency improvement of MPP pruning over various SOTA models.</b> For each dataset (Computers, Photo, WikiCS, CS), accuracy (%) is reported alongside training speed (epoch/s) in parentheses.	77
D.1	<b>Performance comparison of FP and trainable FP under high missing rates.</b> Node classification accuracy (%) on six datasets (Cora, CiteSeer, Pubmed, Computer, Photo, WikiCS) with a uniform feature missing rate of 90%.	78

E.1	<b>Performance under complex extreme missing data conditions.</b> Node classification accuracy (%) on four datasets (Computer, Photo, WikiCS, Physics) with uniform feature missing rate of 99.5% and edge missing rate of 60%. Shaded rows indicate our iPaD variants, compared against SOTA models with feature propagation (FP) preprocessing. . . . .	80
-----	---	----

## List of Figures / 圖片

4.1	<b>The overall architecture of the iPaD framework.</b> The process begins with imputing the incomplete feature matrix via <b>APCFI</b> . In parallel, a powerful GNN teacher is trained on the recovered data, while the <b>MPP</b> module generates a pruned, lightweight MLP student. Finally, the <b>MP-KRD</b> module distills knowledge from the teacher to the student, producing a final model that is both robust and efficient for inference. . . . .	15
4.2	<b>The overall workflow of the APCFI module.</b> It takes a graph with an incomplete feature matrix as input and leverages ASDE and Dynamic Joint Channel Diffusion (DJCD) to output a robustly recovered feature matrix for downstream tasks. . . . .	17
4.3	<b>Illustration of the Approximate Shortest Distance Estimation (ASDE) process.</b> The algorithm begins with observed features (step-0) and propagates shortest-path information through the graph. Each entry in the ASDE matrix reflects the shortest hop-distance to an observed feature within the same channel. . . . .	19
4.4	<b>The iterative process of Dynamic Joint Channel Diffusion.</b> In each of the K steps, a diffusion operation updates the missing values based on neighbors, followed by a “Fixed” step where the original, observed feature values are reset to prevent oversmoothing and preserve data fidelity. . . . .	21
4.5	<b>Comparison of Tuned and Classic GNN layer structures.</b> (a) The classic GNN layer stacks message-passing and activation modules, typically with fewer than 3 layers. (b) The tuned GNN layer incorporates residual connections, normalization, activation, and dropout, allowing for deeper architectures (up to 11 layers) and improved robustness to incomplete data. . . . .	25
4.6	<b>The high-level workflow of Mirror Projection Pruning (MPP).</b> A large GNN architecture is projected to a proxy MP-MLP, which is then trained and pruned. The resulting sparse structure is projected back to define the final, efficient ‘Pruned GNN’. . . . .	27
4.7	<b>Structural illustration of SAGEConv and its MP-MLP projection.</b> Left: The SAGEConv layer consists of two linear branches, one operating on neighborhood-aggregated features (Propagate). Right: The SAGEConvMLP projection replaces neighborhood aggregation with a second self-branch, yielding a pure local MLP structure with identical parameterization. This one-to-one mapping is fundamental to efficient pruning and parameter transfer in MPP. . .	28
4.8	<b>Mirror Projection from GNN to MP-MLP.</b> Each GNN-layer is replaced by a corresponding ConvMLP layer, yielding an untrained MP-MLP that mirrors the GNN architecture. . . . .	29
4.9	<b>Proxy Training and Pruning in MP-MLP.</b> The MP-MLP is trained and pruned to yield a sparse, efficient student model. . . . .	29

4.10	<b>Inverse Projection from Pruned MP-MLP to Pruned GNN.</b> The pruning mask and weights are transferred directly, instantly producing a sparse, efficient GNN. . . . .	30
4.11	<b>Teacher Training Workflow.</b> The pruned GNN obtained from the MPP module is fully trained on the imputed feature matrix $\hat{\mathbf{X}}$ generated by APCFI. The resulting efficient model is used as the Teacher for knowledge distillation. . . . .	32
4.12	<b>The high-level view of the MP-KRD module.</b> It takes the trained GNN Teacher and the Pruned MP-MLP Student as input. A Reliable Label Sampler, guided by a curriculum, selects high-quality knowledge to distill, producing the final, efficient inference model. . . . .	32
4.13	<b>The workflow of the MP-KRD module.</b> The pruned GNN teacher provides sample reliability signals for curriculum-based distillation, guiding the training of the MP-MLP student via cross-entropy, logit KD (KRD), and feature KD (CWD) losses. . . . .	33
4.14	<b>Histograms of teacher-student agreement (a) and disagreement (b) plotted against the normalized entropy change.</b> Easier samples (low entropy change) show high agreement, while harder samples show high disagreement. . . . .	36
4.15	The empirical agreement probability ( $p_k$ ) is fitted with the power curve (red line) to dynamically determine the optimal curriculum “power <sub>fit</sub> ” parameter. . . . .	37
5.1	<b>Total training time (including the MPP step) of the MPP-pruned GCN model on the CS dataset under different prune rates.</b> Higher prune rates lead to a significant reduction in total training time. . . . .	49
5.2	Test accuracy of SAGE and GCN models on the Physics dataset under different MPP prune rates. The accuracy remains highly stable up to a prune rate of 0.7, revealing a wide “sweet spot” for efficiency gains without a significant performance penalty. . . . .	50
5.3	<b>Robustness comparison of different GCN configurations under varying feature missing ratios on the Physics dataset.</b> Both the ‘MPP-GCN(tuned)’ teacher and the final ‘iPaD-GCN’ student model demonstrate superior robustness, maintaining high accuracy even at a 90% missing ratio. In contrast, the baseline ‘GCN w/ Zero’ imputation suffers a catastrophic performance collapse. . . . .	56
5.4	<b>Robustness comparison of tuned and classic GCN/SAGE architectures on the Physics dataset.</b> All models are evaluated using Zero imputation for missing features, providing a uniform baseline for comparison. Both tuned GCN (green) and tuned SAGE (red) with MPP maintain higher accuracy than their classic counterparts as the missing ratio increases, demonstrating the impact of model tuning and pruning on robustness. . . . .	59
A.1	Contour plot for the key hyperparameters (alpha and beta) of the APCFI module on the Cora dataset. . . . .	73
A.2	Contour plot for the key hyperparameters of the MP-KRD module. Darker blue regions indicate higher accuracy. The plot reveals strong dependencies between parameters. . . . .	74

# Chapter 1

## Introduction

Graph Neural Networks (GNNs) have emerged as a fundamental paradigm for learning from graph-structured data [2, 3, 4], demonstrating state-of-the-art performance across diverse domains. By jointly leveraging both the topological structure of graphs and rich node-level features, GNNs are capable of capturing complex relational dependencies that are often inaccessible to traditional neural architectures [2, 3, 4]. Consider, for instance, a large-scale recommender system: GNNs can effectively model the intricate relationships between users and items to provide personalized recommendations. [5, 6]

However, deploying such powerful models in real-world scenarios is severely hindered by challenges. In the context of recommender systems, these challenges manifest as follows [5, 7, 8]: (1) many users have sparse interaction histories, leading to **incomplete data**; (2) training a GNN on billions of user-item interactions incurs **high training costs**; and (3) generating recommendations in real-time demands **low inference latency**, which often exceeds the capabilities of large GNNs. This thesis introduces a framework designed to holistically address these practical yet critical challenges.

### 1.1 Motivation

As the adoption of GNNs in high-stakes applications like recommender systems and social network analysis accelerates, the demand for robust and scalable graph learning becomes increasingly critical. In these practical settings, multiple challenges—including missing node features, incomplete graph structures, and limited computational resources—often co-occur and interact unpredictably [9, 10, 11].

Existing research often addresses these challenges separately, developing targeted solutions for either missing data, training efficiency, or inference cost independently. This siloed strategy, however, leads to a significant practical gap. A solution that excels at data imputation may still be too computationally expensive to train or deploy. Conversely, an efficient, pruned model might fail catastrophically when faced with the noisy and incomplete data typical of real-world environments. Consequently, a unified framework that addresses these challenges **simultaneously** is not merely a convenience, but a necessity for robust, practical

deployment. To fully appreciate the need for this integrated approach, we first dissect the distinct yet interrelated challenges limiting the practical use of GNNs.

## **1.2 Challenges**

Despite recent advances in GNN research, several persistent challenges continue to impede the reliable and efficient deployment of GNNs in real-world scenarios:

### **1.2.1 Incomplete Data**

Real-world graphs are often incomplete, containing absent node features or missing edges. In our recommender system example, this corresponds to new users with no purchase history. This missing information fundamentally disrupts the message-passing mechanism of GNNs, often leading to suboptimal representations and substantially degraded predictive performance [12, 13]. This data incompleteness not only impacts the final accuracy but also complicates the training of deeper, more expressive models.

### **1.2.2 High Training Cost**

The recursive neighborhood aggregation and large parameter space inherent to many powerful GNNs cause training costs to scale rapidly with graph size and model depth [4, 14, 15]. Training a GNN on a graph with millions of users and items, for example, demands significant computational and memory resources, making it a costly and time-consuming endeavor, especially in resource-constrained settings. This high training cost is further exacerbated when models need to be large and deep to handle noisy, incomplete data. Ultimately, even if a powerful model can be trained, its practical utility is limited by its inference efficiency.

### **1.2.3 Inefficient Inference**

At inference time, conventional GNNs require multi-hop neighborhood aggregation for each prediction. To recommend a single item to a user, the model might need to access information from thousands of neighboring nodes, introducing considerable computational overhead [5, 16, 17]. This inefficiency poses a significant obstacle for deployment in latency-sensitive applications like real-time recommendation. This often drives the use of smaller, simpler models, which in turn are more vulnerable to the data incompleteness issues discussed earlier, thus creating a vicious cycle of compromises.

### **1.2.4 Lack of Unified Solutions**

These challenges reveal that most existing solutions focus on addressing individual issues in



isolation. There is a clear lack of a holistic solution capable of addressing the full spectrum of real-world constraints from data imperfection to deployment efficiency. This gap underscores the need for the unified and modular framework proposed in this thesis.

### 1.3 Research Goal

In direct response to the challenges of data incompleteness, high training costs, and inference inefficiency, the primary objective of this research is to design and develop a unified, modular framework. Our central hypothesis is that a framework **synergistically integrating** robust imputation, pre-training pruning, and architecture-aligned distillation could overcome these limitations in a holistic way.

Specifically, this thesis aims to develop **iPaD** (Joint Learning with **Imputation, Pruning, And Distillation** for Robust Lightweight Graph Neural Networks), a framework that provides a comprehensive, end-to-end solution by:

1. Developing a feature imputation method **Approximate Pseudo-Confidence Feature Imputation (APCFI)** that is both robust to severe sparsity and computationally efficient enough for large-scale graphs.
2. Designing a novel proxy-based pruning strategy **Mirror Projection Pruning (MPP)** that reduces model complexity and accelerates training before the most expensive steps, by leveraging the isomorphism between GNNs and MLPs.
3. Creating a reliable knowledge distillation mechanism **Mirror Projection Knowledge-inspired Reliable Distillation (MP-KRD)** that efficiently transfers knowledge from a powerful GNN teacher to a lightweight, fast-inference MLP student, thus completing the pipeline from robust training to efficient deployment.

Through this unified approach, we aim to enable the practical deployment of GNNs across a wide range of real-world applications, ensuring strong performance even in challenging and imperfect environments.

### 1.4 Contributions

This thesis makes the following key contributions:

1. **A Unified, Synergistic Framework:** We introduce **iPaD**, a novel framework that systematically and synergistically addresses the three principal challenges of real-world GNN deployment. Its modules are designed to be complementary, where the output of one stage (e.g., the pruned student model from MPP) serves as the ideal input for the next (e.g., the distillation process in MP-KRD).

2. **A Highly Efficient Imputation Method (APCFI):** We develop **APCFI**, a feature imputation method that achieves comparable accuracy to state-of-the-art methods but with a significant **efficiency breakthrough**, reducing computation time by over 200x on benchmark datasets and making robust imputation practical for large graphs.
3. **A Novel Proxy Pruning Paradigm (MPP):** We propose **MPP**, a pre-training pruning strategy that introduces a new paradigm of proxy-based pruning. By projecting a GNN to an isomorphic MLP, it decouples the complex GNN pruning task from training and effectively leverages mature MLP pruning techniques, significantly reducing computational overhead.
4. **An End-to-End Distillation Pipeline (MP-KRD):** We design **MP-KRD**, an architecture-aligned knowledge distillation mechanism that completes the optimization pipeline. It uniquely uses the pruned MLP generated by MPP as a student, enabling highly efficient inference while preserving the relational reasoning capabilities of the original GNN teacher.
5. **Empirical Validation of Robust Architectures:** We are the first to systematically demonstrate that simple architectural modifications in a **TunedGNN**, such as adding LayerNorm and Dropout, not only improve accuracy on complete data but also serve as a critically effective backbone for enhancing GNN robustness against severe feature missingness.

## 1.5 Thesis Organization

This thesis is organized as follows:

- **Chapter 2: Related Work** reviews prior research on GNNs with incomplete data, feature imputation, model pruning, and knowledge distillation.
- **Chapter 3: Problem Statement** formally defines the key challenges addressed in this work and introduces the necessary notations and settings.
- **Chapter 4: Methodology** details the proposed iPaD framework and its core modules: APCFI, MPP, TunedGNN, and MP-KRD.
- **Chapter 5: Experiments** presents the experimental setup, benchmark datasets, baseline methods, and evaluation metrics, followed by a comprehensive analysis of results.
- **Chapter 6: Conclusion** summarizes the main findings of this study and discusses potential directions for future research on robust and efficient graph neural networks.

## Chapter 2

### Related Work

This chapter reviews prior research in areas crucial to our work. Methods for handling incomplete data in GNNs, particularly feature imputation, are first reviewed. The subsequent sections discuss the landscape of GNN model pruning and knowledge distillation. For each area, the limitations of existing approaches are analyzed to precisely situate the contributions of the proposed iPaD framework.

#### 2.1 Feature Imputation in Graph Neural Networks

Feature imputation is a critical component for enhancing the robustness of Graph Neural Networks (GNNs) in the presence of missing node features—a common challenge in real-world graph-structured data [18, 19, 20, 21]. Existing approaches can be broadly classified into heuristic methods, learning-based models, and feature propagation techniques.

Heuristic methods, such as Zero, Mean, or Label Propagation (LP) Imputation, are computationally efficient but often fail to capture the complex feature distributions present in graphs [22], leading to suboptimal performance. Learning-based approaches, including GCNMF and PaGNN, employ deep neural network architectures to model and reconstruct missing features [12, 20]. Despite their theoretical expressiveness, these models can be limited by substantial computational overhead and may not outperform simpler methods in high-missingness regimes.

More recently, non-learning-based feature propagation methods have gained traction. Feature Propagation (FP) exploits the underlying graph structure to propagate known features, offering a strong balance of performance and efficiency [19]. Pseudo-Confidence Feature Imputation (PCFI) further advances this by introducing channel-wise pseudo-confidence to guide a two-stage diffusion process, achieving state-of-the-art imputation quality [21]. However, this high performance comes at a significant cost: PCFI’s channel-by-channel calculation leads to extremely high computational complexity, hindering its scalability to large graphs. This reveals a critical trade-off between imputation performance and computational efficiency in the current literature.

To address this efficiency-accuracy trade-off, **APCFI** is introduced to retain the high imputation quality of confidence-aware methods such as PCFI, while drastically reducing compu-

tational burden via Approximate Shortest Distance Estimation (ASDE) and Dynamic Joint Channel Diffusion (DJCD).

Table 2.1 provides a clear comparison of these imputation strategies.

Method	Imputation Performance	Computational Cost	Key Idea	Key Limitation
Zero/Mean	Low	Very Low	Constant Filling	Unable to capture characteristic distribution
LP [22]	Low/Medium	Very Low	Feature Replacement	Ignore existing features
PaGNN [12]	Medium	Very High	Learnable Method	Sensitive to noise, poor scalability
GCNMF [20]	Medium	Very High	Matrix Factorization, Learnable Method	Sensitive to noise, poor scalability
FP [19]	Medium	Low	Feature Propagation-based	Insensitive to differences between features
PCFI [21]	High	High	Confidence-aware Feature Propagation	Channel-by-channel calculation, poor scalability
<b>APCFI(Ours)</b>	<b>High</b>	<b>Low</b>	Approximate Confidence & Joint Channel Diffusion	Not applicable to heterogeneous graphs

**Table 2.1: Comparison of representative feature imputation methods.**

## 2.2 Model Pruning in Graph Neural Networks

As GNNs scale to handle increasingly large and complex graphs, model pruning has emerged as a crucial strategy to reduce computational costs during both training and inference. Existing GNN pruning paradigms include channel pruning (e.g., GCNP [23]), hardware-aware optimization (e.g., PruneGNN [24]), and methods that simultaneously prune graph structures and model weights during training (e.g., UGS [25], ICPG [26]).

However, these methods often suffer from two major limitations. First, they typically entangle the pruning process with the GNN training pipeline, introducing significant computational overhead and complex, model-specific algorithmic designs that can substantially **increase the total training cost**. Second, they are often proposed as specific, monolithic algorithms,

lacking the **scalability and flexibility** to incorporate a wide array of mature pruning techniques from other domains, such as the MLP literature.

To overcome these challenges of high overhead and low flexibility, we introduce **MPP** (Mirror Projection Pruning). Instead of being another monolithic algorithm, MPP proposes a new **proxy-based pruning paradigm**. Its core idea is to decouple pruning from GNN training by projecting the GNN to an isomorphic MLP, where mature, off-the-shelf pruning algorithms can be applied efficiently *before* the main training commences. The advantages of this paradigm are summarized in Table 2.2.

Method	Total Training Cost	Scalability	Key Idea
GCNP [23]	↑	✗	Using a novel LASSO regression formulation to prune model channels
PruneGNN [24]	↑	✗	Introduces novel SIMD-aware kernels and accelerate structured GNN pruning on parallel hardware
UGS [25]	↑	✗	Prune graphs and weights simultaneously during training for specific ensemble algorithms
ICPG [26]	↑	✗	Iterative co-pruning during training for specific ensemble algorithms
<b>MPP(Ours)</b>	↓	✓	A GNN proxy pruning framework that off-the-shelf MLP-compatible algorithms

**Table 2.2: Comparison of representative GNN pruning methods.**

### 2.3 Knowledge Distillation for Graph Models

Knowledge distillation (KD) is a pivotal technique for compressing GNNs, enabling the deployment of lightweight yet high-performing models. In the context of graph data, KD approaches can be broadly categorized into GNN-to-GNN and GNN-to-MLP distillation.

GNN-to-GNN methods, such as LSP [27] and TinyGNN [28], focus on transferring knowledge from a large GNN teacher to a smaller GNN student. While effective, they do not eliminate the graph dependency at inference time. To achieve maximum inference speed, GNN-to-MLP distillation has emerged as a promising direction. Methods like GLNN [29] show

that an MLP student, once distilled, can achieve competitive performance without access to the graph structure. More advanced techniques like KRD [30] introduce reliability-aware curriculum learning to handle potentially noisy knowledge from the teacher.

Despite these advances, existing GNN-to-MLP methods exhibit several gaps. First, the student MLP is typically an **arbitrary MLP**, lacking any intrinsic structural alignment with the GNN teacher. Second, the distillation process is designed **independently of any model pruning stage**, missing a key opportunity for synergistic optimization. Finally, while some methods address teacher noise, a more robust, pipeline-aware approach is needed.

Our proposed **MP-KRD** (Mirror Projection Knowledge-inspired Reliable Distillation) is designed specifically to fill these gaps. It creates a powerful synergy by using the **structurally-aligned, isomorphic MLP** generated by our MPP stage as the student. This ensures architectural consistency and maximizes the effectiveness of the knowledge transfer, which is further enhanced by a reliability-aware curriculum. The unique characteristics of MP-KRD are summarized in Table 2.3.

Method	Student Source	Feature Alignment	Handles Noisy Teacher	Key Idea
LSP [27]	Small GNN	✗	✗	Transfer the topological structure
GNN-SD [31]	Teacher GNN	✓	✗	GNNs Self-Distillation framework
TinyGNN [28]	Small GNN	✗	✗	Distillation with fewer parameter GNNs
GLNN [29]	Arbitrary MLP	✗	✗	Learning predictions from GNNs
KRD [30]	Arbitrary MLP	✗	✓	Reliability-based learning
<b>MP-KRD(Ours)</b>	<b>Isomorphic MLP</b>	✓	✓	Structural Alignment + Reliability Course Learning

**Table 2.3: Comparison of representative knowledge distillation methods for graph models.**

## 2.4 High-Level Comparison of SOTA GNNs

Beyond the specific areas of imputation and compression, it is also instructive to compare our approach against the broader landscape of recent state-of-the-art (SOTA) GNNs, particularly concerning the trade-off between performance and efficiency.

Table 2.4 provides a qualitative comparison of prominent SOTA models. As the table illustrates, while most recent Transformer-based models like GraphGPS[32] and NAGphormer[33] achieve high predictive performance, they often do so at the cost of low computational efficiency. Our **MP-TunedGNN** (the teacher model in our framework), by integrating the MPP module, stands out by maintaining high performance while significantly boosting efficiency compared to both its “TunedGNN” backbone and other SOTA models. This highlights a key contribution of our work in addressing the prevalent performance-efficiency trade-off in modern GNN architectures.

Method	Efficiency	Robustness
GraphGPS [32]	Low	-
NAGphormer [33]	Medium	-
Exphormer [34]	Medium	-
GOAT [35]	Medium	-
NodeFormer [36]	Medium	-
SGFormer [37]	Medium	-
Polynormer [38]	Medium	-
TunedGNN [1]	High	✓*

**Table 2.4: High-level comparison of recent SOTA GNN models on efficiency, and robustness.** The asterisk (\*) indicates that robustness is not the primary focus of the original paper but is evaluated in our study.

## 2.5 Chapter Summary

In this chapter, we have reviewed prior art across feature imputation, model pruning, and knowledge distillation, identifying the key limitations of existing methods in each domain. To synthesize these findings and crystallize the motivation for our work, we present a holistic comparison in Table 2.5. This table evaluates a wide spectrum of methods against the three core challenges identified in this thesis: handling data **Completeness**, achieving high **Training Efficiency**, and ensuring model **Scalability** for inference.

The table clearly illustrates the “siloe” nature of existing research. Imputation methods like PCFI excel at “Completeness” but do not address efficiency. High-performance SOTA models like GraphGPS [32] have low “Training Efficiency” and “Scalability”. Distillation methods like KRD offer “Scalability” but do not handle data “Completeness”.

Method	Completeness	Training Efficiency	Scalability
PaGNN [12]	✓	-	-
GCNMF [20]	✓	-	-
FP [19]	✓	-	-
PCFI [21]	✓	-	-
GraphGPS [32]	-	✓	-
NAGphormer [33]	-	✓	-
Expormer [34]	-	✓	-
GOAT [35]	-	✓	-
NodeFormer [36]	-	✓	-
SGFormer [37]	-	✓	-
Polynormer [38]	-	✓	-
TunedGNN [1]	✓*	✓	-
LSP [27]	-	-	✓
GNN-SD [31]	-	-	-
TinyGNN [28]	-	-	✓
GLNN [29]	-	-	✓
KRD [30]	-	-	✓
<b>iPaD(Ours)</b>	✓	✓	✓

**Table 2.5: Holistic comparison of representative methods across the core challenges of Completeness, Training Efficiency, and Scalability.** Our proposed iPaD is the only framework designed to excel in all three aspects.



## Chapter 3

### Problem Statement

#### 3.1 Problem Background

Graph Neural Networks (GNNs) have achieved remarkable success in a wide range of real-world applications, such as social network analysis, recommender systems, and bioinformatics [16, 39, 40]. However, their practical performance is highly contingent on the completeness of the input graph data [41, 42]. In realistic scenarios, graphs are frequently afflicted by missing node features or absent edges, stemming from factors like privacy constraints, sensor failures, or data collection limitations [43, 44]. Such missing information fundamentally disrupts the message-passing mechanisms central to GNN architectures, often resulting in suboptimal node representations and diminished performance on downstream tasks.

Missingness in graph data typically manifests in two primary forms:

- Uniform missingness refers to randomly distributed missing features throughout the graph, often caused by unpredictable device malfunctions or sporadic data loss.
- Structural missingness arises when missing data is concentrated within specific sub-graphs or clusters of nodes, commonly due to systematic issues such as privacy restrictions or selective sampling policies.

In addition to data incompleteness, the scaling of GNNs to large and complex graphs introduces substantial challenges in terms of training cost (computational resources and time) and inference cost (the need for extensive neighborhood aggregation at test time) [45, 46]. Existing approaches tend to focus on addressing a single aspect—such as feature imputation, model pruning, or architectural tuning—without providing a unified and modular solution that jointly addresses the intertwined challenges of robustness, efficiency, and scalability under incomplete data.

To overcome these practical challenges, a comprehensive framework capable of robust feature imputation, resource-efficient training, and cost-effective inference is essential for deploying GNNs in real-world settings characterized by high levels of missingness and limited computational resources [47].

### 3.2 Notations

Notation	Description
$\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathbf{A})$	An undirected graph, where $\mathcal{V}$ is the set of nodes, $\mathcal{E}$ is the set of edges, and $\mathbf{A}$ is the adjacency matrix.
$\mathcal{V} = \{v_i\}_{i=1}^N$	The set of $N$ nodes in the graph.
$\mathbf{A} \in \{0, 1\}^{N \times N}$	The adjacency matrix; $\mathbf{A}_{ij} = 1$ if there is an edge between $v_i$ and $v_j$ , otherwise 0.
$\mathbf{X} = [x_{i,d}] \in \mathbb{R}^{N \times F}$	Node feature matrix; $x_{i,d}$ is the $d$ -th feature of node $v_i$ .
$\mathcal{N}(v_i)$	The set of neighbors of node $v_i$ .
$\mathbf{M} \in \{0, 1\}^{N \times F}$	Mask matrix: 1 for observed features, 0 for missing.
$V_k^{(d)}$	Set of nodes with known $d$ -th feature (source nodes).
$V_u^{(d)}$	Set of nodes with unknown $d$ -th feature (missing nodes).
$x^{(d)} = \begin{bmatrix} x_k^{(d)} \\ x_u^{(d)} \end{bmatrix}$	Graph signal for $d$ -th channel, partitioned into known and unknown.
$\mathbf{A}^{(d)} = \begin{bmatrix} \mathbf{A}_{kk}^{(d)} & \mathbf{A}_{ku}^{(d)} \\ \mathbf{A}_{uk}^{(d)} & \mathbf{A}_{uu}^{(d)} \end{bmatrix}$	Adjacency matrix reordered for known/unknown nodes for channel $d$ .
$\hat{\mathbf{X}} = [\hat{x}_{i,d}]$	The recovered (imputed) feature matrix.

**Table 3.1: Summary of main notations used throughout the paper.** This table lists the symbols, definitions, and descriptions for graph structure, feature matrices, and related variables.

### 3.3 Formal Statement

Given a graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathbf{A}, \mathbf{X})$  with node feature matrix  $\mathbf{X} \in \mathbb{R}^{N \times F}$  and a missingness mask  $\mathbf{M} \in \{0, 1\}^{N \times F}$ , the objective of this work is to develop a unified framework that addresses the following intertwined challenges in graph representation learning under incomplete data:

#### 1. Robust Feature Imputation:

Given that a subset of features  $x_{i,d}$  is missing (i.e.,  $\mathbf{M}_{i,d} = 0$ ), the goal is to recover the missing values  $\hat{x}_{i,d}$  such that the imputed feature matrix  $\hat{\mathbf{X}}$  maximizes the downstream task performance (e.g., node classification accuracy) under varying missingness patterns (uniform/structural).

## 2. Resource-Efficient Training:

Let  $C_{\text{train}}(f)$  denote the computational cost (e.g. Time or FLOPs) of training a GNN-based model  $f$  on  $\mathcal{G}$ . The challenge is to minimize  $C_{\text{train}}(f)$  without sacrificing the performance on incomplete graphs, ideally through effective pruning or architectural optimization.

## 3. Cost-Effective Inference:

Define  $C_{\text{infer}}(f)$  as the inference cost for deploying  $f$  on new, potentially incomplete graphs. The target is to minimize  $C_{\text{infer}}(f)$  (e.g., reducing message passing or memory overhead) while ensuring high accuracy and reliability, possibly by leveraging knowledge distillation or graph-to-MLP transformation [29, 30, 48].

The joint optimization problem underlying our framework can be formalized as follows:

$$\begin{aligned} & \max_{f, \hat{\mathbf{X}}, \hat{\mathcal{A}}} \quad \text{TaskAcc}(f(\hat{\mathbf{X}}, \hat{\mathcal{A}})) \\ & \min_f \quad C_{\text{train}}(f) \\ & \min_f \quad C_{\text{infer}}(f) \\ & \text{s.t.} \quad \begin{cases} \mathbf{M}_{i,d} = 0 \implies \hat{x}_{i,d} \text{ is imputed} \\ \text{MissingRate}(\mathbf{M}) \leq \rho_X \\ \text{MissingRate}(\mathbf{A}) \leq \rho_A \end{cases} \end{aligned}$$

This joint optimization objective seeks to maximize the downstream task accuracy,  $\text{TaskAcc}(f(\hat{\mathbf{X}}, \hat{\mathcal{A}}))$ , by simultaneously optimizing the model  $f$ , the imputed feature matrix  $\hat{\mathbf{X}}$ , and the imputed adjacency matrix  $\hat{\mathcal{A}}$ . The framework also aims to minimize both the training cost  $C_{\text{train}}(f)$  and the inference cost  $C_{\text{infer}}(f)$ , ensuring efficiency during both model development and deployment.

Subject to the following constraints: - Every missing feature entry  $(i, d)$ , indicated by  $\mathbf{M}_{i,d} = 0$ , must be properly imputed in  $\hat{x}_{i,d}$ . - The missing rate of the feature matrix,  $\text{MissingRate}(\mathbf{M})$ , must not exceed the predefined threshold  $\rho_X$ . - The missing rate of the adjacency matrix,  $\text{MissingRate}(\mathcal{A})$ , must not exceed the threshold  $\rho_A$ .

Here,  $\mathbf{M}$  denotes the feature mask indicating missing entries, while  $\mathcal{A}$  and  $\hat{\mathcal{A}}$  represent the original and imputed adjacency matrices, respectively. The parameters  $\rho_X$  and  $\rho_A$  define the maximum tolerable missing rates for features and edges, guaranteeing robust model performance under severe data incompleteness. In our experiments, we set  $\rho_X = 99.5\%$  for

feature imputation in APCFI, and evaluate iPaD under edge missing rates up to  $\rho_A = 90\%$ , demonstrating the framework’s strong resilience to extremely high levels of missing data.

### 3.4 Summary

The central research question addressed in this work is as follows: **How can a modular and unified learning framework be designed to achieve robust prediction performance under severe data incompleteness, while simultaneously reducing training and inference costs for practical deployment?**

This framework tackles the joint objective by seamlessly integrating feature imputation, pruning, architectural tuning, and knowledge distillation into a single pipeline. Such integration not only enhances robustness and efficiency, but also provides a principled foundation for the subsequent methodological developments presented in this paper.

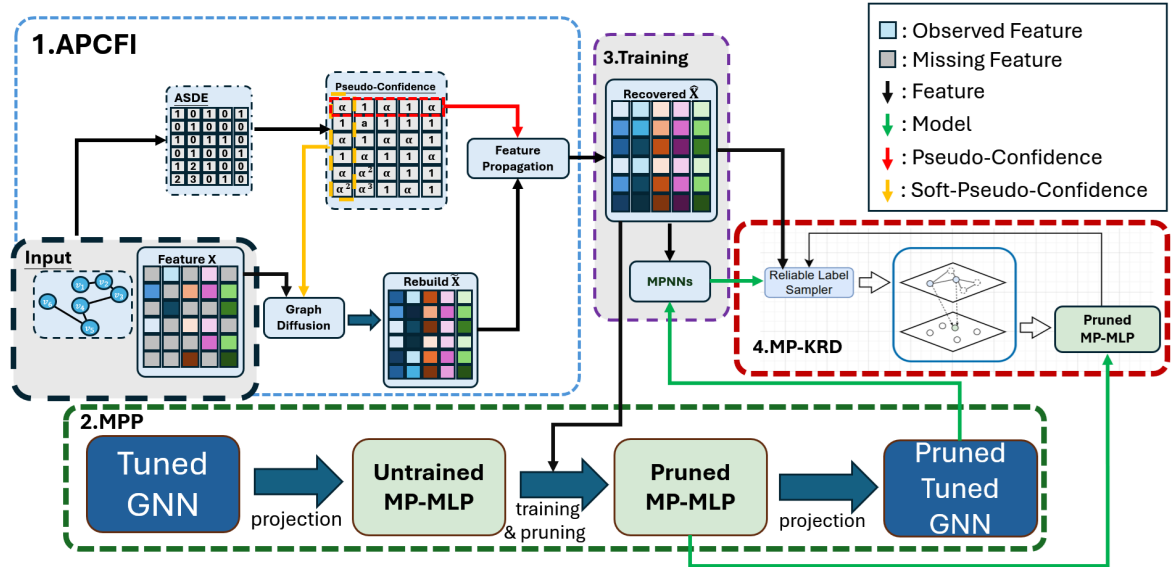
# Chapter 4

## Methodology

This chapter details the architecture and workflow of our proposed **iPaD** framework. As illustrated in Figure 4.1, iPaD is an end-to-end pipeline designed to address the challenges of incomplete data, high training costs, and inefficient inference in a unified and synergistic manner. The entire framework can be understood as a four-stage process:

- (1) Robust Data Imputation
- (2) Parallel Teacher-Student Preparation
- (3) Teacher Training
- (4) Architecture-Aligned Reliable Knowledge Distillation.

The subsequent sections will delve into the technical details of each core module.



**Figure 4.1: The overall architecture of the iPaD framework.** The process begins with imputing the incomplete feature matrix via APCFI. In parallel, a powerful GNN teacher is trained on the recovered data, while the MPP module generates a pruned, lightweight MLP student. Finally, the MP-KRD module distills knowledge from the teacher to the student, producing a final model that is both robust and efficient for inference.

### 4.1 Framework Overview

The workflow of our proposed iPaD framework, as depicted in Figure 4.1, is a systematic

pipeline that transforms a graph with incomplete features into a lightweight, robust, and fast inference model. The process unfolds across three distinct yet interconnected stages:

1. **Stage 1: Robust Data Imputation with APCFI.** The process begins with an input graph whose feature matrix,  $\mathbf{X}$ , contains missing values. To ensure a high-quality foundation for subsequent learning, the graph is first passed through our **APCFI** (Approximate Pseudo-Confidence Feature Imputation) module. APCFI leverages graph diffusion and an **ASDE** to generate pseudo-confidence scores, which guide a feature propagation process. The output of this stage is a fully recovered feature matrix,  $\hat{\mathbf{X}}$ , where missing entries have been robustly imputed.
2. **Stage 2: Synergistic Teacher and Student Generation via MPP.** The second stage efficiently prepares both the teacher and student models in a sequential, highly efficient process driven by our **MPP** module. The workflow is as follows:
  - (a) **Proxy Pre-training and Pruning:**  
We begin with our powerful GNN architecture (e.g., **TunedGNN**) and project it to its isomorphic MP-MLP. This proxy MLP is then pre-trained and pruned on the target dataset. This efficient step determines the optimal sparse structure and directly yields our lightweight **Student Model** (the “Pruned MP-MLP”).
  - (b) **Teacher Architecture Definition:**  
The pruned structure from the student model is then inversely projected back to the original GNN architecture. This creates the final, optimized architecture for our teacher: the “Pruned GNN”.

This synergistic process is a core innovation of our framework: the lightweight student’s architecture is defined in the same efficient step that determines the optimized architecture for the powerful teacher, all before the main, computationally expensive GNN training commences.

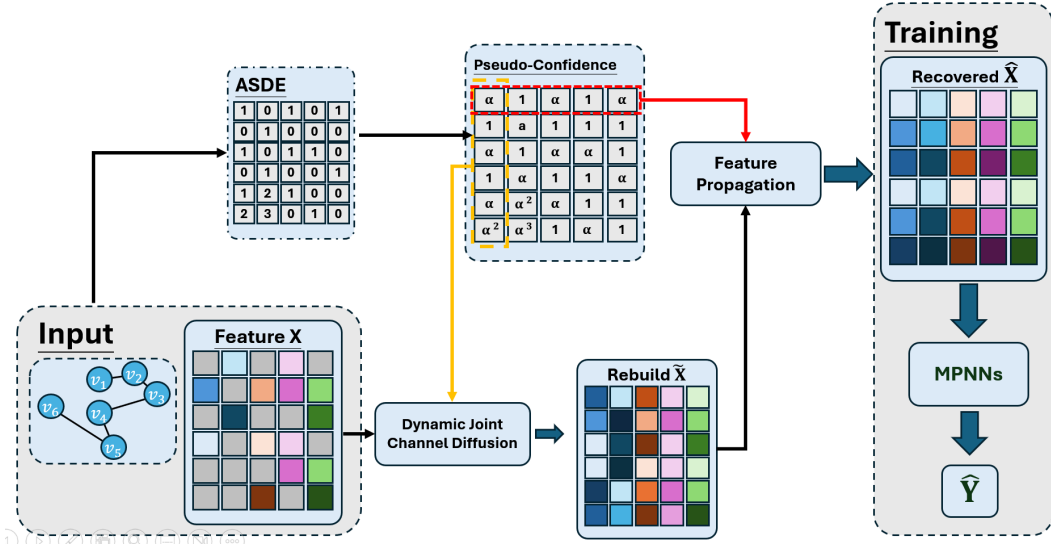
### 3. **Stage 3: Teacher Training.**

The “Pruned GNN” architecture, defined in the previous stage, is fully trained on the imputed dataset  $(\hat{\mathbf{X}}, \mathcal{G})$  produced by the APCFI module. This dedicated training phase enables the pruned model to achieve strong predictive performance while maintaining high computational efficiency. The resulting model, referred to as the **Teacher Model**, serves as a robust and lightweight source of knowledge for the subsequent distillation process. By leveraging both the optimized sparse architecture and the high-quality imputed features, the teacher model is well-suited to guide the training of the student model in the final stage.

4. **Stage 4: Knowledge Distillation with MP-KRD.** In the final stage, knowledge transfer is orchestrated. The **MP-KRD** module takes the high-performance “Pruned GNN” Teacher and the lightweight “Pruned MP-MLP” Student. It employs a reliability-aware sampling strategy to distill the rich, relational knowledge from the teacher into the student. The final output of the entire iPaD pipeline is this pruned, distilled MLP student—a model that is highly efficient, robust against incomplete data, and ready for fast deployment in real-world applications.

## 4.2 APCFI: Approximate Pseudo-Confidence Feature Imputation

As the first stage of the iPaD pipeline, the foundational challenge of data incompleteness is addressed. To this end, **APCFI** is proposed as a robust and scalable feature completion strategy, advancing beyond prior art by introducing a more sophisticated, two-stage normalization process to enhance imputation quality. Figure 4.2 illustrates the overall workflow of the APCFI module.



**Figure 4.2: The overall workflow of the APCFI module.** It takes a graph with an incomplete feature matrix as input and leverages ASDE and Dynamic Joint Channel Diffusion (DJCD) to output a robustly recovered feature matrix for downstream tasks.

### 4.2.1 Motivation and Core Idea

Previous high-performance imputation methods like PCFI [21] successfully utilize channel-wise confidence to guide feature diffusion. However, this approach faces two key limitations: the calculation of exact shortest paths is computationally expensive, and the raw confidence scores may not optimally reflect the relative importance of nodes.

APCFI is designed to overcome both limitations. Its core philosophy is to improve both the

efficiency of confidence calculation and the quality of the diffusion weights. This is achieved through three key innovations:

1. **Approximate Shortest Distance Estimation (ASDE):** Replaces PCFI’s costly k-hop traversal with a highly parallelizable message-passing scheme to efficiently approximate shortest-path distances.
2. **Soft-Pseudo Confidence (Soft-PC):** Introduces a “softmax” layer before the standard row-stochastic normalization to transform raw confidence scores into context-aware, relative reliability weights.
3. **Dynamic Joint Channel Diffusion:** Replaces PCFI’s sequential diffusion with a parallel process, reducing complexity from  $O(d \times k)$  to  $O(k)$ .

#### 4.2.2 Technical Workflow

The APCFI workflow consists of three main steps to transform an incomplete feature matrix  $\mathbf{X}$  into a recovered matrix  $\hat{\mathbf{X}}$ .

##### Step 1: Approximate Shortest Distance Estimation (ASDE)

To efficiently quantify the reliability of feature imputation, the raw pseudo-confidence (PC) is computed for each node-feature pair based on the shortest-path distance (SPD) from a missing-feature node to its nearest observed-feature node in the same channel. The ASDE algorithm, illustrated in Figure 4.3, iteratively propagates distance information through the graph, enabling rapid approximation of these SPD values.

Formally, the shortest-path distance for node  $v_i$  and feature channel  $d$  is defined as:

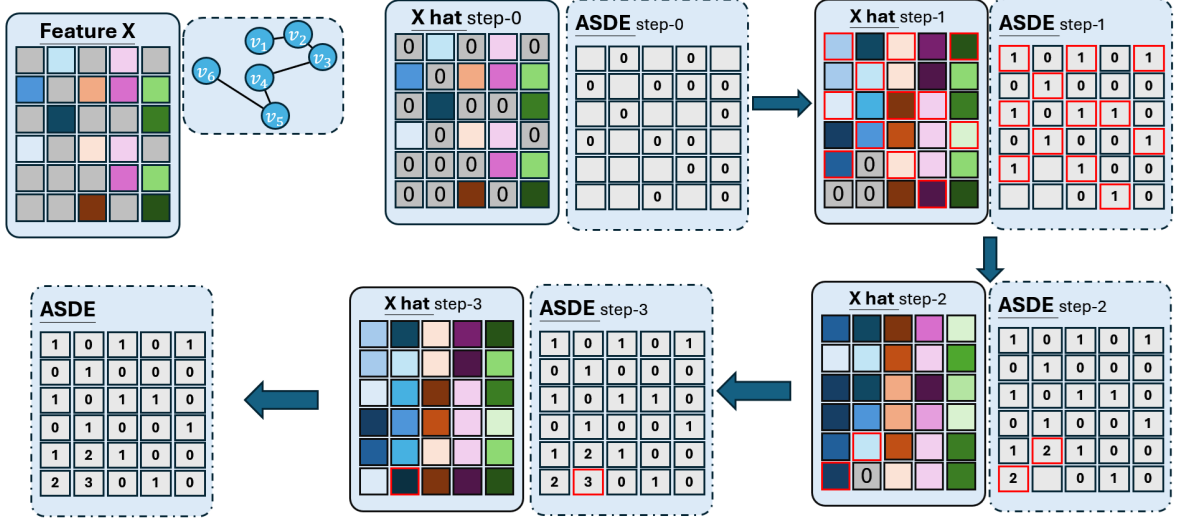
$$\text{SPD}_{i,S_d} = \begin{cases} 1, & \text{if } x_{i,d} \text{ is observed} \\ \min \left\{ t \mid \hat{x}_{i,d}^{(t)} \neq 0 \right\}, & \text{if } x_{i,d} \text{ is missing} \end{cases} \quad (4.1)$$

The raw pseudo-confidence is then given by:

$$pc_i^{(d)} = \alpha^{\text{SPD}_{i,S_d}} \quad (4.2)$$

where  $\alpha \in (0, 1)$  is a decay factor that controls the influence of distance on confidence.





**Figure 4.3: Illustration of the Approximate Shortest Distance Estimation (ASDE) process.** The algorithm begins with observed features (step-0) and propagates shortest-path information through the graph. Each entry in the ASDE matrix reflects the shortest hop-distance to an observed feature within the same channel.

## Step 2: Two-Stage Normalization

This step is the core of our imputation mechanism, involving a two-stage normalization process to generate the final diffusion weights.

### Stage 1: Soft-Pseudo Confidence

The raw pseudo-confidence scores  $pc_i^{(d)}$  are first passed through a “softmax” operation to produce the “soft confidence matrix”. The goal is to generate context-aware, relative reliability scores, which is a key innovation over PCFI. The resulting weights  $w_{i,d}$  are calculated as:

$$w_{i,d} = \frac{\exp(pc_i^{(d)})}{\sum_{j=1}^N \exp(pc_j^{(d)})} \quad (4.3)$$

To illustrate the effectiveness of this mechanism, consider the following example. Given a raw pseudo-confidence matrix **PC**:

$$\mathbf{PC} = \begin{bmatrix} 0.5 & 0.5 \\ 1.0 & 0.7 \\ 0.25 & 1.0 \end{bmatrix}$$

The resulting soft-pseudo confidence weights **soft\_PC** are:

$$\mathbf{soft\_PC} = \begin{bmatrix} 0.292 & 0.258 \\ 0.481 & 0.316 \\ 0.227 & 0.426 \end{bmatrix}$$

Notably, even though the first node has the same pseudo-confidence value (0.5) for both feature channels, its final soft-pseudo confidence weights differ (0.292 vs. 0.258). This is because softmax normalization takes into account the entire pseudo-confidence distribution for each channel, allowing identical values to have different relative importance depending on their context. This makes the weighting process more flexible and robust for the subsequent diffusion step.

### Stage 2: Row-Stochastic Transition

To ensure a stable and mathematically sound diffusion process, the “soft confidence matrix” (denoted as  $\mathbf{W}'$ ) undergoes a final row-stochastic normalization. This step is crucial for preserving the feature scale during propagation [21]. It is achieved by normalizing each weight by its corresponding row sum:

$$\mathbf{W} = \mathbf{D}^{-1}\mathbf{W}' \quad (4.4)$$

where  $\mathbf{D}$  is the diagonal degree matrix with  $D_{ij} = \sum_j W'_{ij}$ . The resulting matrix  $\mathbf{W}$  is row-stochastic, meaning each row sums to 1.

### Step 3: Dynamic Joint Channel Diffusion (DJCD)

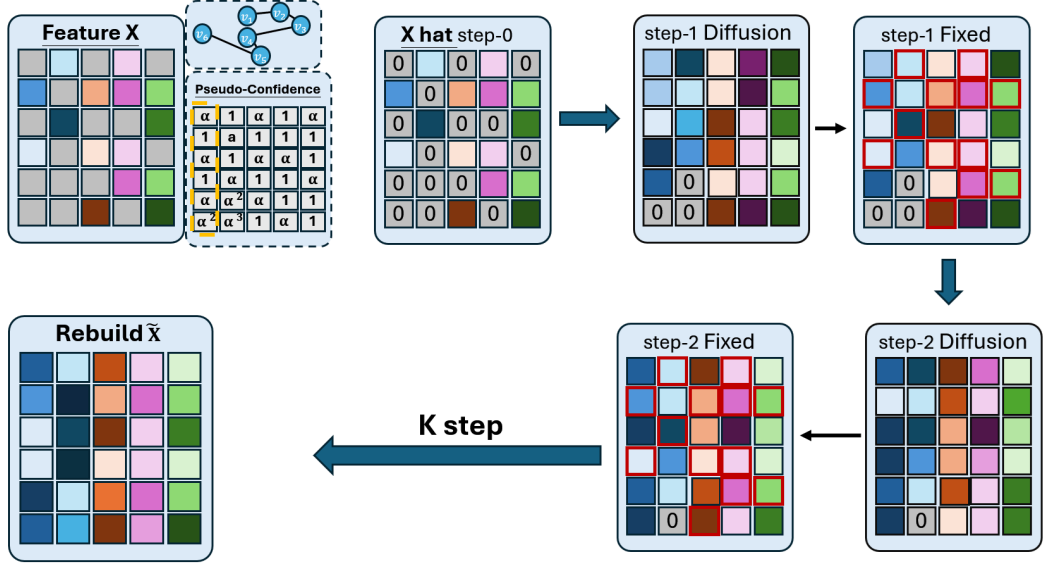
In this step, joint diffusion is performed based on the normalized weights, enabling efficient feature propagation across multiple channels.

After obtaining the normalized, row-stochastic weight matrix  $\mathbf{W}$ , the final step involves performing the actual feature diffusion. This iterative process, conducted for a total of  $K$  steps, propagates information from observed nodes to missing nodes across the graph.

The update rule for each node  $v_i$  and feature channel  $d$  at iteration  $t + 1$  is defined as:

$$\hat{x}_{i,d}^{(t+1)} = \begin{cases} x_{i,d}, & \text{if } (i, d) \text{ is observed} \\ \frac{\sum_{j \in \mathcal{N}(i)} \bar{w}_{ij} \cdot \hat{x}_{j,d}^{(t)}}{\sum_{j \in \mathcal{N}(i)} \bar{w}_{ij}}, & \text{otherwise} \end{cases} \quad (4.5)$$

For any feature that is already **observed**, its value is kept fixed. This is a crucial step to prevent the loss of ground-truth information and mitigate the “over-smoothing” problem, a



**Figure 4.4: The iterative process of Dynamic Joint Channel Diffusion.** In each of the  $K$  steps, a diffusion operation updates the missing values based on neighbors, followed by a “Fixed” step where the original, observed feature values are reset to prevent oversmoothing and preserve data fidelity.

technique also validated in the Feature Propagation (FP) [19] paper.

For any **missing** feature, its value is updated by computing a **weighted average** of the features of its neighboring nodes from the previous iteration. The weights  $\bar{w}_{ij}$  are precisely the ones we derived from our two-stage normalization process, ensuring that more reliable neighbors contribute more significantly to the imputation.

Notably, due to this design, the diffusion process can be conducted jointly and in parallel across all feature channels, which constitutes a key efficiency advantage of APCFI over PCFI. The entire iterative process is illustrated in Figure 4.4.

#### Step 4: Node-wise Inter-Channel Propagation

As a final refinement, following the primary feature diffusion in the DJCD step, a node-wise inter-channel propagation mechanism is introduced. This step is designed to leverage latent correlations between different feature channels, further fine-tuning the imputed values for each node. The core design is guided by the **Humility Principle**: nodes with lower confidence in a certain feature value should be more receptive to “suggestions” from other highly correlated features.

The process consists of the following computational steps. Let  $\hat{\mathbf{X}}_{\text{dcd}}$  denote the feature matrix recovered from the DJCD step, and let  $\mathbf{PC} \in \mathbb{R}^{N \times F}$  be the matrix of pseudo-confidence scores, where each entry  $pc_{i,d}$  represents the confidence in the  $d$ -th feature of node  $v_i$ .

1. **Inter-Channel Correlation:** The feature correlation matrix  $\mathbf{C} \in \mathbb{R}^{F \times F}$  is computed to capture the global linear relationships between all feature channels across the graph:

$$\mathbf{C} = \text{Corrcoef}(\tilde{\mathbf{X}}_{\text{djcd}}^T)$$

2. **Confidence-Weighted Signal:** For each node, an initial "suggestion" signal is generated by taking the deviation of each feature from its global mean, weighted by its pseudo-confidence score. This ensures that more reliable features contribute more strongly to the suggestion signal. Let  $\tilde{\mathbf{X}}_{\text{djcd}}$  be the feature mean vector:

$$\mathbf{M}_\alpha = \mathbf{PC} \odot (\tilde{\mathbf{X}}_{\text{djcd}} - \bar{\tilde{\mathbf{X}}}_{\text{djcd}})$$

where  $\odot$  denotes element-wise multiplication.

3. **Suggestion Aggregation:** The signal  $\mathbf{M}_\alpha$  is propagated through the correlation matrix  $\mathbf{C}$ , aggregating weighted suggestions from all other channels and enabling features to "advise" each other based on global correlation:

$$\mathbf{M}_{\text{conf}} = \mathbf{M}_\alpha \cdot \mathbf{C}$$

4. **Humility-based Update:** The update term  $\Delta\mathbf{X}$  is computed by weighting the aggregated suggestions  $\mathbf{M}_{\text{conf}}$  with a "humility" factor  $(1 - \mathbf{PC})$  and a hyperparameter  $\beta$ . The humility factor ensures that features with low initial confidence are updated more significantly:

$$\Delta\mathbf{X} = \beta \cdot (1 - \mathbf{PC}) \odot \mathbf{M}_{\text{conf}}$$

The final imputed feature matrix is obtained by applying this update to the result of the DJCD step:

$$\hat{\mathbf{X}} = \tilde{\mathbf{X}}_{\text{djcd}} + \Delta\mathbf{X}$$

This two-tiered process, combining broad graph-based diffusion with a fine-grained, correlation-aware adjustment, enables APCFI to produce a more robust and nuanced feature representation, providing a solid foundation for downstream GNN models.

#### 4.2.3 Implementation and Pseudocode

The complete APCFI workflow is detailed in the algorithms below (Algorithms 1 to 4).

In summary, by introducing a sophisticated two-stage normalization process and a highly efficient joint diffusion mechanism, APCFI provides a robust and scalable solution for data

incomplete. The recovered feature matrix,  $\hat{\mathbf{X}}$ , serves as a high-quality data foundation for subsequent model training. Specifically, the enhanced completeness and reliability of  $\hat{\mathbf{X}}$  enable the downstream GNN backbone to better exploit graph structure and node attributes, thereby improving overall learning stability and predictive performance, even in challenging real-world scenarios with severe missingness. The next section introduces the robust GNN backbone, which is designed to fully leverage this improved feature representation.

---

**Algorithm 1** PyTorch-style Pseudo-code of APCFI

---

```

1: function apcfi_fill(feature, edge_index)
2:     mask = feature.isnan()

3:     Step 1: ASDE and Dynamic Joint Channel Diffusion
4:     diffusion_x, asde_matrix = diffusion(feature, edge_index, mask)

5:     Step 2: Node-wise Inter-Channel Propagation
6:     output = node_propagation(diffusion_x, asde_matrix)

7:     Return the imputed feature matrix
8:     return output

```

---



---

**Algorithm 2** PyTorch-style Pseudo-code of Function diffusion

---

```

1: function diffusion(feature, edge_index, mask)

2:     Step 1: Calculate ASDE matrix
3:     asde_matrix = ASDE(edge_index, mask)

4:     Step 2: Compute pseudo-confidence
5:     confidence_matrix = pseudo_confidence(edge_index, asde_matrix)
6:     soft_confidence_matrix = F.softmax(confidence_matrix, dim=1)

7:     Step 3: Normalize by row transition
8:     normalize_confidence = row_transition(soft_confidence_matrix, edge_index[0])

9:     Step 4: Perform channel-wise diffusion
10:    diffusion_x = approximate_channel_diffusion(feature, normalize_confidence,
                                                edge_index, mask)

11:    Return: Channel-Diffusion(imputed) feature matrix
12:    return diffusion_x, asde_matrix

```

---

---

**Algorithm 3** PyTorch-style Pseudo-code of Function `row_transition`

---

```
1: function row_transition(soft_confidence_matrix, edge_index[0])

2:   Step 1: row-wise sum
3:   deg_w = torch_scatter.scatter_add(confidence_matrix, row_index, dim=0)

4:   Step 2: Inverse degree normalization
5:   deg_w_inv = deg_w.pow(-1.0)

6:   Step 3: Fill diagonal with zeros
7:   deg_w_inv.masked_fill(deg_w_inv == float(inf), 0)

8:   Step 4: Normalize confidence matrix
9:   norm_confidence_matrix = confidence_matrix * deg_w_inv[row_index]

10:  Return: Normalized confidence matrix
11:  return norm_confidence_matrix
```

---

---

**Algorithm 4** PyTorch-style Pseudo-code of Function `node_propagation`

---

```
1: function node_propagation(diffusion_x, asde_matrix)

2:   Step 1: Compute inter-channel correlation (conf)
3:   conf = torch.corrcoef(diffusion_x.T).nan_to_num().fill_diagonal_(0)

4:   Step 2: Compute channel mean difference
5:   alpha_mean = (alpha ** asde_matrix) * (diffusion_x - torch.mean(diffusion_x,
dim=0))

6:   Step 3: Aggregate with channel correlation
7:   alpha_mean_conf = torch.matmul(alpha_mean, conf)

8:   Step 4: Update node features
9:   update_x = beta * (1 - (alpha ** asde_matrix)) * alpha_mean_conf
10:  updated_x = diffusion_x + update_x

11:  Return: Node-updated feature matrix
12:  return updated_x
```

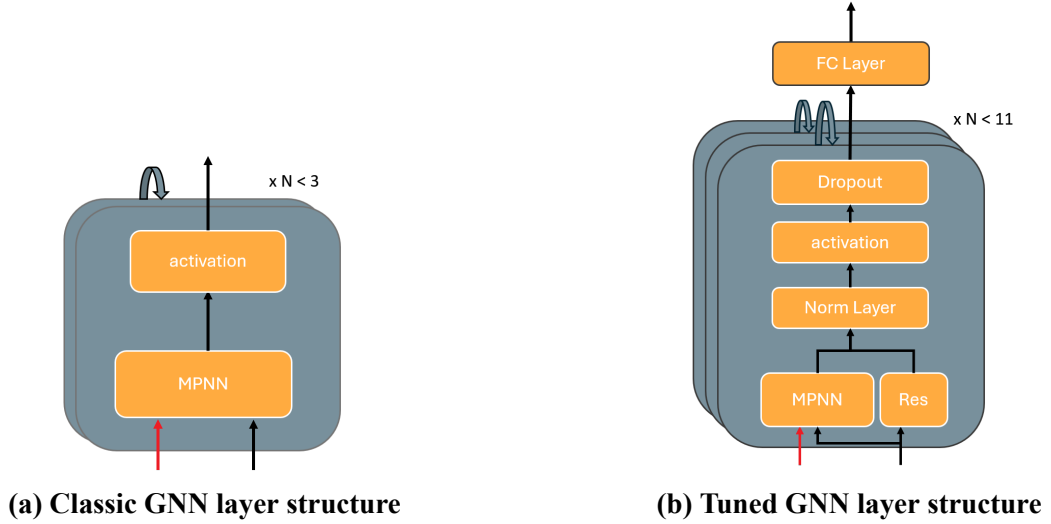
---

### 4.3 TunedGNN: Architecture Optimization for Robustness

Having established a method to obtain a complete and robust feature matrix, the next logical step is to define a GNN architecture capable of effectively learning from this data. While standard GNNs like GCN or GraphSAGE are widely used, they are often architecturally shallow and can be sensitive to the subtle noise and imperfections that may remain even in imputed data. To build a powerful and reliable “teacher” model for our framework, a more resilient foundational architecture is required.

To this end, the **TunedGNN** is introduced as the backbone model. It is important to note that **TunedGNN** is not a novel architecture proposed in this thesis. Rather, it is a systematic application of several well-established, state-of-the-art architectural best practices from recent GNN research [1]. The goal is to construct a model with enhanced stability, expressiveness, and, most critically, robustness against data perturbations. As illustrated in Figure 4.5b, the **TunedGNN** structure is significantly deeper and more complex than classic GNN layers.

#### 4.3.1 Key Design Components



**Figure 4.5: Comparison of Tuned and Classic GNN layer structures.** (a) The classic GNN layer stacks message-passing and activation modules, typically with fewer than 3 layers. (b) The tuned GNN layer incorporates residual connections, normalization, activation, and dropout, allowing for deeper architectures (up to 11 layers) and improved robustness to incomplete data.

The key design components systematically integrated into our **TunedGNN** backbone are as follows:

- **Deeper Layer Stacking:**

A deeper architecture (e.g., up to 11 layers in the experiments) is employed compared to traditional shallow GNNs. This allows for a larger receptive field, enabling the model to capture more complex, higher-order relational patterns within the graph, while carefully managing the risk of over-smoothing through the components below.

- **Normalization (e.g., LayerNorm):**

Each GNN layer is followed by a normalization layer, typically LayerNorm, to stabilize the distribution of intermediate feature representations. This practice is crucial for promoting faster, more stable convergence and is particularly effective in mitigating the impact of noise in high-missingness scenarios.

- **Residual Connections:**

To enable effective training of such a deep architecture, skip (or residual) connections are added between layers. As shown in Figure 4.5b, the output of the message-passing block is added to the original input before subsequent transformations. This mitigates the vanishing gradient problem and ensures a smoother gradient flow throughout the network.

- **Dropout:**

Dropout layers are applied after each non-linearity or normalization layer. This acts as a powerful regularization technique, enhancing the model’s generalization capabilities and improving its robustness to perturbations in both node and edge features.

The overall layer structure for **TunedGNN** can be formally summarized by the following sequence of operations for each layer, transforming the input representation  $h_i$  to the output  $h_{i+1}$ :

$$h'_i = \text{Aggregate}(h_i, \{h_j \mid j \in \mathcal{N}(i)\}) \quad (4.6)$$

$$h''_i = h'_i + \text{Linear}(h_i) \quad (4.7)$$

$$h'''_i = \text{Norm}(\sigma(h''_i)) \quad (4.8)$$

$$h_{i+1} = \text{Dropout}(h'''_i) \quad (4.9)$$

where Equation 4.6 represents the core message-passing step, and Equation 4.7 represents the residual connection.

In conclusion, the **TunedGNN** architecture provides a highly performant and robust foundation for our framework. However, this power and resilience come at the cost of a large parameter count and significant computational complexity. This directly motivates the need for the next module in our pipeline: an effective and efficient pruning strategy to make this



powerful model practical for real-world training and deployment.

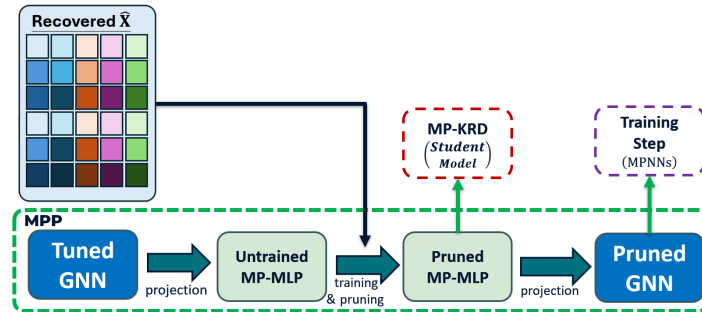
#### 4.4 MPP: Mirror Projection Pruning

The TunedGNN [1] architecture, as described in the previous section, provides a powerful and robust foundation for our framework. However, its depth and complexity lead to high computational costs during training. To address this challenge, we introduce **MPP**, a novel paradigm designed to efficiently reduce model complexity without sacrificing representational capacity.

Existing GNN pruning techniques are often tightly intertwined with the training process, resulting in significant computational overhead and a lack of flexibility. In contrast, MPP introduces a **proxy-based, pre-training pruning** strategy. Its core idea is to decouple the complex pruning task by projecting the GNN to an isomorphic MLP, applying mature pruning techniques in the simpler MLP space, and then mapping the pruned structure back. This design, inspired by studies like MLPinit [49] that demonstrate GNN-MLP parameter alignment, allows us to leverage a rich ecosystem of existing pruning tools, such as DepGraph (torch-pruning) [50]. Crucially, MPP serves a dual purpose in our pipeline: (1) it produces an efficient, pruned GNN to act as our **Teacher model**, and (2) it simultaneously generates a lightweight, structurally-aligned MLP to serve as our **Student model** for the final distillation stage.

##### 4.4.1 Mirror Projection Pruning (MPP): Theory and Workflow

Mirror Projection Pruning (MPP) (See Figure 4.6) is built on the structural isomorphism between common GNN layers and their MLP counterparts, extending the insight of MLPInit [49] that an MLP with the same architecture as a GNN can serve as a parameter proxy.

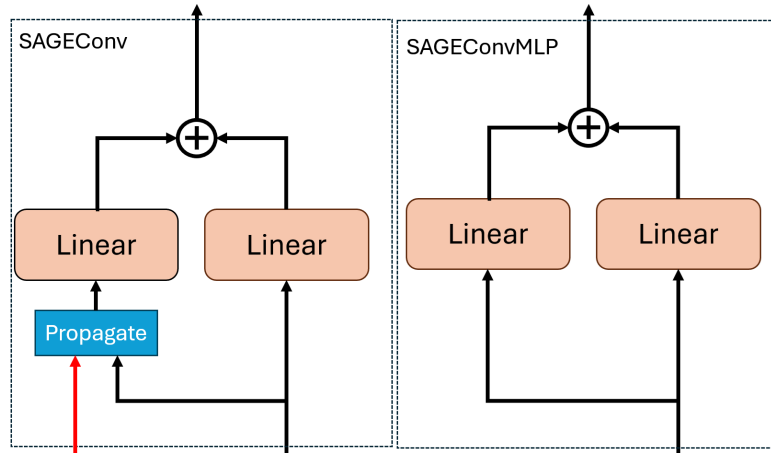


**Figure 4.6: The high-level workflow of Mirror Projection Pruning (MPP).** A large GNN architecture is projected to a proxy MP-MLP, which is then trained and pruned. The resulting sparse structure is projected back to define the final, efficient ‘Pruned GNN’.

Our approach generalizes this concept: an isomorphic MLP can act as an efficient surrogate for a GNN, enabling all training, pruning, and architectural search to be performed in the lightweight MLP domain. Afterward, the resulting sparse structure and weights are seamlessly mapped back to the original GNN, yielding an efficient, pruned GNN with minimal engineering effort.

### The Principle of GNN-MLP Isomorphism

The isomorphism between GNN and MLP layers is achieved by removing the neighborhood aggregation (message-passing) operation from the GNN, resulting in a pure node-level transformation (See Figure 4.7).



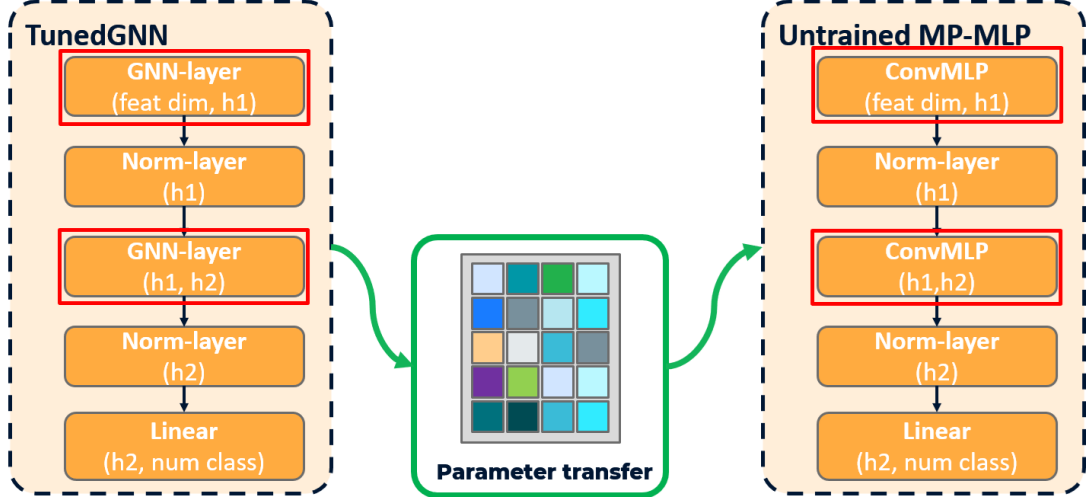
**Figure 4.7: Structural illustration of SAGEConv and its MP-MLP projection.** Left: The SAGEConv layer consists of two linear branches, one operating on neighborhood-aggregated features (Propagate). Right: The SAGEConvMLP projection replaces neighborhood aggregation with a second self-branch, yielding a pure local MLP structure with identical parameterization. This one-to-one mapping is fundamental to efficient pruning and parameter transfer in MPP.

For example, a SAGEConv layer comprises two linear branches—one operating on the node’s own features, the other on the aggregated features from its neighbors. By eliminating the propagation step, we obtain an MLP layer (SAGEConvMLP) with two self-branches, preserving both parameterization and output dimension.

This theoretical foundation ensures that all layers, hidden dimensions, and normalization structures in the original GNN are exactly matched by the MP-MLP, making direct parameter and mask transfer possible in both directions.

### Step 1: Mirror Projection—GNN to MP-MLP

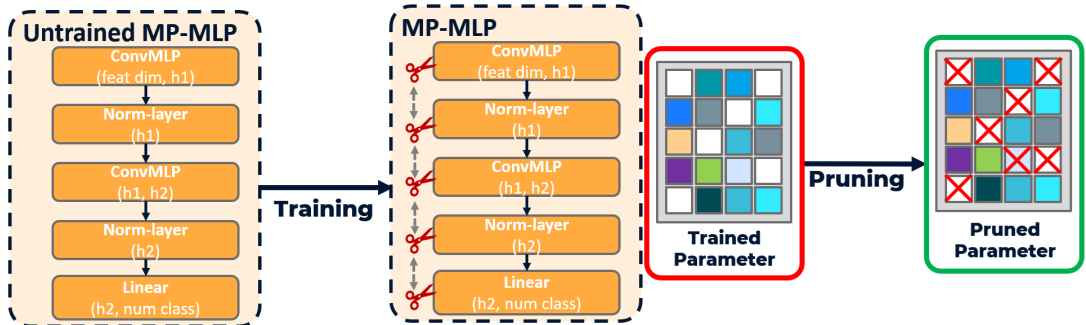
Given a TunedGNN backbone (e.g., TunedGNN [1]), the first step is to construct a mirror-projected MLP (MP-MLP) by replacing each GNN-layer with its isomorphic ConvMLP layer. All message-passing operations are removed, while linear transformations, normalization, and nonlinearity are preserved. This results in an untrained MP-MLP network that is structurally identical to the original GNN, as shown in Figure 4.8.



**Figure 4.8: Mirror Projection from GNN to MP-MLP.** Each GNN-layer is replaced by a corresponding ConvMLP layer, yielding an untrained MP-MLP that mirrors the GNN architecture.

### Step 2: Training and Pruning in MP-MLP Space

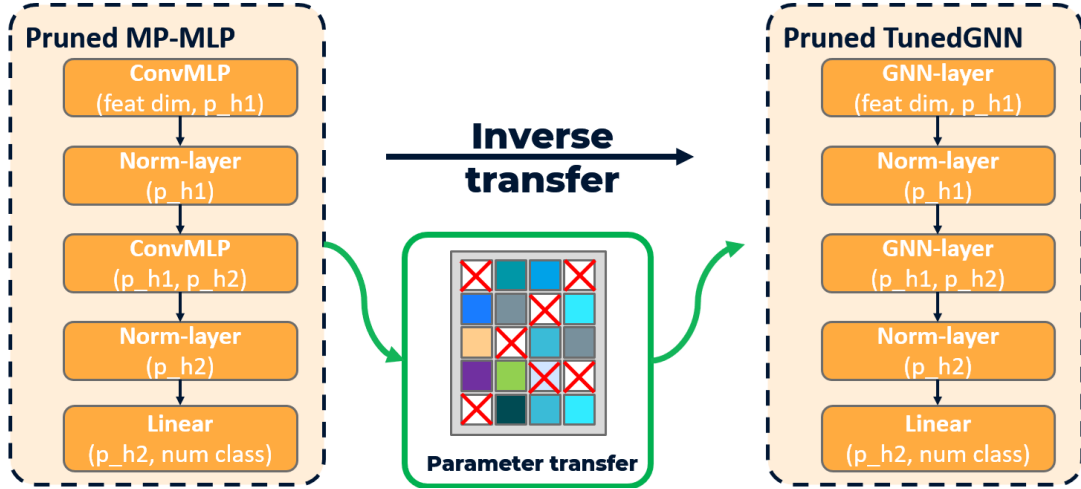
All subsequent training and pruning are performed in the MP-MLP domain. The MP-MLP is first trained on the completed feature set for the target task, then pruned using state-of-the-art MLP pruning algorithms (e.g., LAMP [51], DepGraph [50]), which efficiently remove redundant weights or channels based on importance. This step is illustrated in Figure 4.9.



**Figure 4.9: Proxy Training and Pruning in MP-MLP.** The MP-MLP is trained and pruned to yield a sparse, efficient student model.

### Step 3: Inverse Projection—Parameter and Mask Transfer Back to GNN

Once pruning is completed, the sparse weights and binary pruning mask from the MP-MLP are mapped directly back to the original GNN, thanks to the strict structural alignment. Each pruned ConvMLP layer's parameters and mask are copied to the corresponding GNN-layer, transforming the large, dense GNN into a compact, sparse network without modifying the original codebase. See Figure 4.10.



**Figure 4.10: Inverse Projection from Pruned MP-MLP to Pruned GNN.** The pruning mask and weights are transferred directly, instantly producing a sparse, efficient GNN.

#### Advantages of MPP

- **Computational efficiency:** Training and pruning are much faster in the MLP domain, as no neighborhood aggregation is needed.
- **Structural flexibility:** Any MLP pruning or search technique can be applied to the proxy MP-MLP and directly benefit the GNN.
- **Seamless transfer:** Strict one-to-one mapping enables direct synchronization of parameters and masks between MLP and GNN.
- **Minimal engineering overhead:** No modification to the GNN implementation is required; all changes are applied through parameter transfer.

This mirror projection–pruning–inverse projection workflow offers a highly flexible, efficient, and practical GNN compression paradigm, making advanced pruning and architectural optimization accessible for real-world graph learning tasks.

## Implementation and Pseudocode

The entire process is outlined in Algorithm 5. This workflow efficiently produces a pruned teacher GNN and a structurally-aligned student MLP, perfectly setting the stage for our final knowledge transfer module.

---

**Algorithm 5** MPP Proxy Pruning Pseudocode

---

```
1: function mpp_prune(gnn_model, feature, label)

2:   Step 1: Remove message-passing and project to MLP proxy
3:   mp_mlp = project_to_mlp(gnn_model)

4:   Step 2: Pre-train the proxy model
5:   mp_mlp.train(feature, label)

6:   Step 3: Prune the proxy MLP
7:   pruning_mask = prune_mlp(mp_mlp, method="LAMP")

8:   Step 4: Generate the student model
9:   pruned_mlp = apply_mask(mp_mlp, pruning_mask)

10:  Step 5: Parameter and Mask Transfer back to original GNN (teacher model)
11:  pruned_gnn = apply_mask(gnn_model, pruning_mask)

12:  Return the pruned GNN and pruned MLP
13:  return pruned_gnn, pruned_mlp
```

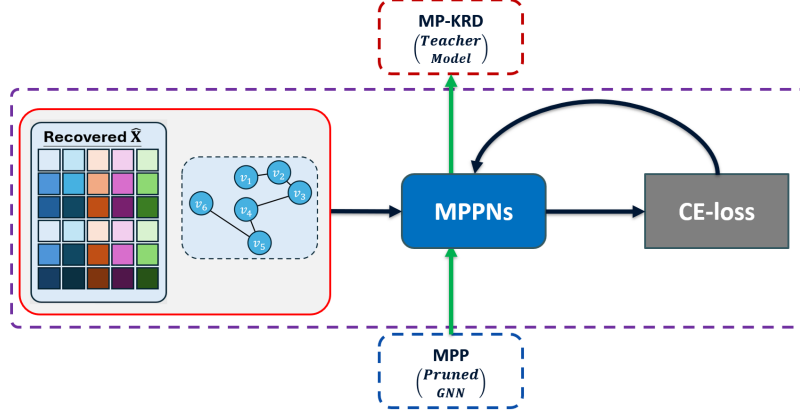
---

### 4.5 Teacher Training

After obtaining the pruned GNN from the MPP workflow, a dedicated training phase is conducted to optimize its parameters on the imputed dataset generated by the APCFI module. Specifically, the pruned GNN is trained using the recovered feature matrix  $\hat{\mathbf{X}}$  and the original graph structure, following standard supervised learning protocols for node classification.

Throughout this stage (see Fig. 4.11), the model benefits from both the lightweight architecture achieved via pruning and the improved data quality resulting from feature imputation. The cross-entropy loss is adopted as the objective function to maximize classification accuracy. Upon completion of training, the resulting high-performance and efficient GNN is designated as the **Teacher Model** for the subsequent knowledge distillation phase.

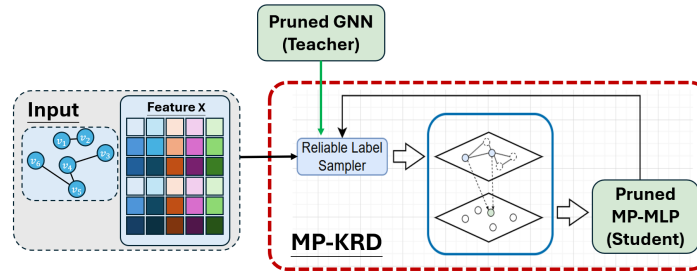
This process ensures that the teacher model not only possesses strong predictive capabilities but also maintains computational efficiency, thereby serving as an ideal source of knowledge for transferring to a student model.



**Figure 4.11: Teacher Training Workflow.** The pruned GNN obtained from the MPP module is fully trained on the imputed feature matrix  $\hat{\mathbf{X}}$  generated by APCFI. The resulting efficient model is used as the Teacher for knowledge distillation.

#### 4.6 MP-KRD: Mirror Projection Knowledge-inspired Reliable Distillation

The final stage of our iPaD pipeline, **MP-KRD**, orchestrates the knowledge transfer from the powerful teacher model to the lightweight student model, both of which were prepared by the preceding MPP module. Specifically, it takes the trained “Pruned TunedGNN” as the Teacher and the “Pruned MP-MLP” as the Student. As depicted in Figure 4.12, the goal is to create a final inference model that is both highly efficient and robust.



**Figure 4.12: The high-level view of the MP-KRD module.** It takes the trained GNN Teacher and the Pruned MP-MLP Student as input. A Reliable Label Sampler, guided by a curriculum, selects high-quality knowledge to distill, producing the final, efficient inference model.

#### 4.6.1 Motivation: The Need for Reliable Distillation

A naive distillation process that simply forces the student to mimic all of the teacher’s outputs can be problematic. A teacher trained on imputed or noisy data may exhibit low confidence or even make mistakes on certain samples. Forcing the student to learn this “unreliable knowledge” can harm its performance. To motivate our approach, consider two samples (A and B) with the same ground-truth label  $[0, 1, 0]$ , where the teacher GNN produces the following output logits:

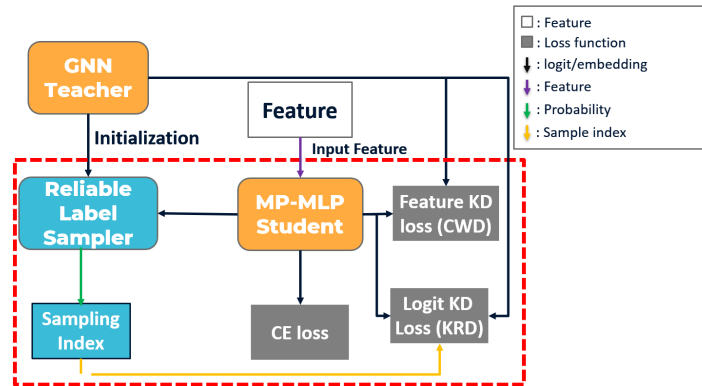
- **Sample A:** GNN output logits  $[0.3, 0.5, 0.2]$
- **Sample B:** GNN output logits  $[0.075, 0.9, 0.025]$

Although both samples are correctly classified, the cross-entropy loss for Sample B is noticeably lower than for Sample A. Intuitively, this means the model is more confident in its prediction for Sample B. Prediction reliability is systematically quantified using a robust mechanism rather than relying solely on raw logit values.

Therefore, instead of treating all teacher knowledge equally, MP-KRD employs a **reliability-aware curriculum**. The core idea is to quantify sample “difficulty” or “unreliability” by the **stability of the teacher’s prediction under noise perturbation**, measured via entropy change. This allows the student to learn from the most stable and reliable knowledge first.

#### 4.6.2 The MP-KRD Workflow

The MP-KRD process is an iterative workflow that adaptively transfers knowledge from the teacher to the student (see Fig. 4.13).



**Figure 4.13: The workflow of the MP-KRD module.** The pruned GNN teacher provides sample reliability signals for curriculum-based distillation, guiding the training of the MP-MLP student via cross-entropy, logit KD (KRD), and feature KD (CWD) losses.

### Step 1: Sample Difficulty Calculation.

The workflow begins by quantifying the reliability of the teacher’s prediction for each training sample  $i$ . This is done by measuring the change in the teacher’s output entropy when a small Gaussian noise,  $\epsilon \sim \mathcal{N}(0, \delta)$ , is added to the input features  $\mathbf{X}$ .

$$\Delta H_i = |H(f_{\text{GNN}}(\mathbf{X}', i)) - H(f_{\text{GNN}}(\mathbf{X}, i))| \quad (4.10)$$

$$\Delta e_i = \frac{\Delta H_i - \min_j \Delta H_j}{\max_j \Delta H_j - \min_j \Delta H_j} \quad (4.11)$$

The normalized entropy change,  $\Delta e_i \in [0, 1]$ , serves as our metric for sample difficulty, where a higher value indicates a less reliable or more difficult sample.

### Step 2: Curriculum-based Sampling.

In each training epoch, we use a curriculum to select a subset of reliable samples. The probability of selecting sample  $i$  is governed by a dynamic power function:

$$p_i = 1 - \Delta e_i^{\text{power}}$$

Based on this probability, a mask is generated via Bernoulli sampling, and only the selected samples ( $S_{\text{epoch}}$ ) are used for training in that epoch. This ensures that easier samples (low  $\Delta e_i$ ) are prioritized in the early stages of training.

### Step 3: Composite Distillation Objective.

The student model is trained on the sampled subset  $S_{\text{epoch}}$  using a composite loss function designed to capture knowledge from multiple perspectives:

$$L_{\text{MP-KRD}} = L_{\text{CE}} + \lambda L_{\text{KRD}} + \alpha L_{\text{CWD}} \quad (4.12)$$

The hyperparameters  $\lambda$  and  $\alpha$  balance the contributions of these different forms of knowledge.

- **Cross-entropy loss** ( $L_{\text{CE}}$ ): The standard supervised cross-entropy loss against the ground-truth labels,

$$L_{\text{CE}} = - \sum_{i \in N} y_i \log(\hat{y}_i^S) \quad (4.13)$$

where  $y_i$  is the true label and  $\hat{y}_i^S$  denotes the student model’s softmax output.



- **Logit-based distillation loss** ( $L_{\text{KRD}}$ ): Encourages the student's output distribution to match the teacher's,

$$L_{\text{KRD}} = \sum_{i \in S_{\text{epoch}}} \text{KL}(\text{Student}_{\text{logit}}(x_i), \text{Teacher}_{\text{logit}}(x_i)) \quad (4.14)$$

where KL denotes the Kullback-Leibler divergence, and  $S_{\text{epoch}}$  is the set of training samples selected for distillation at the current epoch (e.g., the sampled reliable points under the curriculum policy).

- **Feature-based distillation loss** ( $L_{\text{CWD}}$ ): Pushes the student's hidden representations ( $z_i^{\text{student}}$ ) to align with the teacher's ( $z_i^{\text{teacher}}$ ) [52],

$$L_{\text{node-wise}}^l = \text{KL}(\text{Softmax}(\mathbf{Z}^{S,l}), \text{Softmax}(\mathbf{Z}^{T,l})) \quad (4.15)$$

$$L_{\text{channel-wise}}^l = \text{KL}(\text{Softmax}((\mathbf{Z}^{S,l})^\top), \text{Softmax}((\mathbf{Z}^{T,l})^\top)) \quad (4.16)$$

$$L_{\text{CWD}}^l = L_{\text{node-wise}}^l + L_{\text{channel-wise}}^l \quad (4.17)$$

$$L_{\text{CWD}}^{\text{total}} = \sum_{l \in L} L_{\text{CWD}}^l \quad (4.18)$$

where  $l$  indexes the selected layers for distillation, and  $L$  is the set of such layers.

#### Step 4: Adaptive Curriculum Update.

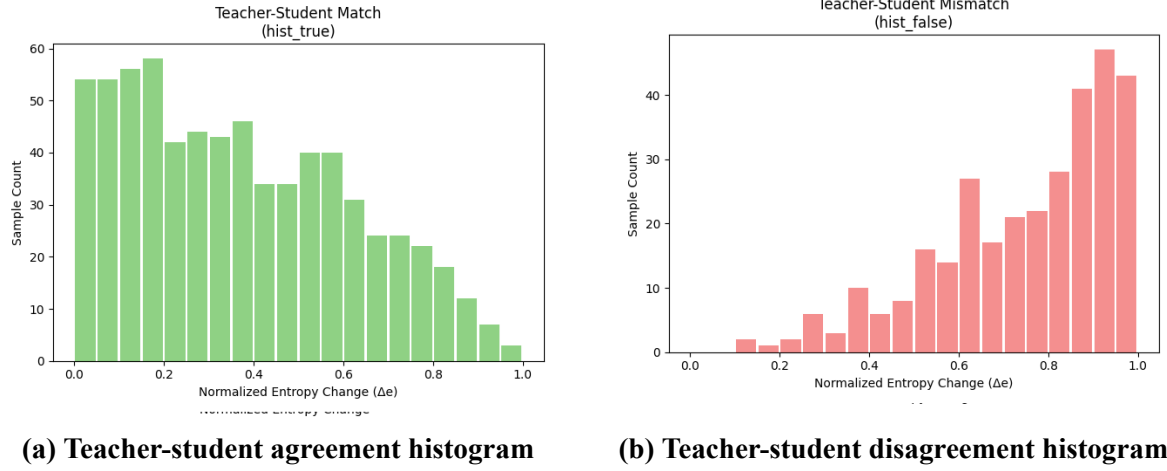
To make the curriculum dynamic, the ‘power’ parameter is updated at the end of each epoch. This is done by analyzing the student's current performance against the teacher's.

In order to investigate how sample difficulty affects the consistency between teacher and student predictions, all training samples are partitioned into discrete bins based on their normalized entropy change ( $\Delta e$ ). For each bin, the frequencies of teacher-student matches and mismatches are computed, and their distributions are illustrated in Figure 4.14. This provides an empirical foundation for the subsequent curriculum adjustment process.

To quantitatively measure the agreement probability within each entropy bin, the following metric is defined:

$$\text{hist} = \frac{\text{hist}_{\text{true}}}{\text{hist}_{\text{true}} + \text{hist}_{\text{false}} + \varepsilon} \quad (4.19)$$

where  $\text{hist}_{\text{true}}$  and  $\text{hist}_{\text{false}}$  denote the number of agreement and disagreement samples in  $\varepsilon$  is a small constant for numerical stability.



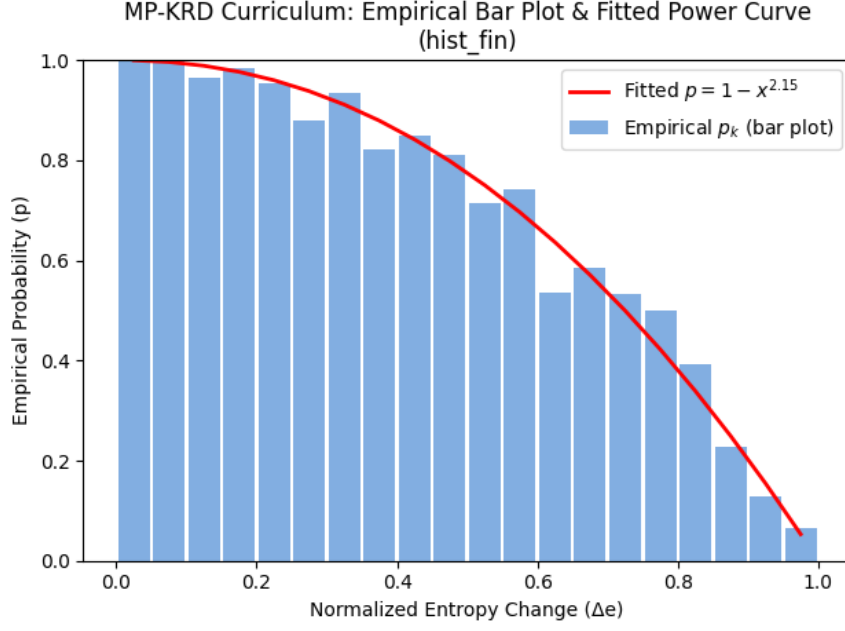
**Figure 4.14: Histograms of teacher-student agreement (a) and disagreement (b) plotted against the normalized entropy change.** Easier samples (low entropy change) show high agreement, while harder samples show high disagreement.

Next, as shown in Figure 4.15, we fit the empirical probabilities  $\text{hist}_{fin}$  using the function:

$$p = 1 - \Delta e^{\text{power}_{fit}} \quad (4.20)$$

$$\text{power}_{fit} = \frac{\log(1 - p)}{\log(\Delta e)} \quad (4.21)$$

where  $\Delta e$  is the normalized entropy change and  $\text{power}_{fit}$  is determined by least-squares fitting. This process dynamically determines the optimal curriculum power parameter for the current epoch.



**Figure 4.15:** The empirical agreement probability ( $p_k$ ) is fitted with the power curve (red line) to dynamically determine the optimal curriculum “ $\text{power}_{fit}$ ” parameter.

Finally, the “power” for the next epoch is updated using an **Exponential Moving Average (EMA)** to ensure a smooth curriculum progression:

$$\text{power}(t) = \beta \cdot \text{power}(t - 1) + (1 - \beta) \cdot \text{power}_{fit}(t) \quad (4.22)$$

where  $\beta$  is a smoothing coefficient(momentum) and  $\text{power}_{fit}(t)$  is the fitted value from the previous step.

### 4.6.3 Deployment

After MP-KRD training, only the sparse pruned mirror-projected MLP student is deployed, achieving both high efficiency and robustness through distilled GNN knowledge.

## 4.7 Innovation Summary

Tables 4.1–4.3 provide a component-wise classification of the techniques integrated into the framework, categorized into **adoption**, **improvement**, or **original innovation**. This classification specifies the methodological sources and identifies the components that constitute original contributions.

Within the **Imputation Module (APCFI)**, the **Parallel Diffusion** mechanism is classified as an **original innovation**, enabling simultaneous multi-channel feature estimation to enhance computational stability and scalability on large-scale graphs. In the **Pruning Module (MPP)**, the **Proxy Pruning** strategy, also categorized as an **original innovation**, employs a low-cost MP-MLP proxy during pruning to accelerate both pruning and GNN training, while improving the effectiveness of knowledge distillation. For the **Distillation Module (MP-KRD)**, the **Homogeneous Student Structure** represents an **original innovation** that resolves structural mismatches between teacher and student models, and the modified **CWD Loss** improves the efficiency and stability of channel-wise knowledge transfer.

Technical Component	Source	Innovation Description	Innovation Category
Feature-Propagation	[19]	Adopted standard feature propagation as the fundamental step for missing feature estimation	Adoption
Pseudo-Confidence	[21]	Adopted existing pseudo-confidence calculation method as the basis for feature reliability estimation	Adoption
Parallel Diffusion	<b>Our</b>	Designed a parallel diffusion strategy to perform feature estimation simultaneously on multiple subgraph channels, improving computational efficiency and stability for large-scale graphs	Original Innovation

**Table 4.1: Imputation Module (APCFI) Innovations.** Detailed breakdown of the APCFI components, specifying the source of each technique and categorizing the level of innovation.

Technical Component	Source	Innovation Description	Innovation Category
Parameter Transfer	[49]	Redesigned parameter transfer to adapt projection and weight mapping for GNNs, improving stability after pruning	Adoption
Channel Pruning	[50, 51]	Adopted conventional channel pruning strategies without modifying the original computation and selection criteria	Adoption
Proxy Pruning	<b>Our</b>	Proposed replacing GNN with a low-cost MLP as a proxy model during pruning. The proxy significantly accelerates pruning and GNN training, while enhancing the effectiveness of knowledge distillation by providing a lightweight yet structure-aware student model	Original Innovation

**Table 4.2: Pruning Module (MPP) Innovations.** Comprehensive summary of pruning techniques, including parameter transfer, channel pruning, and the proposed proxy pruning strategy using MLP to accelerate GNN training and enhance KD.

Technical Component	Source	Innovation Description	Innovation Category
<b>KDMLP</b>	[29]	Adapted existing MLP-based distillation methods to fit GNN architectures, applied to student models generated by MPP pruning	Adoption
<b>Curriculum Learning</b>	[30]	Incorporated curriculum learning strategies into the distillation process, gradually increasing sample difficulty to enhance student model learning capability	Adoption
<b>Homogeneous Student Structure</b>	<b>Our</b>	Proposed mp-MLP-based student model design to align its structure with the teacher's embedding space, addressing teacher-student structure mismatch in distillation	Original Innovation
<b>Channel-wise Distillation Loss</b>	[52]	Modified the existing CWD Loss to adapt weighting and regularization for GNN feature channels, improving the efficiency and stability of knowledge transfer	Improvement

**Table 4.3: Distillation Module (MP-KRD) Innovations.** Overview of MP-KRD components, covering KDMLP, curriculum learning, the proposed homogeneous student structure, and the improved CWDLoss for GNN adaptation.

## Chapter 5

### Experiment

#### 5.1 Experimental Setup

To ensure a thorough, fair, and reproducible evaluation of the proposed iPaD framework, this section details the datasets, compared methods, training configurations, and evaluation protocols used throughout this chapter.

##### 5.1.1 Datasets

Experiments on eight widely-used node classification benchmark dataset are conducted. These include three classic citation networks: **Cora**, **CiteSeer**, and **PubMed** [53]. Additionally, five medium-to-large-scale graphs are utilized: the **Computers** and **Photo** co-purchase networks, the **CS** and **Physics** co-authorship networks [54], and the **WikiCS** dataset [55]. Key statistics for each dataset are summarized in Table 5.1. All tasks are formulated as semi-supervised node classification, with accuracy as the primary evaluation metric.

Datasets	nodes	edges	Features	Classes	Metric
Cora [53]	2,708	5,278	1,433	7	Accuracy
CiteSeer [53]	3,327	4,522	3,703	6	Accuracy
PubMed [53]	19,717	44,324	500	3	Accuracy
Computers	13,752	245,861	767	1	Accuracy
Photo	7,650	119,081	745	8	Accuracy
WikiCS [55]	11,701	216,123	300	6	Accuracy
CS	18,333	81,894	6,805	15	Accuracy
Physics	34,493	247,962	8,415	5	Accuracy

**Table 5.1: Statistics of Benchmark Datasets Used in Node Classification Tasks.** The table summarizes the number of nodes, edges, features, classes, and the evaluation metric for each dataset.

### 5.1.2 Compared Methods

Experiments involve a diverse range of baseline and state-of-the-art models. For imputation-specific analysis (Section II), The proposed **APCFI** is compared against **Zero**, **LP** [22], **sRMGCNN** [56], **GCNMF** [20], **PaGNN** [12], **FP** [19], and **PCFI** [21]. For end-to-end model performance (Section III), **iPaD** is benchmarked against a strong **TunedGNN** [1], baseline and recent state-of-the-art models, including **GraphGPS** [32], **NAGphormer** [33], **Exphormer** [34], **GOAT** [35], **NodeFormer** [36], **SGFormer** [37], and **Polynormer** [38]. Additionally, for knowledge distillation, we directly compare our method with state-of-the-art KD baselines, including **GLNN**[29] and **KRD**[30], to demonstrate the effectiveness of our approach under distillation scenarios.

### 5.1.3 Training Details

All models are trained using the Adam optimizer. To ensure a fair and robust comparison, critical hyperparameters, including the learning rate, are systematically tuned for each model and dataset using the **Optuna** [57]. Unless otherwise specified, the number of training epochs follows the settings from the TunedGNN benchmark [1].

To account for stochasticity, **all reported experimental results are the average of 10 independent runs, each with a different but fixed random seed**. All experiments are conducted on a workstation equipped with an NVIDIA RTX 4090 GPU, an AMD Ryzen 9 7900 CPU, and 64 GB of DDR5 RAM.

### 5.1.4 Missing Data Settings

To evaluate robustness, robustness is evaluated by systematically introducing missing data under two primary settings: **feature missing** and **edge missing**. Feature missing includes both uniform (random) and structural (patterned) masking, with missing rates set to **0.0**, **0.3**, **0.6**, **0.9**, and the extreme **0.995**. Edge missing uses random masking with rates of **0.0**, **0.3**, **0.6**, and **0.9**.

### 5.1.5 Dataset Split Strategy

To ensure fair comparisons with prior work, the specific data split protocols established in the benchmark literature are adopted to ensure fair comparisons with prior work. The experiments use two distinct splitting strategies based on the evaluation’s focus:



## Imputation Performance Analysis (Section II).

The protocol from the PCFI paper [21] is strictly followed. Specifically, for each dataset, 10 different random splits are generated. In each split, 20 nodes per class are allocated for training ( $N_{\text{train}}$ ). The validation set size is then set to  $1500 - N_{\text{train}}$ . All remaining nodes are used for testing.

## All Other Experiments (SOTA, Robustness, Ablations).

For all other experiments (Sections III to VII), the widely-accepted splits from the TunedGNN benchmark [1] are used. This involves several settings:

- For **Cora, CiteSeer, and PubMed**, we use the standard semi-supervised splits as defined in [53].
- For **Computers, Photo, CS, and Physics**, we use a standard 60%/20%/20% split for training, validation, and testing, respectively.
- For **WikiCS**, we use the official public split provided by its original authors [55].

Adhering to these established protocols ensures the reliability and comparability of our results with the broader academic community.

## 5.2 Validating APCFI: A Superior Balance of Accuracy and Efficiency

This section validates our first claim: that **APCFI** effectively resolves the performance-efficiency trade-off present in existing feature imputation methods. Specifically, we evaluate whether APCFI achieves greater efficiency than its high-performance predecessor, PCFI, while maintaining comparable imputation accuracy. The analysis focuses on verifying that the observed efficiency improvements do not compromise the quality of feature imputation.

### 5.2.1 Efficiency Comparison

As shown in Table 5.2, APCFI achieves significant improvements in efficiency. Across all benchmark datasets, APCFI is only marginally slower than the simple FP baseline. Crucially, compared to the high-performance PCFI, APCFI demonstrates a massive efficiency improvement, running on average **over 200 times faster**. This result confirms that APCFI is a practically viable imputation solution for large-scale graphs where PCFI’s complexity would be prohibitive.

Method	Cora	CiteSeer	PubMed	Photo	Computers	Runtime(s)
FP	$0.05 \pm 0.01$	$0.06 \pm 0.00$	$0.06 \pm 0.00$	$0.06 \pm 0.00$	$0.09 \pm 0.00$	0.06
PCFI	$73.40 \pm 0.42$	$252.8 \pm 2.29$	$31.18 \pm 0.43$	$32.45 \pm 0.16$	$39.99 \pm 0.17$	85.97
<b>APCFI(our)</b>	<b><math>0.19 \pm 0.04</math></b>	<b><math>0.39 \pm 0.02</math></b>	<b><math>0.19 \pm 0.01</math></b>	<b><math>0.46 \pm 0.00</math></b>	<b><math>0.72 \pm 0.00</math></b>	<b>0.39</b>

**Table 5.2: Runtime comparison of different feature imputation methods across datasets.** The table reports the runtime (in seconds) for FP, PCFI, and APCFI methods on several benchmark datasets. APCFI achieves comparable efficiency to FP and significantly outperforms PCFI in computational cost.

### 5.2.2 Accuracy under High Missing Rates

Table 5.3 addresses the second question, showing that this immense efficiency gain does not compromise performance. Under both structural and uniform missingness scenarios with up to 99.5% of features absent, APCFI’s accuracy is statistically on par with the much slower PCFI. It consistently and significantly outperforms other baselines, especially in high-missingness regimes.

These results confirm our hypothesis: **APCFI successfully combines the high accuracy of PCFI with the high efficiency of FP**, establishing it as a powerful and practical foundation for the data imputation stage of our framework.

Missing Type	Method	Cora	CiteSeer	PubMed	Photo	Computers	Average $\uparrow$	Avg.rank $\downarrow$
Structural	Zero*	44.15 $\pm$ 8.44	31.68 $\pm$ 4.50	48.20 $\pm$ 3.65	79.68 $\pm$ 2.17	72.03 $\pm$ 1.91	55.15	8.2
	LP*	<b>74.52 <math>\pm</math> 1.60</b>	<b>65.89 <math>\pm</math> 2.29</b>	72.25 $\pm$ 3.78	82.42 $\pm$ 2.57	76.28 $\pm$ 1.43	74.27	5
	sRMGCNN*	29.31 $\pm$ 0.71	24.21 $\pm$ 1.35	OOM	26.10 $\pm$ 1.89	37.15 $\pm$ 0.12	*	10.6
	GCNMF*	29.20 $\pm$ 1.13	24.50 $\pm$ 1.52	40.19 $\pm$ 0.95	26.82 $\pm$ 6.33	30.59 $\pm$ 9.81	30.26	10.4
	PaGNN*	30.55 $\pm$ 8.85	25.69 $\pm$ 3.98	50.82 $\pm$ 4.61	66.91 $\pm$ 3.99	56.50 $\pm$ 3.29	46.10	8.8
	FP*	72.84 $\pm$ 2.85	59.76 $\pm$ 2.47	72.69 $\pm$ 2.66	86.57 $\pm$ 1.50	77.42 $\pm$ 1.59	73.86	5.6
	FP	72.69 $\pm$ 2.84	59.65 $\pm$ 2.61	72.67 $\pm$ 2.88	86.46 $\pm$ 1.50	77.58 $\pm$ 1.55	73.81	6.2
	PCFI*	<b>75.49 <math>\pm</math> 2.10</b>	<b>66.18 <math>\pm</math> 2.75</b>	<b>74.66 <math>\pm</math> 2.26</b>	<b>87.70 <math>\pm</math> 1.29</b>	<b>79.25 <math>\pm</math> 1.19</b>	<b>76.66</b>	<b>1</b>
	PCFI	73.42 $\pm$ 2.74	<b>62.85 <math>\pm</math> 2.83</b>	<b>73.44 <math>\pm</math> 2.32</b>	<b>87.44 <math>\pm</math> 1.37</b>	<b>78.76 <math>\pm</math> 1.45</b>	<b>75.18</b>	<b>2.6</b>
	<b>APCFI(our)</b>	<b>73.44 <math>\pm</math> 2.51</b>	62.02 $\pm$ 2.79	<b>73.39 <math>\pm</math> 2.29</b>	<b>87.24 <math>\pm</math> 1.39</b>	<b>77.77 <math>\pm</math> 1.22</b>	<b>74.77</b>	<b>3.2</b>
Uniform	Zero*	62.63 $\pm$ 2.64	63.19 $\pm$ 1.83	54.70 $\pm$ 3.03	85.40 $\pm$ 1.33	79.49 $\pm$ 1.21	69.08	7.8
	LP*	74.52 $\pm$ 1.60	65.78 $\pm$ 2.29	72.25 $\pm$ 3.78	82.42 $\pm$ 2.57	76.28 $\pm$ 1.43	74.25	7.8
	sRMGCNN*	29.32 $\pm$ 0.74	24.66 $\pm$ 1.90	OOM	26.58 $\pm$ 1.68	37.16 $\pm$ 0.12	*	10.2
	GCNMF*	27.85 $\pm$ 2.27	24.29 $\pm$ 1.47	39.47 $\pm$ 0.76	25.98 $\pm$ 3.90	34.78 $\pm$ 4.69	30.47	10.8
	PaGNN*	53.75 $\pm$ 2.03	44.95 $\pm$ 2.59	60.24 $\pm$ 3.78	85.30 $\pm$ 1.05	78.04 $\pm$ 1.18	64.45	8.4
	FP*	77.55 $\pm$ 2.01	68.00 $\pm$ 2.16	73.88 $\pm$ 2.35	87.75 $\pm$ 1.07	81.47 $\pm$ 0.91	77.73	5.1
	FP	77.55 $\pm$ 2.01	67.47 $\pm$ 2.28	73.42 $\pm$ 2.33	88.05 $\pm$ 1.16	81.25 $\pm$ 1.05	77.55	5.5
	PCFI*	<b>78.53 <math>\pm</math> 1.39</b>	<b>69.40 <math>\pm</math> 1.85</b>	<b>76.44 <math>\pm</math> 1.64</b>	<b>88.60 <math>\pm</math> 1.30</b>	<b>81.79 <math>\pm</math> 0.70</b>	<b>78.95</b>	<b>1.2</b>
	PCFI	<b>78.57 <math>\pm</math> 1.55</b>	<b>68.60 <math>\pm</math> 1.71</b>	<b>75.88 <math>\pm</math> 1.74</b>	<b>88.54 <math>\pm</math> 1.18</b>	<b>81.76 <math>\pm</math> 0.9</b>	<b>78.67</b>	<b>1.8</b>
	<b>APCFI(our)</b>	<b>78.26 <math>\pm</math> 1.79</b>	<b>68.31 <math>\pm</math> 2.01</b>	<b>75.22 <math>\pm</math> 2.24</b>	<b>88.08 <math>\pm</math> 1.25</b>	<b>81.63 <math>\pm</math> 0.98</b>	<b>78.30</b>	<b>3.0</b>

**Table 5.3: Performance comparison of various feature imputation methods under structural and uniform missing scenarios.** The table reports the accuracy (%) on five datasets for each method. Methods marked with \* are directly introduced from the PCFI paper. APCFI achieves the best or near-best performance and average rank under both missing types. The best, second-best, and third-best results are highlighted in green, cyan, and red

### 5.3 Baseline Performance: Establishing Competitiveness with Complete Data

Before stress-testing our framework’s robustness against incomplete data, we must first establish its credibility as a high-performing model in a standard setting. To establish a strong baseline, this section evaluates the fundamental soundness of the iPaD framework and compares its performance with state-of-the-art (SOTA) models under standard, complete-data conditions.

The results, summarized in Tables 5.4 and 5.5, demonstrate a consistently strong and competitive performance, validating our overall design. The nature of this performance, however, shows interesting nuances depending on the GNN backbone.

Model	Cora	CiteSeer	PubMed	Avg Accuracy
GraphGPS (NeurIPS 2022)*	82.84 $\pm$ 1.03	72.73 $\pm$ 1.23	79.94 $\pm$ 0.26	78.50
NAGphormer (ICLR 2023)*	82.12 $\pm$ 1.18	71.47 $\pm$ 1.30	79.73 $\pm$ 0.28	77.77
Expformer (ICML2023)*	82.77 $\pm$ 1.38	71.63 $\pm$ 1.19	79.46 $\pm$ 0.35	77.95
GOAT (ICML 2023)*	83.18 $\pm$ 1.27	71.99 $\pm$ 1.26	79.13 $\pm$ 0.38	78.10
NodeFormer (NeurIPS 2022)*	82.20 $\pm$ 0.90	72.0 $\pm$ 1.10	79.90 $\pm$ 1.00	78.20
SGFormer (NeurIPS 2023)*	<b>84.50 <math>\pm</math> 0.80</b>	72.60 $\pm$ 0.20	80.30 $\pm$ 0.60	79.13
Polynormer (ICLR 2024)*	83.25 $\pm$ 0.89	72.31 $\pm$ 0.83	79.24 $\pm$ 0.73	78.27
GCN	83.77 $\pm$ 0.57	71.83 $\pm$ 0.28	79.89 $\pm$ 0.18	78.50
GLNN-GCN (ICLR 2022)	83.94 $\pm$ 0.77	<b>73.93 <math>\pm</math> 0.62</b>	<b>81.40 <math>\pm</math> 0.69</b>	<b>79.77</b> 1.27 $\uparrow$
KRD-GCN (ICML 2023)	83.96 $\pm$ 1.02	<b>74.26 <math>\pm</math> 0.88</b>	<b>83.08 <math>\pm</math> 0.43</b>	<b>80.43</b> 1.93 $\uparrow$
TunedGCN(NeurIPS 2024)	<b>84.52 <math>\pm</math> 0.62</b>	72.41 $\pm$ 0.71	80.88 $\pm$ 1.20	79.27 <b>0.77</b> $\uparrow$
<b>iPaD-GCN(our)</b>	<b>85.08 <math>\pm</math> 0.77</b>	<b>75.37 <math>\pm</math> 0.58</b>	<b>82.19 <math>\pm</math> 2.27</b>	<b>80.88</b> 2.38 $\uparrow$
GraphSAGE	80.60 $\pm$ 0.82	68.82 $\pm$ 0.92	73.45 $\pm$ 0.77	74.29
GLNN-SAGE (ICLR 2022)	80.89 $\pm$ 0.55	70.36 $\pm$ 0.46	76.23 $\pm$ 1.15	75.83 <b>1.54</b> $\uparrow$
KRD-SAGE (ICML 2023)	83.79 $\pm$ 0.55	72.00 $\pm$ 0.85	79.33 $\pm$ 0.90	77.78 <b>4.08</b> $\uparrow$
TunedGCN(NeurIPS 2024)	83.34 $\pm$ 0.58	70.15 $\pm$ 0.52	76.96 $\pm$ 0.78	76.82 <b>2.53</b> $\uparrow$
<b>iPaD-SAGE(our)</b>	83.16 $\pm$ 1.43	72.98 $\pm$ 0.99	78.59 $\pm$ 3.05	78.24 <b>3.95</b> $\uparrow$

**Table 5.4: SOTA comparison on complete classic datasets.** The best, second-best, and third-best results are highlighted in **green**, **cyan**, and **red**, respectively. Results marked with \* are from TunedGNN’s paper.

Model	Computer	Photo	WikiCS	CS	Physics	Avg Accuracy
GraphGPS (NeurIPS 2022)*	91.19 $\pm$ 0.54	95.06 $\pm$ 0.13	78.66 $\pm$ 0.49	93.93 $\pm$ 0.12	97.12 $\pm$ 0.19	91.19
NAGphormer (ICLR 2023)*	91.22 $\pm$ 0.14	95.49 $\pm$ 0.11	77.16 $\pm$ 0.72	95.75 $\pm$ 0.09	<b>97.34 <math>\pm</math> 0.03</b>	91.39
Expormer (ICML2023)*	91.47 $\pm$ 0.17	95.35 $\pm$ 0.22	78.54 $\pm$ 0.49	94.93 $\pm$ 0.01	96.89 $\pm$ 0.09	91.44
GOAT (ICML 2023)*	90.96 $\pm$ 0.90	92.96 $\pm$ 1.48	77.00 $\pm$ 0.77	94.21 $\pm$ 0.38	96.24 $\pm$ 0.24	90.27
NodeFormer (NeurIPS 2022)*	86.98 $\pm$ 0.62	93.46 $\pm$ 0.35	74.73 $\pm$ 0.94	95.64 $\pm$ 0.22	96.45 $\pm$ 0.28	89.45
SGFormer (NeurIPS 2023)*	91.99 $\pm$ 0.76	95.10 $\pm$ 0.47	73.46 $\pm$ 0.56	94.78 $\pm$ 0.20	96.60 $\pm$ 0.18	90.39
Polynormer (ICLR 2024)*	<b>93.68 <math>\pm</math> 0.21</b>	<b>96.46 <math>\pm</math> 0.26</b>	80.10 $\pm$ 0.67	95.53 $\pm$ 0.16	97.27 $\pm$ 0.08	<b>92.61</b>
GCN	71.77 $\pm$ 0.12	91.43 $\pm$ 0.11	79.60 $\pm$ 0.12	90.28 $\pm$ 0.12	94.50 $\pm$ 0.09	85.52
GLNN-GCN (ICLR 2022)	86.67 $\pm$ 0.45	94.72 $\pm$ 0.31	80.18 $\pm$ 0.19	94.71 $\pm$ 0.13	96.15 $\pm$ 0.17	90.49 <b>4.97</b> $\uparrow$
KRD-GCN (ICML 2023)	84.29 $\pm$ 0.47	93.31 $\pm$ 0.31	79.76 $\pm$ 0.38	94.85 $\pm$ 0.12	96.60 $\pm$ 0.08	89.76 <b>4.12</b> $\uparrow$
TunedGCN(NeurIPS 2024)	<b>93.83 <math>\pm</math> 0.25</b>	95.97 $\pm$ 0.30	<b>80.52 <math>\pm</math> 0.25</b>	95.89 $\pm$ 0.14	<b>97.37 <math>\pm</math> 0.06</b>	<b>92.72</b> <b>7.20</b> $\uparrow$
<b>iPaD-GCN(our)</b>	<b>93.00 <math>\pm</math> 0.34</b>	<b>96.23 <math>\pm</math> 0.22</b>	<b>81.11 <math>\pm</math> 0.40</b>	<b>96.14 <math>\pm</math> 0.11</b>	<b>97.31 <math>\pm</math> 0.10</b>	<b>92.76</b> <b>7.24</b> $\uparrow$
GraphSAGE	91.71 $\pm$ 0.24	95.75 $\pm$ 0.26	80.21 $\pm$ 0.32	95.70 $\pm$ 0.10	96.75 $\pm$ 0.12	92.02
GLNN-SAGE (ICLR 2022)	86.43 $\pm$ 0.56	95.18 $\pm$ 0.42	78.19 $\pm$ 0.41	92.93 $\pm$ 0.23	96.88 $\pm$ 0.10	89.92 <b>2.76</b> $\downarrow$
KRD-SAGE (ICML 2023)	90.21 $\pm$ 0.41	95.22 $\pm$ 0.24	79.88 $\pm$ 0.43	94.74 $\pm$ 0.30	97.14 $\pm$ 0.09	91.44 <b>1.22</b> $\downarrow$
TunedSAGE(NeurIPS 2024)	92.90 $\pm$ 0.26	<b>96.04 <math>\pm</math> 0.24</b>	<b>80.33 <math>\pm</math> 0.37</b>	<b>95.94 <math>\pm</math> 0.11</b>	96.96 $\pm$ 0.08	92.43 <b>0.41</b> $\uparrow$
<b>iPaD-SAGE(our)</b>	91.07 $\pm$ 0.66	95.37 $\pm$ 1.29	79.67 $\pm$ 0.14	<b>96.46 <math>\pm</math> 0.16</b>	96.70 $\pm$ 0.47	91.85 <b>0.17</b> $\downarrow$

**Table 5.5: SOTA comparison on complete medium datasets.** The best, second-best, and third-best results are highlighted in **green**, **cyan**, and **red**, respectively. Results marked with \* are from TunedGNN’s paper.

For the GCN backbone, the proposed iPaD-GCN not only achieves state-of-the-art results on multiple datasets, such as **CiteSeer (75.37%)** and **PubMed (82.09%)**, but also delivers substantial gains on medium-scale datasets, with the average accuracy improving from **85.64%** to **92.72%** (Table 5.5). This improvement is largely attributed to the mitigation of over-parameterization issues common in existing GNNs [25, 26], where excessive model capacity often leads to overfitting. The pruning stage effectively removes uninformative channels, reducing model complexity and improving generalization. Furthermore, the MP-KRD distillation process filters noise from the teacher model while maximizing the transfer of informative knowledge, enabling the student model to retain or even surpass the teacher’s performance.

For the SAGE backbone, iPaD-SAGE maintains a highly competitive performance, achieving an average accuracy of **91.68%** on medium datasets, which is only 0.18% lower than its strong tuned baseline (Table 5.5). This result indicates that the integration of pruning and MP-KRD does not introduce any meaningful performance degradation, while still delivering the benefits of reduced parameters, faster inference, and enhanced robustness.

In conclusion, by demonstrating a significantly improved GCN variant and a competitively on-par SAGE variant, iPaD establishes its credibility as a top-tier framework. Its ability to either substantially boost or maintain SOTA performance validates its fundamental design. Having confirmed this strong baseline, we can now proceed to evaluate the framework on its primary and most critical task: **maintaining robustness in the face of extreme data incompleteness**.

## 5.4 Validating MPP: Quantifying Gains in Training and Inference Efficiency

To rigorously validate the effectiveness of the Mirror Projection Pruning (MPP) mechanism, extensive experiments were conducted to quantify its impact on both training and inference efficiency. The key metrics considered include total training time, per-epoch time, average inference latency, and model parameter count.

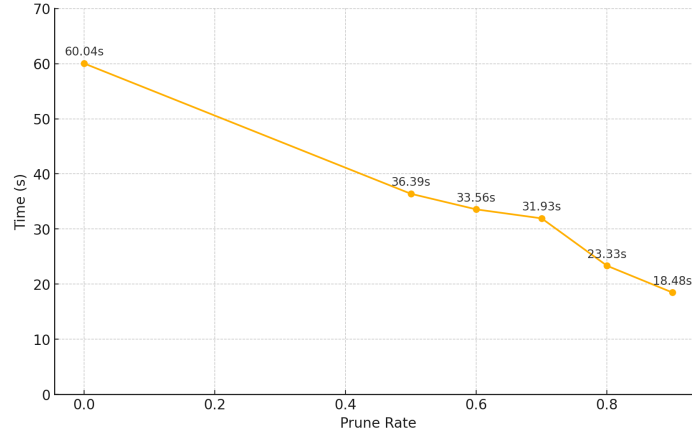
### 5.4.1 Impact on Computational Costs

Our results show that MPP leads to dramatic savings across multiple efficiency metrics. As shown in Table 5.6, increasing the prune rate from 0.0 to 0.9 drastically reduces the computational cost (measured in GFLOPs) and inference time. Specifically, GFLOPs decrease by over **10x** from **275.01** to **25.67**, while the inference time drops by more than **5x** from **15.76 ms** to just **2.79 ms**. This demonstrates the significant resource-saving benefits brought by the MPP pruning mechanism for model deployment.

Prune Rate	GFLOPs↓	Inference Time (ms)↓
0.0	275.01	15.76
0.5	132.70	6.54
0.6	104.96	6.46
0.7	78.15	5.69
0.8	51.72	3.69
0.9	25.67	2.79

**Table 5.6: Impact of MPP pruning on FLOPs and inference time of the GCN model.** Increasing the prune rate substantially reduces both computational cost (GFLOPs) and inference time, demonstrating the efficiency gains enabled by the MPP pruning strategy.

Furthermore, Figure 5.1 demonstrates that this efficiency gain extends to the entire training pipeline. As the prune rate increases, the total training time—which includes the overhead of the MPP process itself—decreases substantially from **60.04s** to just **18.48s**, a more than three-fold reduction. This confirms that MPP is a highly effective strategy for not only accelerating deployment but also shortening the model development cycle.



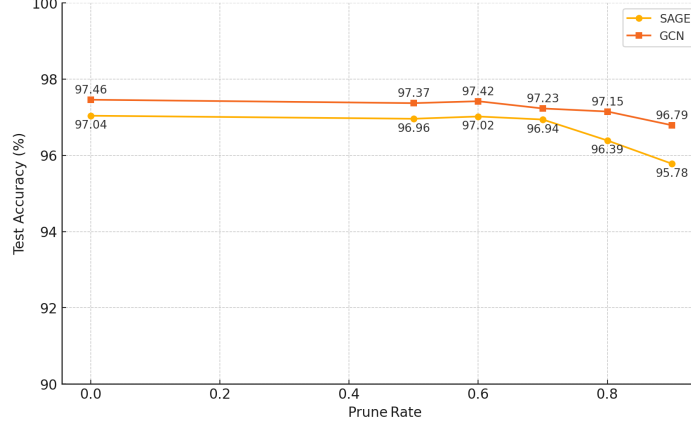
**Figure 5.1: Total training time (including the MPP step) of the MPP-pruned GCN model on the CS dataset under different prune rates.** Higher prune rates lead to a significant reduction in total training time.

#### 5.4.2 Accuracy-Efficiency Trade-off and the “Sweet Spot”

Having established the significant efficiency gains, we now analyze the trade-off with predictive performance. Figure 5.2 plots the test accuracy of both GCN and SAGE models against the prune rate on the Physics dataset.

The results reveal a crucial finding: the existence of a wide “**efficiency sweet spot**”. For both

models, the accuracy remains remarkably stable with negligible degradation up to a prune rate of 0.7. For instance, the GCN model’s accuracy only drops from **97.46%** to **97.23%** at a 70% prune rate, while achieving the substantial computational savings described above. This indicates that a large portion of the model’s parameters can be removed via MPP without harming its core expressive power. While higher prune rates (0.8-0.9) offer further efficiency gains, they come at the cost of a more noticeable, yet still graceful, decline in accuracy, making them suitable for scenarios where efficiency is the absolute priority.



**Figure 5.2:** Test accuracy of SAGE and GCN models on the Physics dataset under different MPP prune rates. The accuracy remains highly stable up to a prune rate of 0.7, revealing a wide “sweet spot” for efficiency gains without a significant performance penalty.

In conclusion, our experiments provide strong evidence for the efficacy of the MPP module. It enables significant reductions in both training and inference costs while maintaining high accuracy up to aggressive pruning rates, confirming it as a practical and powerful component of the **iPaD** pipeline.

## 5.5 Stress Test: Evaluating Robustness under Extreme Data Corruption

The central promise of **iPaD** is its robustness. To rigorously evaluate this property, the framework is tested under extreme levels of data corruption, specifically targeting both the graph structure (edge missing) and node attributes (feature missing).

### 5.5.1 Performance under Edge Missing

We evaluate the framework’s robustness against structural corruption by varying the edge missing ratio from 0% (full edge) to 90%. The results, summarized in Tables 5.7 and 5.8, show that the distilled **iPaD** student consistently matches or surpasses its GNN teacher across a wide range of edge missing rates.



The improvement is especially notable under severe structural degradation (90% edge missing rate). For instance, on the Photo dataset, the “TunedGCN” teacher drops from 95.86% (full edge) to 91.56%, whereas “iPaD-GCN” maintains 93.75%, reducing the accuracy loss by more than half (from  $-4.30$  to  $-2.41$  percentage points). A similar pattern is observed across most datasets, where **iPaD** achieves smaller relative performance drops compared to its teacher. These results indicate that the proposed framework not only preserves competitive accuracy under mild corruption but also provides superior tolerance to extreme edge missing scenarios.

Dataset	Full edge	30% Missing	60% Missing	90% Missing
Cora	83.77	76.76	72.82	64.48
CiteSeer	72.33	68.14	63.44	60.88
PubMed	80.43	76.58	75.36	72.82
Computers	93.54	92.24	91.12	87.36
Photo	95.86	95.25	94.35	91.56
WikiCS	80.52	78.61	77.27	73.50
CS	95.89	94.71	94.27	93.20
Physics	97.36	96.58	96.45	95.92
<b>Average</b>	<b>87.46</b>	<b>84.86</b> 2.60 ↓	<b>83.14</b> 4.32 ↓	<b>79.97</b> 7.49 ↓

**Table 5.7: Test accuracy of TunedGCN [1] on various datasets under different edge missing ratios.** The table reports accuracy under full edge, 30%, 60%, and 90% edge missing conditions. Average accuracy and relative decreases are shown at the bottom.

Dataset	Full edge	30% Missing	60% Missing	90% Missing
Cora	85.02	80.14	76.10	68.36
CiteSeer	75.37	70.64	68.48	65.00
PubMed	82.09	79.46	77.24	74.04
Computers	92.90	92.09	91.58	89.75
Photo	96.16	95.93	95.70	93.75
WikiCS	81.11	79.72	79.35	77.31
CS	96.14	95.54	95.41	94.28
Physics	97.31	96.86	96.47	96.08
<b>Average</b>	<b>88.26</b>	<b>86.29</b> 1.97 ↓	<b>85.04</b> 3.22 ↓	<b>82.32</b> 5.94 ↓

**Table 5.8: Test accuracy of iPaD-GCN on various datasets under different edge missing ratios.** The table reports accuracy under full edge, 30%, 60%, and 90% edge missing conditions. Average accuracy and relative decreases are provided at the bottom.

### 5.5.2 Performance under Feature Missing

The most challenging scenario involves evaluating the framework’s resilience when node features are severely corrupted. This is investigated under two distinct patterns—uniform and structural missingness—at an extreme rate of 99.5%.

Across all four testbeds, a clear and consistent narrative emerges: the performance of the final distilled student model (**iPaD-GCN**) is fundamentally governed by the performance of its teacher (the **MPP-GCN(tuned)** model). This demonstrates a highly successful and faithful knowledge transfer, while also highlighting the inherent limits of distillation when input information is scarce.

#### Under Uniform Feature Missing.

As shown in Table 5.9 (classic datasets) and Table 5.10 (medium datasets), our framework shows remarkable stability. For instance, on the medium datasets, the **MPP-SAGE(tuned)** teacher achieves an average accuracy of 86.69%. **iPaD-SAGE** student closely follows at 86.53%, indicating a near-perfect transfer of knowledge with almost no performance degradation during distillation. This pattern holds true for the GCN backbone as well.

Model	Cora	CiteSeer	PubMed	Accuracy
GCN w Zero	23.29 $\pm$ 0.78	16.70 $\pm$ 1.99	44.04 $\pm$ 0.94	23.29
GCN w FP	68.55 $\pm$ 0.21	44.72 $\pm$ 0.46	70.97 $\pm$ 0.22	68.55 45.3 $\uparrow$
TunedGCN w FP	74.46 $\pm$ 0.47	51.61 $\pm$ 1.29	72.68 $\pm$ 0.74	74.46 51.2 $\uparrow$
MPP-TunedGCN w APCFI(our)	75.55 $\pm$ 0.80	51.20 $\pm$ 0.73	75.22 $\pm$ 0.68	75.22 51.9 $\uparrow$
iPaD-GCN(our)	75.21 $\pm$ 0.36	53.96 $\pm$ 1.27	72.10 $\pm$ 4.90	75.21 51.9 $\uparrow$
GraphSAGE w Zero	34.14 $\pm$ 1.34	23.98 $\pm$ 1.05	39.92 $\pm$ 2.66	34.14
GraphSAGE w FP	71.24 $\pm$ 0.47	50.77 $\pm$ 0.85	74.87 $\pm$ 0.68	71.24 37.1 $\uparrow$
TunedSAGE w FP	74.10 $\pm$ 1.56	53.30 $\pm$ 0.85	73.26 $\pm$ 1.05	74.10 40.0 $\uparrow$
MPP-TunedSAGE w APCFI(our)	74.90 $\pm$ 1.22	54.72 $\pm$ 0.51	76.47 $\pm$ 0.62	74.90 40.8 $\uparrow$
iPaD-SAGE(our)	74.52 $\pm$ 0.56	53.94 $\pm$ 0.82	76.19 $\pm$ 0.75	74.52 40.4 $\uparrow$

**Table 5.9: Comparison under 99.5% uniform feature missing: vanilla GNN, MPP-Tuned, and iPaD with common imputation methods on classic datasets.** Model names indicate the imputation strategy. The best results are highlighted in **green**, the second-best in **cyan**, and the third-best in **red**.

Model	Computer	Photo	WikiCS	CS	Physics	Accuracy
SGFormer w FP	81.98 $\pm$ 3.29	87.18 $\pm$ 2.16	75.02 $\pm$ 0.30	89.21 $\pm$ 1.56	93.65 $\pm$ 0.79	85.41
Polynormer w FP	80.75 $\pm$ 3.15	89.88 $\pm$ 2.79	74.38 $\pm$ 0.50	90.24 $\pm$ 0.25	42.42 $\pm$ 15.97	75.53
GCN w Zero	38.50 $\pm$ 0.04	28.09 $\pm$ 0.19	51.03 $\pm$ 0.26	22.74 $\pm$ 0.00	50.41 $\pm$ 0.00	38.15
GCN w FP	38.51 $\pm$ 0.38	65.25 $\pm$ 4.78	71.36 $\pm$ 0.31	23.13 $\pm$ 1.17	79.25 $\pm$ 3.33	55.50 17.4 $\uparrow$
TunedGCN w FP	88.81 $\pm$ 0.34	71.48 $\pm$ 1.38	74.51 $\pm$ 0.35	84.00 $\pm$ 0.51	94.60 $\pm$ 0.21	82.68 44.5 $\uparrow$
MPP-TunedGCN w APCFI(our)	88.19 $\pm$ 0.29	86.47 $\pm$ 1.13	76.64 $\pm$ 0.34	91.09 $\pm$ 0.15	94.52 $\pm$ 0.16	87.38 49.2 $\uparrow$
iPaD-GCN(our)	86.72 $\pm$ 0.36	86.25 $\pm$ 0.29	74.49 $\pm$ 0.32	90.00 $\pm$ 0.15	94.87 $\pm$ 0.08	86.47 48.3 $\uparrow$
GraphSAGE w Zero	55.26 $\pm$ 3.42	72.07 $\pm$ 1.56	42.94 $\pm$ 0.63	60.72 $\pm$ 3.67	50.41 $\pm$ 0.00	56.28
GraphSAGE w FP	79.74 $\pm$ 0.45	86.48 $\pm$ 0.41	71.47 $\pm$ 0.48	89.77 $\pm$ 0.35	86.67 $\pm$ 1.73	82.83 26.6 $\uparrow$
TunedSAGE w FP	86.05 $\pm$ 0.35	74.72 $\pm$ 3.51	90.23 $\pm$ 0.55	87.08 $\pm$ 0.50	85.18 $\pm$ 2.54	85.18 28.9 $\uparrow$
MPP-TunedSAGE w APCFI(our)	86.69 $\pm$ 0.35	86.20 $\pm$ 3.39	75.12 $\pm$ 0.22	90.80 $\pm$ 0.16	92.69 $\pm$ 1.42	86.30 30.0 $\uparrow$
iPaD-SAGE(our)	84.57 $\pm$ 0.68	87.73 $\pm$ 2.15	73.63 $\pm$ 0.41	91.69 $\pm$ 0.27	94.59 $\pm$ 0.25	86.44 30.2 $\uparrow$

**Table 5.10: Comparison under 99.5% uniform feature missing: vanilla GNN, SGFormer, Polynormer, TunedGNN, MPP-Tuned, and iPaD with common imputation methods on medium datasets.** Model names indicate the imputation strategy. The best results are highlighted in **green**, the second-best in **cyan**, and the third-best in **red**.

### Under Structural Feature Missing.

This same trend is observed under the challenging structural missingness scenario, as shown in Table 5.11 and Table 5.12. The **MPP-Tuned** models, equipped with our APCFI imputer, consistently set the performance benchmark. The corresponding **iPaD** student models again track this performance closely. For example, on the medium datasets, the **MPP-GCN(tuned)** teacher scores 84.25%, while the **iPaD-GCN** student achieves a comparable 83.24%.

Model	Cora	CiteSeer	PubMed	Accuracy
GCN w Zero	23.29 $\pm$ 1.00	16.70 $\pm$ 0.34	44.04 $\pm$ 0.08	28.01
GCN w FP	68.55 $\pm$ 0.90	44.72 $\pm$ 0.24	70.97 $\pm$ 0.36	61.41 33.4 $\uparrow$
TunedGCN w FP	<b>70.84 <math>\pm</math> 1.06</b>	46.80 $\pm$ 0.70	<b>73.13 <math>\pm</math> 1.18</b>	<b>63.59</b> 35.6 $\uparrow$
<b>MPP-TunedGCN w APCFI(our)</b>	<b>70.46 <math>\pm</math> 0.75</b>	50.62 $\pm$ 0.57	71.77 $\pm$ 0.53	<b>64.28</b> 36.3 $\uparrow$
<b>iPaD-GCN(our)</b>	60.94 $\pm$ 2.51	<b>52.22 <math>\pm</math> 0.53</b>	72.35 $\pm$ 0.35	61.84 33.8 $\uparrow$
GraphSAGE w Zero	34.14 $\pm$ 0.16	23.98 $\pm$ 0.74	41.28 $\pm$ 0.36	32.53
GraphSAGE w FP	<b>71.24 <math>\pm</math> 0.89</b>	<b>50.77 <math>\pm</math> 0.75</b>	<b>74.97 <math>\pm</math> 0.48</b>	<b>65.66</b> 32.5 $\uparrow$
<b>TunedSAGE w FP</b>	51.72 $\pm$ 5.41	48.39 $\pm$ 1.41	<b>74.52 <math>\pm</math> 0.97</b>	58.21 25.1 $\uparrow$
<b>MPP-TunedSAGE w APCFI(our)</b>	65.96 $\pm$ 0.95	50.57 $\pm$ 2.00	69.50 $\pm$ 3.91	62.01 28.9 $\uparrow$
<b>iPaD-SAGE(our)</b>	59.56 $\pm$ 4.11	<b>52.08 <math>\pm</math> 1.60</b>	64.31 $\pm$ 4.63	58.65 25.5 $\uparrow$

**Table 5.11: Comparison under 99.5% structural feature missing: vanilla GNN, MPP-Tuned, and iPaD with common imputation methods on classic datasets.** Model names indicate the imputation strategy. The best results are highlighted in **green**, the second-best in **cyan**, and the third-best in **red**.

Model	Computer	Photo	WikiCS	CS	Physics	Avg Accuracy
SGFormer w FP	70.04 $\pm$ 3.33	80.10 $\pm$ 2.96	74.54 $\pm$ 0.42	74.94 $\pm$ 2.48	90.56 $\pm$ 1.46	78.04
Polynormer w FP	81.70 $\pm$ 2.81	86.13 $\pm$ 0.63	70.22 $\pm$ 0.57	87.81 $\pm$ 0.40	41.96 $\pm$ 16.89	73.56
GCN w Zero	38.54 $\pm$ 0.02	28.27 $\pm$ 0.20	39.32 $\pm$ 0.63	22.74 $\pm$ 0.00	50.41 $\pm$ 0.00	35.86
GCN w FP	38.64 $\pm$ 0.69	59.77 $\pm$ 0.39	63.12 $\pm$ 0.32	25.83 $\pm$ 0.07	83.35 $\pm$ 3.18	54.14 18.3 $\uparrow$
TunedGCN w FP	84.88 $\pm$ 0.31	65.65 $\pm$ 4.22	75.46 $\pm$ 0.37	71.24 $\pm$ 2.05	92.74 $\pm$ 0.14	77.99 42.1 $\uparrow$
MPP-TunedGCN w APCFI(our)	87.97 $\pm$ 0.26	77.92 $\pm$ 4.17	74.13 $\pm$ 0.39	87.57 $\pm$ 0.20	93.68 $\pm$ 0.13	84.25 48.4 $\uparrow$
iPaD-GCN(our)	86.47 $\pm$ 0.11	79.55 $\pm$ 4.53	71.63 $\pm$ 0.39	85.20 $\pm$ 0.26	93.33 $\pm$ 0.06	83.24 47.4 $\uparrow$
GraphSAGE w Zero	47.75 $\pm$ 1.38	45.49 $\pm$ 4.53	52.71 $\pm$ 0.36	33.22 $\pm$ 0.64	50.41 $\pm$ 0.03	45.92
GraphSAGE w FP	79.73 $\pm$ 1.00	86.05 $\pm$ 0.60	70.71 $\pm$ 0.26	83.79 $\pm$ 0.24	84.66 $\pm$ 2.51	80.99 35.1 $\uparrow$
TunedSAGE w FP	83.46 $\pm$ 1.04	80.93 $\pm$ 3.96	72.50 $\pm$ 0.65	83.79 $\pm$ 0.30	90.88 $\pm$ 0.27	82.31 36.4 $\uparrow$
MPP-TunedSAGE w APCFI(our)	85.32 $\pm$ 0.51	87.35 $\pm$ 0.49	73.35 $\pm$ 0.34	88.71 $\pm$ 0.24	89.31 $\pm$ 1.44	84.81 38.9 $\uparrow$
iPaD-SAGE(our)	83.36 $\pm$ 0.39	86.46 $\pm$ 0.44	74.74 $\pm$ 0.34	87.93 $\pm$ 0.18	93.05 $\pm$ 0.08	84.31 38.4 $\uparrow$

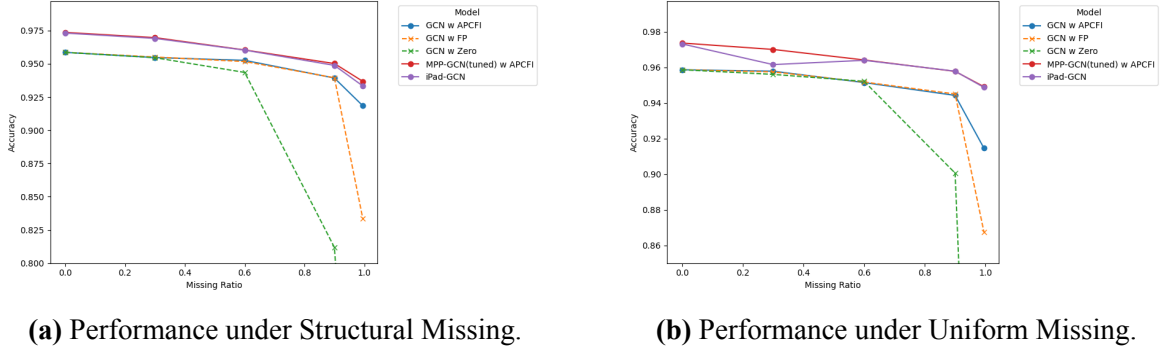
**Table 5.12: Comparison under 99.5% structural feature missing: vanilla GNN, SGFormer, Polynormer, TunedGNN, MPP-Tuned, and iPaD with common imputation methods on medium datasets.** Model names indicate the imputation strategy. The best results are highlighted in **green**, the second-best in **cyan**, and the third-best in **red**.

### Interpretation of Results.

Figure 5.3 presents the robustness of various GCN configurations under increasing feature missing ratios on the Physics dataset. Across both structural and uniform missing settings, the **iPaD-GCN** (student) and **MPP-GCN(tuned)** (teacher) models consistently achieve high accuracy, even at an extreme missing rate of 90%. This demonstrates the effectiveness and robustness of the proposed framework in handling severe feature incompleteness.

Notably, the performance of the distilled student model is always on par with, but never surpasses, that of its teacher. This observation highlights a key boundary of knowledge distillation in the context of feature missing: while the student can faithfully inherit the relational knowledge captured by the teacher, it cannot recover information that was never present in the input. Consequently, the quality of feature imputation and the upper bound of the teacher’s performance together define a hard ceiling for the student. In contrast, as discussed in the

following section, in the case of edge (structural) missing, the student occasionally surpasses the teacher by filtering out noise through distillation. This distinction underlines the fundamental differences in the limits of distillation under feature versus structural incompleteness.



**Figure 5.3: Robustness comparison of different GCN configurations under varying feature missing ratios on the Physics dataset.** Both the ‘MPP-GCN(tuned)’ teacher and the final ‘iPaD-GCN’ student model demonstrate superior robustness, maintaining high accuracy even at a 90% missing ratio. In contrast, the baseline ‘GCN w/ Zero’ imputation suffers a catastrophic performance collapse.

## 5.6 Dissecting the Framework: An Ablation Study on Core Modules

To validate the central claim that the strength of iPaD arises from the synergistic interaction of its components, a comprehensive ablation study is conducted. Each of the three core modules—the **tuned** backbone, **Mirror Projection Pruning (MPP)**, and **Mirror Projection Knowledge-inspired Reliable Distillation (MP-KRD)**—is systematically enabled to analyze their individual and combined contributions. Experiments are performed on the Photo and CS datasets using both GCN and SAGE backbones. This analysis assesses whether all modules are essential for achieving the optimal balance of accuracy, training efficiency, and inference speed.

### 5.6.1 Ablation Results with GCN Backbone

Table 5.13 presents the ablation results of the **Tuned**, **MPP**, and **MP-KRD** modules on the Photo and CS datasets (GCN backbone), evaluating node classification accuracy, training speed, and inference latency under various module combinations.

The baseline model (row 1, all modules disabled) demonstrates the standard performance-efficiency trade-off for GCNs. Enabling the Tuned module leads to a substantial increase in accuracy but incurs higher inference time and reduced training speed, highlighting the cost of improved expressiveness. The introduction of the MPP module markedly improves computational efficiency, both in training and inference, while the integration of MP-KRD yields

Tuned	MPP	MP-KRD	Photo			CS		
			Accuracy $\uparrow$	Training Speed (epoch/s) $\uparrow$	Inference Time (ms) $\downarrow$	Accuracy $\uparrow$	Training Speed (epoch/s) $\uparrow$	Inference Time (ms) $\downarrow$
$\times$	$\times$	$\times$	$91.43 \pm 0.10$	$118.2 \pm 6.33$	$3.09 \pm 0.04$	$90.28 \pm 0.12$	$45.1 \pm 0.18$	$8.45 \pm 0.41$
$\checkmark$	$\times$	$\times$	<b><math>95.97 \pm 0.30</math></b>	<b><math>28.8 \pm 2.24</math></b>	<b><math>12.69 \pm 1.37</math></b>	<b><math>95.89 \pm 0.14</math></b>	<b><math>23.3 \pm 0.42</math></b>	<b><math>16.55 \pm 1.00</math></b>
$\times$	$\checkmark$	$\times$	$90.69 \pm 0.17$	$166.1 \pm 0.69$	$2.47 \pm 0.20$	$90.04 \pm 0.06$	$77.5 \pm 0.27$	$5.50 \pm 0.17$
$\times$	$\times$	$\checkmark$	$93.94 \pm 0.39$	$165.5 \pm 6.33^*$	$1.27 \pm 0.25$	$93.89 \pm 0.13$	$52.1 \pm 0.18^*$	$6.03 \pm 0.61$
$\checkmark$	$\checkmark$	$\times$	$95.92 \pm 0.12$	$47.2 \pm 1.15$	$7.52 \pm 0.54$	$96.01 \pm 0.08$	$47.7 \pm 0.20$	$8.16 \pm 1.01$
$\checkmark$	$\times$	$\checkmark$	<b><math>96.21 \pm 0.23</math></b>	<b><math>89.0 \pm 2.24^*</math></b>	<b><math>2.59 \pm 0.34</math></b>	<b><math>96.34 \pm 0.12</math></b>	<b><math>27.2 \pm 0.42^*</math></b>	<b><math>10.36 \pm 0.58</math></b>
$\times$	$\checkmark$	$\checkmark$	$93.46 \pm 0.31$	$165.9 \pm 0.69^*$	$1.20 \pm 0.07$	$93.79 \pm 0.15$	$75.9 \pm 0.27^*$	$3.95 \pm 0.34$
$\checkmark$	$\checkmark$	$\checkmark$	<b><math>96.23 \pm 0.22</math></b>	<b><math>118.4 \pm 1.15^*</math></b>	<b><math>1.65 \pm 0.07</math></b>	<b><math>96.14 \pm 0.10</math></b>	<b><math>48.9 \pm 0.20^*</math></b>	<b><math>6.34 \pm 0.47</math></b>

**Table 5.13: Ablation study of tuned, MPP, and MP-KRD modules on the Photo and CS datasets (GCN Model).** The table reports accuracy, training speed (epoch/s), and inference time (ms) for different combinations of modules. When MP-KRD is used, (\*) indicates the speed measured during knowledge distillation.

further gains in model compactness and inference speed through knowledge distillation.

Of particular note, the configuration that combines all three modules (Tuned, MPP, and MP-KRD) achieves the optimal balance: it consistently yields the highest or near-highest accuracy, fastest inference, and robust training efficiency across both datasets. Intermediate configurations confirm that each module addresses the limitations introduced by the others, and their joint integration delivers the best overall trade-off.

### 5.6.2 Ablation Results with SAGE Backbone

Table 5.14 presents the ablation results for the combination of Tuned, MPP, and MP-KRD modules on the Photo and CS datasets, using the SAGE backbone. Each configuration is evaluated in terms of accuracy, training speed (epoch/s), and inference time (ms).

The baseline SAGE model, with all modules disabled, establishes a reference for both performance and computational efficiency. Introducing the Tuned module significantly enhances accuracy but results in increased inference latency and reduced training speed, illustrating the classic accuracy-efficiency trade-off. The application of the MPP module effectively recovers computational efficiency, improving both training speed and inference time, although the accuracy does not surpass that of the tuned backbone.

Tuned	MPP	MP-KRD	Photo			CS		
			Accuracy $\uparrow$	Training Speed (epoch/s) $\uparrow$	Inference Time (ms) $\downarrow$	Accuracy $\uparrow$	Training Speed (epoch/s) $\uparrow$	Inference Time (ms) $\downarrow$
$\times$	$\times$	$\times$	$95.14 \pm 0.20$	$103.9 \pm 0.46$	$4.31 \pm 0.27$	$95.67 \pm 0.10$	$20.1 \pm 0.16$	$20.60 \pm 0.51$
$\checkmark$	$\times$	$\times$	<b><math>96.04 \pm 0.24</math></b>	<b><math>73.4 \pm 2.25</math></b>	<b><math>5.00 \pm 0.41</math></b>	<b><math>95.94 \pm 0.11</math></b>	<b><math>15.7 \pm 0.08</math></b>	<b><math>25.08 \pm 1.56</math></b>
$\times$	$\checkmark$	$\times$	$95.14 \pm 0.32$	$118.8 \pm 1.15$	$3.53 \pm 0.24$	$95.59 \pm 0.10$	$25.3 \pm 0.75$	$18.20 \pm 0.41$
$\times$	$\times$	$\checkmark$	$94.70 \pm 0.39$	$139.6 \pm 0.46^*$	$1.51 \pm 0.33$	$95.75 \pm 0.14$	$32.0 \pm 0.16^*$	$9.32 \pm 0.50$
$\checkmark$	$\checkmark$	$\times$	<b><math>95.83 \pm 0.64</math></b>	<b><math>90.49 \pm 0.58</math></b>	<b><math>4.63 \pm 0.10</math></b>	$95.87 \pm 0.15$	$21.3 \pm 0.52$	$20.58 \pm 2.10$
$\checkmark$	$\times$	$\checkmark$	<b><math>96.01 \pm 0.21</math></b>	<b><math>81.2 \pm 2.25^*</math></b>	<b><math>2.13 \pm 0.14</math></b>	<b><math>96.45 \pm 0.12</math></b>	<b><math>20.2 \pm 0.08^*</math></b>	<b><math>14.3 \pm 1.08</math></b>
$\times$	$\checkmark$	$\checkmark$	$93.82 \pm 1.07$	$150.5 \pm 1.15^*$	$1.31 \pm 0.18$	$95.67 \pm 0.16$	$53.8 \pm 0.75^*$	$5.77 \pm 0.53$
$\checkmark$	$\checkmark$	$\checkmark$	<b><math>95.33 \pm 1.29</math></b>	<b><math>96.8 \pm 1.03^*</math></b>	<b><math>1.92 \pm 0.03</math></b>	<b><math>96.46 \pm 0.16</math></b>	<b><math>35.8 \pm 0.52^*</math></b>	<b><math>8.35 \pm 0.61</math></b>

**Table 5.14: Ablation study of tuned, MPP, and MP-KRD modules on the Photo and CS datasets (GraphSAGE Model).** The table reports accuracy, training speed (epoch/s), and inference time (ms) for different combinations of modules. When MP-KRD is used, (\*) indicates the speed measured during knowledge distillation.

Incorporating the MP-KRD module, either alone or in combination with other modules, results in a lightweight student model, characterized by improved inference speed and competitive accuracy. Notably, the complete pipeline—integrating **Tuned**, **MPP**, and **MP-KRD**—achieves the best overall trade-off. This configuration delivers the highest or near-highest accuracy and the fastest inference speed across both datasets. For example, on the CS dataset, this combination yields an accuracy of **96.46%** with an inference time of only **8.35 ms**.

These findings further validate the necessity of a holistic integration of all modules. The joint design not only addresses the individual limitations of each component but also provides a synergistic effect, resulting in a lightweight, accurate, and efficient GNN model suitable for real-world deployment.

### 5.6.3 Conclusion of Ablation Study

These ablation studies provide compelling evidence that the joint design of the **tuned backbone**, **MPP**, and **MP-KRD** is critical. Each module systematically addresses the limitations introduced by the previous one, demonstrating a powerful synergistic effect. The results confirm that it is the holistic integration of these components, rather than any single module alone, that enables iPaD to achieve its final, optimal balance of accuracy, training efficiency, and inference speed.

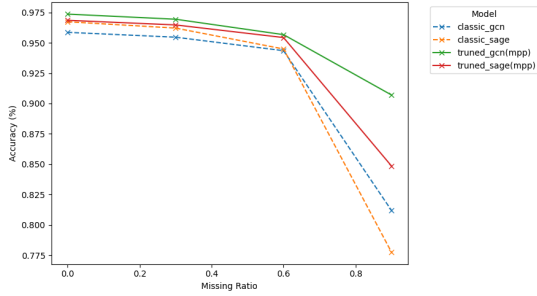


## 5.7 Component-Level Validation: Internal Ablation Studies

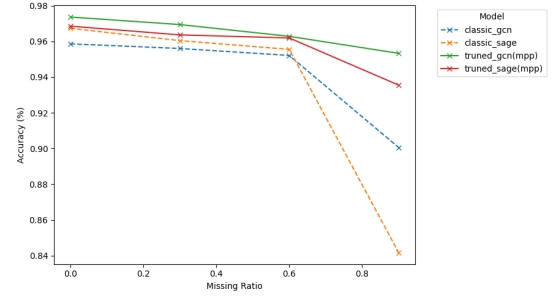
Having demonstrated the synergistic success of the overall iPaD pipeline, we now zoom in to validate the specific design choices within our key modules. Targeted ablation studies are conducted to assess the robustness of the TunedGNN architecture, the contribution of the Soft-PC mechanism to APCFI, and the necessity of all components within the MP-KRD module.

### 5.7.1 Validating the Robustness of the TunedGNN Backbone

This experiment validates our choice of **TunedGNN** as the foundational teacher architecture. We compare its performance against classic GNN counterparts under varying degrees of feature missingness. As illustrated in Figures 5.4, the **TunedGNN** variants consistently and substantially outperform the classic GNNs, especially at high missing rates. The accuracy of classic GNNs often collapses, while the **TunedGNN** models exhibit a much more graceful degradation. This confirms that our chosen architecture provides the essential resilience needed to serve as a reliable teacher in imperfect data conditions.



(a) Comparison under Structural missing.



(b) Comparison under Uniform missing.

**Figure 5.4: Robustness comparison of tuned and classic GCN/SAGE architectures on the Physics dataset.** All models are evaluated using Zero imputation for missing features, providing a uniform baseline for comparison. Both tuned GCN (green) and tuned SAGE (red) with MPP maintain higher accuracy than their classic counterparts as the missing ratio increases, demonstrating the impact of model tuning and pruning on robustness.

### 5.7.2 Validating the APCFI Design: A Component-wise Ablation Study

To justify the specific design of the APCFI module, a detailed component-wise ablation study is conducted. This analysis isolates the individual contributions of the three key innovations —**ASDE**, **DJCD**, and **Soft-PC**—in terms of both computational efficiency and final imputation accuracy.

### Analysis of Efficiency Contributions

Table 5.15 presents a runtime breakdown of the APCFI components, demonstrating how each innovation contributes to the overall efficiency breakthrough.

The baseline configuration (row 1), which mimics the slow, sequential nature of PCFI, suffers from extremely high runtimes. The results clearly show that:

- Enabling our proposed **ASDE** module (row 2) reduces the SPD runtime catastrophically, for instance, from 148.89s down to a mere 0.021s on the CiteSeer dataset. This validates ASDE as the key to solving the confidence calculation bottleneck.
- Subsequently, enabling **DJCD** (row 3) provides a similar dramatic speedup for the diffusion step, reducing its runtime from 101.7s to just 0.057s. This confirms that our joint channel diffusion is the key to efficient feature propagation.

ASDE	DJCD	Soft-PC	CiteSeer		Photo	
			SPD-RunTime(s)↓	Diffusion-RunTime(s)↓	SPD-RunTime(s)↓	Diffusion-RunTime(s)↓
✗	✗	✗	148.89	101.75	10.18	23.39
✓	✗	✗	0.021	101.7	0.021	23.39
✓	✓	✗	0.021	0.057	0.021	0.23
✓	✓	✓	<b>0.021</b>	<b>0.061</b>	<b>0.021</b>	<b>0.24</b>

**Table 5.15:** Ablation study of APCFI components on runtime performance. The table reports the runtime (in seconds) for the two main stages of imputation (SPD calculation and Diffusion). Each row corresponds to enabling a key innovation, demonstrating how each component contributes to the final efficiency.

### Analysis of Accuracy Contribution

Having established the efficiency gains from ASDE and DJCD, we now investigate the role of the final component, **Soft-PC**. Notably, Table 5.15 (comparing rows 3 and 4) shows that adding Soft-PC introduces **negligible computational overhead**. Its contribution lies in enhancing accuracy.

The results of this targeted ablation study, as presented in Table 5.16, reveal a nuanced impact of the Soft-PC mechanism. While the inclusion of Soft-PC results in a marginal performance decrease on three out of five datasets (Cora, CiteSeer, and Photo), it yields a substantial accuracy improvement on the Computers dataset (from 80.65% to 81.63%). This notable gain is sufficient to increase the overall average accuracy from 78.23% to 78.40%. These findings

Method	Cora	CiteSeer	PubMed	Photo	Computers	Average $\uparrow$
APCFI wo soft-PC	78.35	68.37	75.69	88.09	80.65	78.23
<b>APCFI</b>	<b>78.29</b>	<b>68.31</b>	<b>75.69</b>	<b>88.08</b>	<b>81.63</b>	<b>78.40</b>

**Table 5.16: Ablation study of soft-PC in APCFI.** The table reports test accuracy on Cora, CiteSeer, PubMed, Photo, and Computers datasets for the APCFI method with and without soft-PC. Avg Accuracy indicates the average accuracy across all datasets.

indicate that the adaptive, context-aware weighting provided by Soft-PC is particularly advantageous for specific graph structures, such as the co-purchase network in the Computers dataset. Although the effect is not uniformly positive across all benchmarks, the potential for significant improvement on certain graph topologies, combined with a slight overall performance increase, supports the inclusion of Soft-PC as a valuable component to enhance the general robustness of the **APCFI** model.

## Conclusion

In conclusion, this component-wise analysis validates our design. **ASDE** and **DJCD** are responsible for the massive efficiency breakthroughs that make APCFI practical for large-scale graphs, while **Soft-PC** provides a final layer of performance refinement at virtually no additional cost. This justifies the inclusion of all three innovations in our final APCFI model.

### 5.7.3 Validating the MP-KRD Design

Finally, the complex MP-KRD module is dissected to clarify the contribution of its core components: the reliability-aware curriculum (**KRD**), the structurally-aligned student (**MP-MLP**), and the feature-based loss (**CWD-Loss**).

#### For classic datasets

Tables 5.17 and 5.18 present the ablation results of the MP-KRD module on Cora, CiteSeer, and PubMed using GCN and SAGE backbones. Across both architectures, integrating all three components—KRD, MP-MLP, and CWD-loss—consistently yields the highest average accuracy (**80.83%** for GCN and **78.18%** for SAGE), surpassing all partial configurations. Even on smaller and more homogeneous graphs, each component makes a distinct positive contribution, and their combination produces a clear synergistic effect in accuracy and robustness. These findings indicate that, even on smaller and more homogeneous graphs, each component makes a distinct and beneficial contribution to the overall distillation process.

KRD	MP-MLP	CWD-Loss	Cora	CiteSeer	PubMed	Average Accuracy
			Accuracy (Inference(ms))	Accuracy (Inference(ms))	Accuracy (Inference(ms))	
-	-	-	83.76 (2.61)	72.45 (2.26)	80.43 (2.30)	78.88
✗	✗	✗	79.22 (1.49)	70.13 (1.63)	78.78 (1.46)	76.04
✓	✗	✗	84.65 (1.49)	74.17 (1.63)	<b>83.02</b> (1.46)	80.61 4.57↑
✓	✓	✗	<b>84.93</b> (1.29)	<b>75.18</b> (1.37)	<b>82.67</b> (1.32)	<b>80.93</b> 4.89↑
✓	✓	✓	<b>85.02</b> (1.29)	<b>75.37</b> (1.37)	82.09 (1.32)	<b>80.83</b> 4.79↑

**Table 5.17: Ablation study of KRD, MP-MLP, and CWD-loss components in the MP-KRD module on classic datasets using the GCN model.** The table reports test accuracy on Cora, CiteSeer, and PubMed for different component combinations. Avg Accuracy indicates the average accuracy and relative improvement. The best, and second-best results are highlighted in **green**, and **red**

KRD	MP-MLP	CWD-Loss	Cora	CiteSeer	PubMed	Average Accuracy
			Accuracy (Inference(ms))	Accuracy (Inference(ms))	Accuracy (Inference(ms))	
-	-	-	82.25 (1.90)	69.18 (3.00)	76.85 (4.55)	76.09
✗	✗	✗	73.74 (1.57)	62.05 (1.73)	73.58 (2.47)	69.79
✓	✗	✗	80.96 (1.57)	70.39 (1.73)	77.59 (2.47)	76.31 6.52↑
✓	✓	✗	<b>81.27</b> (1.38)	<b>72.64</b> (1.82)	<b>78.37</b> (2.30)	<b>77.43</b> 7.64↑
✓	✓	✓	<b>83.16</b> (1.38)	<b>72.78</b> (1.82)	<b>78.59</b> (2.30)	<b>78.18</b> 8.39↑

**Table 5.18: Ablation study of KRD, MP-MLP, and CWD-loss components in the MP-KRD module on classic datasets using the SAGE model.** The table reports test accuracy on Cora, CiteSeer, and PubMed for different component combinations. Avg Accuracy indicates the average accuracy and relative improvement. The best, and second-best results are highlighted in **green**, and **red**

### For Medium-scale datasets

Tables 5.19 and 5.20 report the results on Computers, Photo, WikiCS, CS, and Physics.

KRD	MP-MLP	CWD-Loss	Computers	Photo	WikiCS	CS	Physics	Average Accuracy
			Accuracy (Inference(ms))	Accuracy (Inference(ms))	Accuracy (Inference(ms))	Accuracy (Inference(ms))	Accuracy (Inference(ms))	
-	-	-	93.73 (9.95)	95.78 (7.80)	80.52 (6.40)	96.05 (7.94)	97.44 (8.65)	92.70
✗	✗	✗	91.30 (1.88)	95.21 (2.24)	79.62 (1.67)	95.97 (6.17)	<b>97.25 (7.55)</b>	91.82
✓	✗	✗	<b>91.78 (1.88)</b>	95.64 (2.24)	<b>80.58 (1.67)</b>	96.12 (6.17)	96.95 (7.55)	<b>92.21</b> 0.39↑
✓	✓	✗	91.72 (1.65)	<b>95.71 (2.43)</b>	79.74 (1.55)	<b>96.32 (6.34)</b>	97.16 (7.29)	92.13 0.31↑
✓	✓	✓	<b>93.00 (1.65)</b>	<b>96.23 (2.43)</b>	<b>81.11 (1.55)</b>	<b>96.14 (6.34)</b>	<b>97.31 (7.29)</b>	<b>92.76</b> 0.90↑

**Table 5.19: Ablation study of KRD, MP-MLP, and CWD-loss components in the MP-KRD module on medium-scale datasets using the GCN model.** The table reports test accuracy on Computers, Photo, WikiCS, CS, and Physics datasets for different component combinations. Avg Accuracy indicates the average accuracy and relative change. The best, and second-best results are highlighted in **green**, and **red**

For the GCN backbone, the complete MP-KRD setup achieves the highest average accuracy of **92.76%**, outperforming all ablated variants and showing notable gains on datasets such as Photo (+0.61%) and CS (+0.09%) compared to the KRD-only setting.

For the SAGE backbone, the full configuration reaches **91.68%**, slightly surpassing the KRD teacher baseline (**91.51%**) while maintaining comparable inference efficiency.

### Summary

These results confirm that each component—knowledge distillation, MLP-based message passing, and curriculum-weighted loss—contributes positively to the distilled student’s performance. More importantly, their integration not only improves robustness under various dataset characteristics but also enables efficient, graph-agnostic inference without compromising accuracy. Notably, the proposed MP-KRD design allows the student model to achieve performance that is consistently close to, and in several cases surpasses, that of its GNN teacher, demonstrating its strong potential for real-world deployment where both efficiency and predictive quality are critical.

KRD	MP-MLP	CWD-Loss	Computers	Photo	WikiCS	CS	Physics	Average Accuracy
			Accuracy (Inference(ms))	Accuracy (Inference(ms))	Accuracy (Inference(ms))	Accuracy (Inference(ms))	Accuracy (Inference(ms))	
-	-	-	92.95 (7.94)	95.80 (4.82)	80.22 (3.59)	95.84(20.51)	96.81(50.77)	92.32
✗	✗	✗	85.18 (2.39)	93.37 (2.42)	77.28 (1.42)	93.68 (6.17)	96.50(10.15)	89.20
✓	✗	✗	<b>91.37</b> (2.39)	<b>95.10</b> (2.42)	<b>80.18</b> (1.42)	93.75 (6.17)	<b>97.14</b> (10.15)	<b>91.51</b> 2.31↑
✓	✓	✗	87.07 (2.39)	<b>94.27</b> (2.12)	<b>80.52</b> (1.69)	<b>96.46</b> (6.34)	<b>96.99</b> (8.49)	91.06 1.86↑
✓	✓	✓	<b>90.89</b> (2.39)	94.92 (2.12)	79.60 (1.69)	<b>96.44</b> (6.34)	96.50 (8.49)	<b>91.68</b> 2.48↑

**Table 5.20: Ablation study of KRD, MP-MLP, and CWD-loss components in the MP-KRD module on medium-scale datasets using the SAGE model.** The table reports test accuracy on Computers, Photo, WikiCS, CS, and Physics datasets for different component combinations. Avg Accuracy indicates the average accuracy and relative change. The best, and second-best results are highlighted in **green**, and **red**

## Chapter 6

### Discussion and Conclusion

#### 6.1 Summary of Findings and Contributions

This thesis addressed three interconnected challenges that hinder the practical deployment of Graph Neural Networks: incomplete data, high training costs, and inefficient inference. We argued that isolated solutions are insufficient and proposed that a **holistic, synergistic framework** is necessary to overcome these obstacles simultaneously.

To this end, we presented **iPaD**, a unified framework that seamlessly integrates the entire GNN pipeline. The journey begins with **APCFI**, which robustly and efficiently imputes missing features. It then proceeds to **MPP**, a novel proxy-pruning paradigm that creates a lightweight yet powerful **TunedGNN** teacher and a structurally-aligned MLP student before expensive training commences. Finally, **MP-KRD** completes the pipeline by efficiently distilling the teacher’s knowledge into the student, producing a final model optimized for real-world deployment.

Our extensive experiments have validated the effectiveness of this synergistic design, leading to several key findings:

- **APCFI successfully resolves the accuracy-efficiency trade-off**, achieving imputation quality comparable to state-of-the-art methods like PCFI but with a more than 200x speedup in computation.
- **MPP provides a clear “efficiency sweet spot”**, capable of reducing training and inference costs by over 3x and 5x respectively, with negligible impact on model accuracy at moderate prune rates.
- Most significantly, our framework exhibits nuanced robustness properties. The distilled student model can **outperform its GNN teacher under severe edge missingness**, showcasing resilience to structural corruption. Conversely, its performance is **capped by the teacher under feature missingness**, revealing a fundamental boundary condition of knowledge distillation for GNNs.

These contributions collectively demonstrate that a holistic design philosophy can lead to GNN systems that are at once robust, efficient, and high-performing.

## 6.2 Limitations

Despite its demonstrated strengths, this work has several limitations that open avenues for future research:

- **Applicability to Graph Types:**

The current designs of APCFI and MPP are primarily tailored for homogeneous graphs. Their effectiveness and generalization capability on highly heterogeneous graphs or graphs with diverse node and edge types have not yet been comprehensively validated.

- **Model and Task Generalization:**

While iPaD shows strong results on node classification, its applicability to other critical graph tasks, such as graph classification and link prediction, remains to be explored. Furthermore, its performance on industrial-scale graph datasets with billions of nodes and edges has not been tested.

- **Automation and Hyperparameters:**

The framework, particularly the MP-KRD module, involves several hyperparameters that currently require manual tuning. A fully adaptive, automated mechanism for hyperparameter selection is still lacking.

- **“TunedGNN” Cost:**

Although MPP effectively mitigates its cost, the underlying **TunedGNN** architecture is inherently more computationally expensive than simpler GNNs, which presents a baseline overhead.

## 6.3 Future Work

The limitations of this study naturally point toward several exciting directions for future research:

- **Unlocking Heterogeneous Graphs:**

A primary goal is to extend the principles of APCFI and MPP to be applicable to heterogeneous graphs, which are common in real-world industrial scenarios.

- **Towards a Fully Automated Pipeline:**

Integrating AutoML techniques, especially for automated hyperparameter optimization (e.g., for the distillation loss weights  $\lambda_1, \lambda_2$ ), would significantly lower the barrier to entry and enhance the framework’s practical usability.

- **Expanding to New Tasks and Scales:**

A crucial next step is to adapt and validate the iPaD pipeline for other graph learning tasks, such as link prediction and graph classification, and to benchmark its perfor-



mance on industrial-scale datasets.

- **Deepening the Understanding of Distillation Robustness:**

Building upon our key finding regarding edge vs. feature missingness, a promising theoretical direction is to further investigate the boundary conditions of knowledge distillation for GNNs, formally characterizing when and why a distilled student can be more robust than its teacher.

## **6.4 Concluding Remarks**

In conclusion, this thesis presented iPaD, a comprehensive solution to enhance the practicality of Graph Neural Networks. More than just a collection of algorithms, iPaD embodies a design philosophy centered on synergy and end-to-end optimization. By demonstrating that it is possible to achieve robustness, training efficiency, and inference speed simultaneously without compromising performance, this work provides a valuable blueprint and a solid foundation for the next generation of GNN systems designed for the complexities of the real world.

## References

- [1] Y. Luo, L. Shi, and X.-M. Wu, “Classic gnns are strong baselines: Reassessing gnns for node classification,” *arXiv preprint arXiv:2406.08993*, 2024.
- [2] L. Waikhom and R. Patgiri, “A survey of graph neural networks in various learning paradigms: methods, applications, and challenges,” *Artificial Intelligence Review*, vol. 56, no. 7, pp. 6295–6364, 2023.
- [3] Z. Wu, S. Pan, F. Chen, G. Long, C. Zhang, and P. S. Yu, “A comprehensive survey on graph neural networks,” *IEEE transactions on neural networks and learning systems*, vol. 32, no. 1, pp. 4–24, 2020.
- [4] B. Khemani, S. Patil, K. Kotecha, and S. Tanwar, “A review of graph neural networks: concepts, architectures, techniques, challenges, datasets, applications, and future directions,” *Journal of Big Data*, vol. 11, no. 1, p. 18, 2024.
- [5] S. Wu, F. Sun, W. Zhang, X. Xie, and B. Cui, “Graph neural networks in recommender systems: a survey,” *ACM Computing Surveys*, vol. 55, no. 5, pp. 1–37, 2022.
- [6] R. Ying, R. He, K. Chen, P. Eksombatchai, W. L. Hamilton, and J. Leskovec, “Graph convolutional neural networks for web-scale recommender systems,” in *Proceedings of the 24th ACM SIGKDD international conference on knowledge discovery & data mining*, pp. 974–983, 2018.
- [7] J. Jiang, C. Wilson, X. Wang, W. Sha, P. Huang, Y. Dai, and B. Y. Zhao, “Understanding latent interactions in online social networks,” *ACM Transactions on the Web (TWEB)*, vol. 7, no. 4, pp. 1–39, 2013.
- [8] H. Zhang, B. Wu, X. Yuan, S. Pan, H. Tong, and J. Pei, “Trustworthy graph neural networks: Aspects, methods, and trends,” *Proceedings of the IEEE*, 2024.
- [9] K. Huang, J. Zhai, Z. Zheng, Y. Yi, and X. Shen, “Understanding and bridging the gaps in current gnn performance optimizations,” in *Proceedings of the 26th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*, pp. 119–132, 2021.
- [10] A. Mohi ud din and S. Qureshi, “A review of challenges and solutions in the design and implementation of deep graph neural networks,” *International Journal of Computers and Applications*, vol. 45, no. 3, pp. 221–230, 2023.
- [11] X. Chen, S. Chen, J. Yao, H. Zheng, Y. Zhang, and I. W. Tsang, “Learning on attribute-missing graphs,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 44, no. 2, pp. 740–757, 2020.
- [12] B. Jiang and Z. Zhang, “Incomplete graph representation and learning via partial graph neural networks,” *arXiv preprint arXiv:2003.10130*, 2020.

- [13] J. Zhou, G. Cui, S. Hu, Z. Zhang, C. Yang, Z. Liu, L. Wang, C. Li, and M. Sun, “Graph neural networks: A review of methods and applications,” *AI open*, vol. 1, pp. 57–81, 2020.
- [14] Z. Huan, Y. Quanming, and T. Weiwei, “Search to aggregate neighborhood for graph neural network,” in *2021 IEEE 37th International Conference on Data Engineering (ICDE)*, pp. 552–563, IEEE, 2021.
- [15] H. Ma, Y. Rong, and J. Huang, “Graph neural networks: Scalability,” *Graph Neural Networks: Foundations, Frontiers, and Applications*, pp. 99–119, 2022.
- [16] C. Gao, Y. Zheng, N. Li, Y. Li, Y. Qin, J. Piao, Y. Quan, J. Chang, D. Jin, X. He, *et al.*, “A survey of graph neural networks for recommender systems: Challenges, methods, and directions,” *ACM Transactions on Recommender Systems*, vol. 1, no. 1, pp. 1–51, 2023.
- [17] X. Zheng, Y. Wang, Y. Liu, M. Li, M. Zhang, D. Jin, P. S. Yu, and S. Pan, “Graph neural networks for graphs with heterophily: A survey,” *arXiv preprint arXiv:2202.07082*, 2022.
- [18] J. You, X. Ma, Y. Ding, M. J. Kochenderfer, and J. Leskovec, “Handling missing data with graph representation learning,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 19075–19087, 2020.
- [19] E. Rossi, H. Kenlay, M. I. Gorinova, B. P. Chamberlain, X. Dong, and M. M. Bronstein, “On the unreasonable effectiveness of feature propagation in learning on graphs with missing node features,” in *Learning on graphs conference*, pp. 11–1, PMLR, 2022.
- [20] H. Taguchi, X. Liu, and T. Murata, “Graph convolutional networks for graphs containing missing features,” *Future Generation Computer Systems*, vol. 117, pp. 155–168, 2021.
- [21] D. Um, J. Park, S. Park, and J. young Choi, “Confidence-based feature imputation for graphs with partially known features,” in *The Eleventh International Conference on Learning Representations*, 2023.
- [22] X. Zhu and Z. Ghahramani, “Learning from labeled and unlabeled data with label propagation,” *ProQuest number: information to all users*, 2002.
- [23] H. Zhou, A. Srivastava, H. Zeng, R. Kannan, and V. Prasanna, “Accelerating large scale real-time gnn inference using channel pruning,” *arXiv preprint arXiv:2105.04528*, 2021.
- [24] D. Gurevin, M. Shan, S. Huang, M. A. Hasan, C. Ding, and O. Khan, “Prunegnn: Algorithm-architecture pruning framework for graph neural network acceleration,” in *2024 IEEE International Symposium on High-Performance Computer Architecture (HPCA)*, pp. 108–123, IEEE, 2024.
- [25] T. Chen, Y. Sui, X. Chen, A. Zhang, and Z. Wang, “A unified lottery ticket hypothesis for graph neural networks,” in *International conference on machine learning*, pp. 1695–1706, PMLR, 2021.

- [26] Y.-D. Sui, X. Wang, T. Chen, M. Wang, X.-N. He, and T.-S. Chua, “Inductive lottery ticket learning for graph neural networks,” *Journal of Computer Science and Technology*, vol. 39, no. 6, pp. 1223–1237, 2024.
- [27] Y. Yang, J. Qiu, M. Song, D. Tao, and X. Wang, “Distilling knowledge from graph convolutional networks,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 7074–7083, 2020.
- [28] B. Yan, C. Wang, G. Guo, and Y. Lou, “Tinygnn: Learning efficient graph neural networks,” in *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pp. 1848–1856, 2020.
- [29] S. Zhang, Y. Liu, Y. Sun, and N. Shah, “Graph-less neural networks: Teaching old mlps new tricks via distillation,” *arXiv preprint arXiv:2110.08727*, 2021.
- [30] L. Wu, H. Lin, Y. Huang, , and S. Z. Li, “Quantifying the knowledge in gnns for reliable distillation into mlps,” in *International Conference on Machine Learning*, PMLR, 2023.
- [31] Y. Chen, Y. Bian, X. Xiao, Y. Rong, T. Xu, and J. Huang, “On self-distilling graph neural network,” *arXiv preprint arXiv:2011.02255*, 2020.
- [32] L. Rampásek, M. Galkin, V. P. Dwivedi, A. T. Luu, G. Wolf, and D. Beaini, “Recipe for a general, powerful, scalable graph transformer,” *Advances in Neural Information Processing Systems*, vol. 35, pp. 14501–14515, 2022.
- [33] J. Chen, K. Gao, G. Li, and K. He, “Nagphormer: A tokenized graph transformer for node classification in large graphs,” *arXiv preprint arXiv:2206.04910*, 2022.
- [34] H. Shirzad, A. Velingker, B. Venkatachalam, D. J. Sutherland, and A. K. Sinop, “Exphormer: Sparse transformers for graphs,” in *International Conference on Machine Learning*, pp. 31613–31632, PMLR, 2023.
- [35] K. Kong, J. Chen, J. Kirchenbauer, R. Ni, C. B. Bruss, and T. Goldstein, “Goat: A global transformer on large-scale graphs,” in *International Conference on Machine Learning*, pp. 17375–17390, PMLR, 2023.
- [36] Q. Wu, W. Zhao, Z. Li, D. P. Wipf, and J. Yan, “Nodeformer: A scalable graph structure learning transformer for node classification,” *Advances in Neural Information Processing Systems*, vol. 35, pp. 27387–27401, 2022.
- [37] Q. Wu, W. Zhao, C. Yang, H. Zhang, F. Nie, H. Jiang, Y. Bian, and J. Yan, “Sgformer: Simplifying and empowering transformers for large-graph representations,” *Advances in Neural Information Processing Systems*, vol. 36, pp. 64753–64773, 2023.
- [38] C. Deng, Z. Yue, and Z. Zhang, “Polynormer: Polynomial-expressive graph transformer in linear time,” *arXiv preprint arXiv:2403.01232*, 2024.
- [39] K. Sharma, Y.-C. Lee, S. Nambi, A. Salian, S. Shah, S.-W. Kim, and S. Kumar, “A survey of graph neural networks for social recommender systems,” *ACM Computing Surveys*, vol. 56, no. 10, pp. 1–34, 2024.
- [40] X.-M. Zhang, L. Liang, L. Liu, and M.-J. Tang, “Graph neural networks and their current applications in bioinformatics,” *Frontiers in genetics*, vol. 12, p. 690049, 2021.

- [41] H. Zhang, B. Wu, X. Yuan, S. Pan, H. Tong, and J. Pei, “Trustworthy graph neural networks: Aspects, methods, and trends,” *Proceedings of the IEEE*, vol. 112, no. 2, pp. 97–139, 2024.
- [42] Z. Wu, S. Pan, F. Chen, G. Long, C. Zhang, and P. S. Yu, “A comprehensive survey on graph neural networks,” *IEEE Transactions on Neural Networks and Learning Systems*, vol. 32, no. 1, pp. 4–24, 2021.
- [43] X. Jiang, Z. Tian, and K. Li, “A graph-based approach for missing sensor data imputation,” *IEEE Sensors Journal*, vol. 21, no. 20, pp. 23133–23144, 2021.
- [44] D. Adhikari, W. Jiang, J. Zhan, Z. He, D. B. Rawat, U. Aickelin, and H. A. Khorshidi, “A comprehensive survey on imputation of missing data in internet of things,” *ACM Computing Surveys*, vol. 55, no. 7, pp. 1–38, 2022.
- [45] Z. Lin, C. Li, Y. Miao, Y. Liu, and Y. Xu, “Pagraph: Scaling gnn training on large graphs via computation-aware caching,” in *Proceedings of the 11th ACM Symposium on Cloud Computing*, pp. 401–415, 2020.
- [46] D. Zheng, C. Ma, M. Wang, J. Zhou, Q. Su, X. Song, Q. Gan, Z. Zhang, and G. Karypis, “Distdgl: Distributed graph neural network training for billion-scale graphs,” in *2020 IEEE/ACM 10th Workshop on Irregular Applications: Architectures and Algorithms (IA3)*, pp. 36–44, IEEE, 2020.
- [47] Z. Xue, Y. Yang, and R. Marculescu, “Sugar: Efficient subgraph-level training via resource-aware graph partitioning,” *IEEE Transactions on Computers*, vol. 72, no. 11, pp. 3167–3177, 2023.
- [48] W. Lu, Z. Guan, W. Zhao, and Y. Yang, “Adagmlp: Adaboosting gnn-to-mlp knowledge distillation,” in *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pp. 2060–2071, 2024.
- [49] X. Han, T. Zhao, Y. Liu, X. Hu, and N. Shah, “MLPInit: Embarrassingly simple GNN training acceleration with MLP initialization,” in *International Conference on Learning Representations*, 2023.
- [50] G. Fang, X. Ma, M. Song, M. B. Mi, and X. Wang, “Depgraph: Towards any structural pruning,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 16091–16101, 2023.
- [51] J. Lee, S. Park, S. Mo, S. Ahn, and J. Shin, “Layer-adaptive sparsity for the magnitude-based pruning,” *arXiv preprint arXiv:2010.07611*, 2020.
- [52] C. Shu, Y. Liu, J. Gao, Y. Zheng, and C. Shen, “Channel-wise knowledge distillation for dense prediction,” in *ICCV*, 2021.
- [53] P. Sen, G. Namata, M. Bilgic, L. Getoor, B. Galligher, and T. Eliassi-Rad, “Collective classification in network data,” *AI magazine*, vol. 29, no. 3, pp. 93–93, 2008.
- [54] O. Shchur, M. Mumme, A. Bojchevski, and S. Günnemann, “Pitfalls of graph neural network evaluation,” *arXiv preprint arXiv:1811.05868*, 2018.

- [55] P. Mernyei and C. Cangea, “Wiki-cs: A wikipedia-based benchmark for graph neural networks,” *arXiv preprint arXiv:2007.02901*, 2020.
- [56] F. Monti, M. Bronstein, and X. Bresson, “Geometric matrix completion with recurrent multi-graph neural networks,” *Advances in neural information processing systems*, vol. 30, 2017.
- [57] T. Akiba, S. Sano, T. Yanase, T. Ohta, and M. Koyama, “Optuna: A next-generation hyperparameter optimization framework,” in *The 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pp. 2623–2631, 2019.

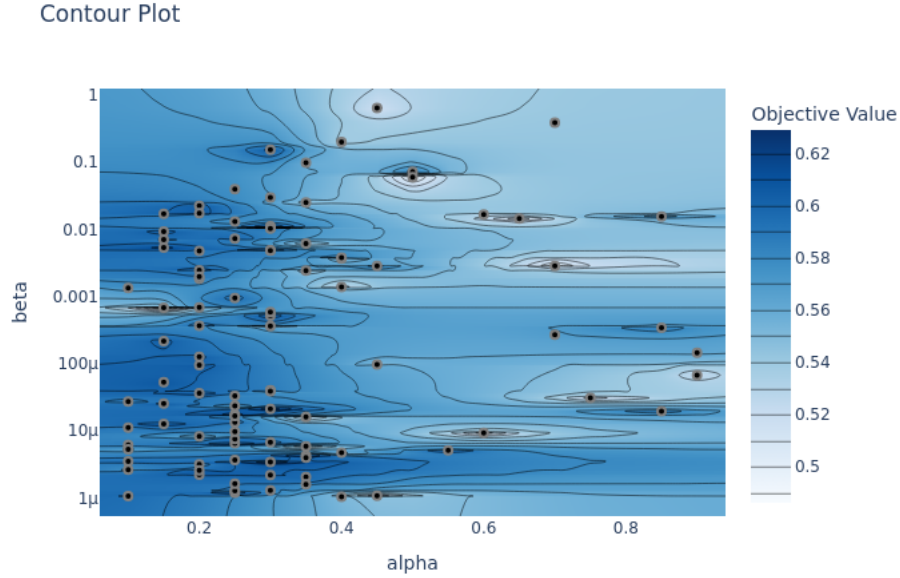
# Appendix A

## Hyperparameters

This appendix provides a detailed analysis of the key hyperparameters for our proposed modules, justifying the settings used in our main experiments. To ensure a systematic and unbiased selection, we utilized the **Optuna** [57] to perform automated hyperparameter optimization (HPO). For each major experiment, we ran over 100 trials, optimizing for the validation set accuracy. The following sections present the insights gained from this process.

### A.1 Analysis of APCFI Hyperparameters

For our **APCFI** module, we analyzed the two most sensitive hyperparameters: the confidence decay factor  $\alpha$  and the inter-channel propagation strength  $\beta$ . The contour plot in Figure A.1 visualizes their complex interaction.



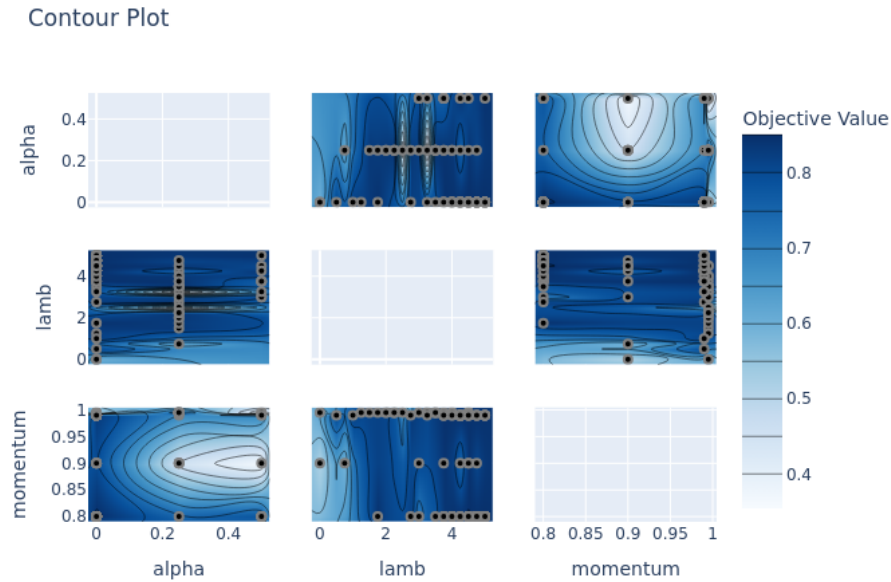
**Figure A.1:** Contour plot for the key hyperparameters ( $\alpha$  and  $\beta$ ) of the APCFI module on the Cora dataset.

The key insight from this plot is the **strong interaction between  $\alpha$  and  $\beta$** . The optimal value for one parameter is highly dependent on the value of the other, leading to several

**islands of high performance.** For instance, when **beta** is very low (around  $10^{-5}$ ), a lower **alpha** (around 0.2-0.4) can achieve good results. However, the largest and most stable high-performance region occurs when **alpha is in a mid-to-high range (0.5 to 0.8)** and **beta is in a low-to-mid range (around  $10^{-3}$ )**. This analysis demonstrates that while the module is sensitive to this parameter combination, its behavior is predictable and optimizable, and it provides a clear guideline for selecting a robust set of default parameters.

## A.2 Analysis of MP-KRD Hyperparameters

For our **MP-KRD** module, we identified **alpha**, **lamb**, and **momentum** as three critical hyperparameters influencing the distillation curriculum and training dynamics. Figure A.2 shows the contour plot illustrating the pairwise interactions of these parameters on the final model accuracy.



**Figure A.2:** Contour plot for the key hyperparameters of the MP-KRD module. Darker blue regions indicate higher accuracy. The plot reveals strong dependencies between parameters.

The analysis reveals several clear trends. The most decisive parameter is **momentum**, where higher values (approximately 0.85 to 1.0) consistently correlate with superior performance. Conversely, the **lamb** parameter performs best in a lower range (approximately 0 to 2.5). The impact of **alpha** is more subtle and interacts with the other two. Overall, the plot indicates that the optimal performance **sweet spot** is achieved with a combination of **high momentum, low lamb, and mid-to-low alpha**. This data-driven insight guided our final parameter selection for the MP-KRD experiments.



## Appendix B

### Generalization Ability

#### B.1 Generalization of iPaD to Various GNN Architectures

The proposed **iPaD** framework is inherently designed to be model-agnostic, enabling its integration with a wide range of graph neural network (GNN) architectures, including both classical message-passing models and recent state-of-the-art (SOTA) designs. Experimental evaluations (Table B.1) confirm that iPaD can be successfully applied to diverse backbones such as GCN, GraphSAGE, Polynormer, and SGFormer, consistently delivering substantial inference speed improvements.

However, the performance gains in terms of accuracy are less pronounced when applied to attention-based architectures (e.g., SGFormer, Polynormer). This limitation may stem from the intrinsic characteristics of attention mechanisms, where the knowledge distilled from the teacher is encoded in highly adaptive, pairwise attention weights. The current iPaD distillation design primarily focuses on message-passing structures and may not effectively capture or transfer the fine-grained relational patterns inherent to attention-based models. Despite this, the efficiency enhancement remains significant, demonstrating the framework’s potential for deployment in scenarios where inference speed is prioritized over marginal accuracy improvements.

<b>Model</b>	<b>Computers</b> Accuracy (Inference(ms))	<b>Photo</b> Accuracy (Inference(ms))	<b>WikiCS</b> Accuracy (Inference(ms))	<b>CS</b> Accuracy (Inference(ms))	<b>Average</b> Accuracy
SGFormer	92.00 (5.95)	96.34 (4.42)	80.05 (12.20)	95.31 (12.69)	90.93
<b>iPaD-SGFormer(Our)</b>	<b>66.99</b> (3.08)	<b>67.27</b> (2.99)	<b>76.06</b> (3.95)	<b>93.76</b> (6.88)	<b>76.02</b>
Polynormer	93.00 (37.94)	96.31 (27.61)	79.66 (43.80)	94.81 (30.84)	90.95
<b>iPaD-Polynormer(Our)</b>	<b>87.41</b> (4.83)	<b>93.37</b> (2.91)	<b>75.53</b> (3.67)	<b>95.93</b> (4.99)	<b>88.06</b>
TunedGCN	93.64 (17.64)	95.97 (12.69)	80.43 (8.29)	95.89 (16.55)	91.48
<b>iPaD-GCN(Our)</b>	<b>93.00</b> (1.53)	<b>96.23</b> (1.65)	<b>81.11</b> (1.50)	<b>96.14</b> (5.45)	<b>91.62</b>
TunedSAGE	93.14 (7.77)	96.04 (5.00)	80.01 (4.27)	95.94 (25.08)	91.28
<b>iPaD-SAGE(Our)</b>	<b>91.07</b> (1.62)	<b>95.33</b> (1.92)	<b>79.67</b> (1.26)	<b>96.46</b> (7.04)	<b>90.63</b>

**Table B.1: Efficiency improvement of the iPaD framework over various SOTA models across four datasets.** For each dataset (Computers, Photo, WikiCS, CS), accuracy (%) is reported alongside inference time (ms) in parentheses.

## Appendix C

### Effect of MPP Pruning on Teacher Models

#### C.1 Rationale for Applying MPP Pruning to Teacher Models

To improve the efficiency of the distillation process, the teacher models are first pruned using the proposed Multi-Path Proxy (MPP) strategy prior to knowledge transfer. While parameter pruning inherently reduces the model size, an appropriately configured MPP preserves the most informative channels and removes redundant computations, thereby maintaining the predictive capability of the teacher. As shown in Table C.1, the pruned teacher models exhibit negligible or no accuracy degradation compared to their unpruned counterparts, while delivering substantial improvements in training speed (e.g., from 19.35 to 54.49 epoch/s for TunedGCN). Furthermore, the MPP strategy is not limited to a specific backbone and can be effectively extended to various classic or state-of-the-art (SOTA) architectures, such as SGFormer, Polynormer, GCN, and GraphSAGE, consistently providing both training efficiency gains and competi

Model	Computers Accuracy (Train(epcoh/s))	Photo Accuracy (Train(epcoh/s))	WikiCS Accuracy (Train(epcoh/s))	CS Accuracy (Train(epcoh/s))	Average Accuracy
SGFormer	92.00 (58.05)	96.34 (75.48)	80.05 (27.22)	95.31 (25.57)	90.93
<b>MPP-SGFormer(Our)</b>	<b>91.87</b> (73.00)	<b>95.46</b> (81.12)	<b>78.81</b> (45.99)	<b>95.15</b> (39.81)	<b>90.32</b>
Polynormer	93.00 (6.90)	96.31 (9.54)	79.66 (6.19)	94.81 (8.72)	90.95
<b>MPP-Polynormer(Our)</b>	<b>93.22</b> (15.94)	<b>96.33</b> (19.10)	<b>79.79</b> (14.58)	<b>93.08</b> (19.34)	<b>90.61</b>
TunedGCN	93.64 (19.35)	95.97 (28.84)	80.43 (41.64)	95.89 (23.32)	91.48
<b>MPP-TunedGCN(Our)</b>	<b>93.62</b> (54.49)	<b>95.92</b> (47.20)	<b>80.68</b> (69.58)	<b>96.01</b> (47.73)	<b>91.56</b>
TunedSAGE	93.14 (55.97)	96.04 (73.37)	80.01 (89.41)	95.94 (15.71)	91.28
<b>MPP-TunedSAGE(Our)</b>	<b>92.87</b> (90.49)	<b>95.83</b> (90.49)	<b>80.27</b> (113.37)	<b>95.87</b> (21.29)	<b>91.21</b>

**Table C.1: Training efficiency improvement of MPP pruning over various SOTA models.** For each dataset (Computers, Photo, WikiCS, CS), accuracy (%) is reported alongside training speed (epoch/s) in parentheses.

## Appendix D

### Joint Training Imputer

#### D.1 Joint Training of FP-based Imputer with MPP and KRD

The APCFI module in the proposed framework is currently applied as an independent pre-processing step, separate from the subsequent MPP pruning and MP-KRD distillation stages. A potential extension is to jointly train APCFI alongside these modules, enabling the feature imputation process to be directly optimized for the downstream node classification task. This could, in principle, yield more task-adaptive imputation results.

However, experimental attempts at such joint training were unsuccessful, as shown in Table D.1. The FP-based imputation methods employed in APCFI are inherently non-differentiable, which restricts gradient-based optimization of their internal parameters. To enable end-to-end training, an additional learnable neural network layer was appended after the imputed feature matrix. In practice, the absence of a well-defined loss function to guide this learnable module towards producing more informative feature representations led to performance degradation. Instead of enhancing the quality of the imputed features, the learnable layer often disrupted the inherent structural and statistical properties of the imputation output, thereby impairing the downstream model’s training.

Method	Cora	CiteSeer	Pubmed	Computer	Photo	WikiCS	Average $\uparrow$
FP	78.44	59.36	79.08	92.33	93.41	78.15	80.13
trainable FP	11.72	20.72	42.80	39.03	37.23	22.38	28.98

**Table D.1: Performance comparison of FP and trainable FP under high missing rates.** Node classification accuracy (%) on six datasets (Cora, CiteSeer, Pubmed, Computer, Photo, WikiCS) with a uniform feature missing rate of 90%.

These findings suggest that the design of a differentiable, loss-guided imputation mechanism remains an open challenge. Developing such methods could be a promising direction for future work, potentially enabling effective end-to-end optimization of the APCFI module together with MPP and KRD.

## Appendix E

### Complex Missing Conditions

#### E.1 Performance under Combined Feature and Edge Missing Conditions

To further assess the robustness of the proposed **iPaD** framework, an additional experiment was conducted under complex extreme missing data scenarios, where both node features and graph structure were heavily corrupted. Specifically, the evaluation considered a uniform feature missing rate of 99.5% in conjunction with an edge missing rate of 60%, simulating an exceptionally challenging setting for graph learning. Results are summarized in Table E.1.

Across four benchmark datasets (Computers, Photo, WikiCS, Physics), **iPaD**-enhanced models were compared against their corresponding SOTA counterparts augmented with standard feature propagation (FP) preprocessing. While the accuracy gap between **iPaD** and FP-based baselines was reduced under such severe dual corruption, *iPaD* still maintained competitive performance levels, particularly for the GCN backbone, which achieved an average accuracy of 78.43% compared to 79.41% for its FP-based variant.

It is worth noting that, unlike the single **edge missing** scenario reported in Section V.I where the **iPaD** architecture consistently outperformed TunedGNN, such a performance advantage was not observed under the combined corruption setting. The primary reason lies in the fact that **edge missing** also degrades the performance of the imputer module, which in turn limits the student model’s ability to identify and exploit reliable feature signals to replace suboptimal teacher model information. As a result, the overall performance becomes more dependent on the teacher model’s reliability, leading to a reduced advantage for the distilled student.

These findings suggest that, although extreme simultaneous feature and edge removal significantly challenges all evaluated methods, **iPaD** retains a comparable level of predictive capability while preserving its inherent efficiency advantages. This reinforces the framework’s applicability in scenarios where graph data is subject to multi-faceted degradation.

<b>Model</b>	<b>Computer</b>	<b>Photo</b>	<b>WikiCS</b>	<b>Physics</b>	<b>Average Accuracy</b>
SGFormer w FP	$81.82 \pm 0.81$	$86.81 \pm 0.62$	$63.61 \pm 0.33$	$83.01 \pm 0.67$	78.13
Polynormer w FP	$80.67 \pm 4.66$	$85.01 \pm 2.21$	$61.38 \pm 0.58$	$77.82 \pm 4.49$	76.22
TunedGCN w FP	$83.54 \pm 0.49$	$88.44 \pm 0.49$	$62.06 \pm 0.27$	$83.59 \pm 0.13$	79.41
<b>iPaD-GCN(Our)</b>	$82.88 \pm 0.76$	$83.69 \pm 0.41$	$63.43 \pm 0.38$	$83.61 \pm 0.10$	78.43
TunedSAGE w FP	$85.74 \pm 0.58$	$88.46 \pm 0.37$	$63.00 \pm 0.34$	$82.75 \pm 0.08$	80.00
<b>iPaD-SAGE(Our)</b>	$80.47 \pm 0.12$	$83.69 \pm 0.41$	$63.43 \pm 0.38$	$82.81 \pm 0.72$	77.61

**Table E.1: Performance under complex extreme missing data conditions.** Node classification accuracy (%) on four datasets (Computer, Photo, WikiCS, Physics) with uniform feature missing rate of 99.5% and edge missing rate of 60%. Shaded rows indicate our iPaD variants, compared against SOTA models with feature propagation (FP) preprocessing.