



Accelerator for LLM-Enhanced GNN with Product Quantization and Unified Indexing

<u>Jiaming Xu*1,2</u>, Jinhao Li*1, Jun Liu1, Hao Zhou1, Guohao Dai^{1,2} Equal contribution*, Shanghai Jiao Tong University1, Infinigence-Al², Correspondence to: Guohao Dai <<u>daiguohao@sjtu.edu.cn</u>>

Outline

Background and Challenge

- ≻Techniques
 - Overview
 - PQ-based MatMul
 - Unified Architecture
 - Extensible GFM-ISA Design

Experiment Results

Graph Neural Network for Graph Data

Graph neural networks (GNNs) are designed for graph data problem.



Problem of Traditional GNN

Vulnerability of traditional GNN on unseen graphs.

[1] Zheng X, Zhang M, Chen C, et al. Gnnevaluator: Evaluating gnn performance on unseen graphs without labels[J]. Advances in Neural Information Processing Systems, 2024, 36.

Era of Large Language Models (LLMs)

LLMs have attracted much attention.

[1] Qin L, Chen Q, Zhou Y, et al. Multilingual large language model: A survey of resources, taxonomy and frontiers[J]. arXiv preprint arXiv:2404.04925, 2024.

[2] https://www.21jingji.com/article/20240516/herald/c7c069827f1f79b182604827b494511b.html

[3] Zhao W X, Zhou K, Li J, et al. A survey of large language models[J]. arXiv preprint arXiv:2303.18223, 2023.

LLM Improves GNN

Utilize the strong generalization of LLM to improve GNN

[1] Qin L, Chen Q, Zhou Y, et al. Multilingual large language model: A survey of resources, taxonomy and frontiers[J]. arXiv preprint arXiv:2404.04925, 2024.

Related Works

Combination Methods of LLM and GNN^[1]

[1] Li Y, Li Z, Wang P, et al. A survey of graph meets large language model: Progress and future directions[J]. arXiv preprint arXiv:2311.12399, 2023.

[2] Liu H, Feng J, Kong L, et al. One for all: Towards training one graph model for all classification tasks[J]. arXiv preprint arXiv:2310.00149, 2023.

[3] Tang J, Yang Y, Wei W, et al. Graphgpt: Graph instruction tuning for large language models[C]//Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval. 2024: 491-500.

[4] Zhao J, Qu M, Li C, et al. Learning on large-scale text-attributed graphs via variational inference[J]. arXiv preprint arXiv:2210.14709, 2022.

Related Works

[4] Zhao J, Qu M, Li C, et al. Learning on large-scale text-attributed graphs via variational inference[J]. arXiv preprint arXiv:2210.14709, 2022.

LLM-Enhanced GNN Dataflow

LLM-Enhanced GNN Dataflow

Challenge

Outline

Background and Challenge

≻Techniques

- Overview
- PQ-based MatMul
- Unified Architecture
- Extensible GFM-ISA Design

Experiment Results

Techniques Overview

Motivation The intensive GEMMs in LLM need to be alleviated while ensuring the accuracy for end-to-end acceleration.

Approach **Product Quantization** is used to approximate the linear operation in LLMs (PQ-based MatMul).

(b) Online Inference

Reduce >70% overall computation workload

Technique 2: Unified Architecture

Motivation The two types of memory access in the LLM-enhanced GNN lower hardware utilization.

Technique 2: Unified Architecture

Approach A hardware architecture with unified indexing unit is designed to support both types of computation.

Technique 3: Extensible GFM-ISA Design

Motivation Existing software frameworks lacking low-level primitives for LLM-enhanced GNNs with PQ-based MatMul.

Technique 3: Extensible GFM-ISA Design

Approach **Extensible GFM-ISA design** is used to offer low-level software primitives based on the unified architecture.

Outline

Background and Challenge

- ≻Techniques
 - Overview
 - PQ-based MatMul
 - Unified Architecture
 - Extensible GFM-ISA Design

Experiment Results

Experiment Setup

➤Tools for Evaluation: Ramulator 2.0 and TSMC 28nm process library

Benchmark Models:

- LLMs: BERT, Sentence Transformer, E5-large-v2, Llama2-7B/13B
- GNN: R-GCN

► Baselines: NVIDIA A100 GPU, SGCN, MEGA, FACT

Datasets

Dataset	Task	Nodes	Avg. Edges				
Cora (CR)	Node/Link	2,708	10,556				
PubMed (PM)	Node/Link	19,717	44,338				
ogbn-arxiv (AX)	Node	169,343	1,166,243				
ogbn-products (PR)	Node	2,449,029	61,859,140				
Wiki-CS (WK)	Node	11,701	216,123				
FB15K237 (FB)	Link	14,541	310,116				
WN18RR (WN)	Link	40,943	93,003				

Table 1: Detailed Information of Datasets

		Freq.	Compute Unit	On-chip Memory	Bandwidth
GNN.	SGCN [27]	1GHz	32x32 SA	512KB	256GB/s HBM2
	MEGA [28]	1GHz	4x8x32 BSEs	392KB	256GB/s HBM1
LLM	FACT [18]	0.5GHz	16x32 SA	256KB	256GB/s HBM2
	GFMEngine	1GHz	16PEs (each with 4x16 SA 2x8 AT)	256KB	256GB/s HBM2

Table 2: System Configurations

Accuracy Evaluation

	BERT*	ST*	E5*	Llama2-7B*	Llama2-13B*
CR-Link	1.97	-0.81	0.92	0.51	-0.02
CR-Node	0.95	-3.64	-0.13	-0.86	0.58
PM-Link	0.41	0.09	0.19	0.16	0.47
PM-Node	2.39	0.54	1.40	1.37	0.36
AX	0.63	-0.92	0.02	-0.09	0.22
PR	1.19	0.28	0.11	-0.79	1.18
WK	2.62	2.41	0.61	0.31	0.62
FB	0.73	-0.06	0.6	0.33	1.2
WN	0.12	0.12	0.16	0.96	-1.69
Avg. Loss (↓)	1.22	-0.22	0.43	0.21	0.32

Table 3: Accuracy Loss

* represents the accuracy loss: the accuracy of original LLM-enhanced GNNs minus the accuracy of those models with PQ-based MatMul)

The accuracy loss of the LLM-enhanced GNNs with PQ-based MatMul can be negligible in most cases.

Hardware Evaluation

Compared with A100 GPU and other DSA, we achieve up to $3.93 \times$, and $38.66 \times$ speedup and up to $102.52 \times$, $37.82 \times$ energy efficiency.

Accelerator for LLM-Enhanced GNN

with Product Quantization and Unified Indexing

Jiaming Xu, supervised by Prof. Guohao Dai

jiamingxu@sjtu.edu.cn, daiguohao@sjtu.edu.cn

