

GOFA: A Generative One-For-All Model for Joint Graph Language Modeling

Presenter: Zhaokun Wang | July 7, 2025

Content

UNIVERSITĂ HEIDELBERC ZUKUNFT SEIT 1386

- 1. Introduction
- 2. Background
- 3. Methodology
- 4. Experiment
- 5. Conclusion

The Age of Foundation Models

- GPT-4, LLaMA, Claude... We've all heard the buzz.
- These models understand and generate natural language impressively.
- But what about non-text data? Specifically, graphs?

Example

How does ChatGPT handle a social network or a molecular structure?

Introduction • 00000

July 7, 2025

Background 0000000 Methodology 00000000 Experiment 000000000



Why Graphs Are Special

Key Difference

Text is linear. Graphs are not.

Graphs represent diverse structures:

- Social networks (e.g., who follows whom)
- Knowledge bases (e.g., "Paris is capital of France")
- Molecules, traffic networks, file systems...

Structure Matters

Who connects to whom? How? Think: Facebook's friend graph vs. Wikipedia articles.

Introduction 00000

Background

Methodology

Experiment



Why It's Hard for LLMs

- LLMs need linear input (sequences of text).
- Graphs are *permutation invariant* (reordering nodes doesn't change the graph's meaning).
- Flattening graphs to text inevitably loses structural information.

Analogy

It's like trying to explain a complex subway map by just listing every station's name in one long sentence. The crucial connections are lost.

Introduction 00000

July 7, 2025

Background 0000000

Methodology

Experiment 000000000



Motivation for a Graph Foundation Model (GFM)

To be useful, it must:

- Learn from vast amounts of **unlabeled** graphs (i.e., be self-supervised).
- Generalize across diverse downstream tasks (classification, question answering, generation...).
- Understand both node/edge **content** and network **structure**.

The Goal

Not just "what is this node?" but "what is its role in the network?"

Introduction 00000

Background 0000000

Methodology

Experiment 000000000



Existing Approaches (and Their Limitations)

1. LLM as Predictor	2. LLM as Enhancer
\blacksquare Flatten graph \rightarrow feed to LLM.	Use LLM to add features to GNNs.
+ Easy to generate text.	+ Good structure modeling.
- Structure is weakly modeled.	- Cannot handle diverse tasks.

Conclusion

Neither approach can fully serve as a true Graph Foundation Model.

Introduction

Background 0000000 Methodology

Experiment 000000000



Enter GOFA

Generative One-For-All

A New Architecture

A novel model that **interleaves** Graph Neural Networks (GNNs) directly into a Large Language Model (LLM).

- Learns to model graphs in a generative, self-supervised way.
- Maintains the LLM's flexibility and zero-shot capabilities.
- Adds the GNN's strong sense of structure and topology.

Analogy

It's like inserting a specialized subway map reader directly into a text-generating machine.

Introduction 00000

Background

Methodology 0000000 Experiment 000000000



GOFA = GNN + LLM (at token level!)

Architectural Overview

■ Input: Text-Attributed Graph (TAG)

■ GOFA Encoder: LLM + interleaved GNN layers

■ GOFA Decoder: LLM-based text generator

Key Idea

GNN layers inject structure awareness into the LLM's representations.

Introduction

Background

•000000

Methodology

Experiment 000000000



A Unified Framework

Core Idea

Treat all graph tasks as graph completion.

- Just as LLMs complete a sentence: The cat sat on the ____
- GOFA completes a graph:

 Given this graph of papers, a promising research direction is _____

Introduction

July 7, 2025

Background
0 • 0 0 0 0 0 0

Methodology

Experiment 000000000



TAG = Text-Attributed Graph

Each Node / Edge has text:

■ "Node A": Paper title, abstract

■ "Edge A-B": Co-citation relationship

TAG = {Nodes, Edges, Text at Nodes, Text at Edges}

Analogy

Like an annotated map with a narrative at every stop.

Introduction

Background 000000

Methodology

Experiment



Text-Attributed Graph (TAG)

The **Text-Attributed Graph (TAG)** converts everything into natural language.

BEFORE (Chaos)

- Node Q76
- Node with feature [0.9, 0.2]
- Node user_123

AFTER (Clarity with TAG)

- "This is the entity for Barack Obama..."
- "This paper is about AI and has a citation count of 500."
- "This is User 123. Their interests include hiking."

Now, the model only needs to do one thing: read.

Introduction

Background 000 000

Methodology

Experiment 000000000



Node of Generation (NOG)

The **Node of Generation (NOG)** is a temporary node we add to the graph to guide the model.

- 1. Its text is the **user's prompt** ("Find the shortest path...").
- 2. It **connects** to the relevant nodes (A and D).
- 3. The model's job is to generate the answer at the NOG.

Sentence Completion

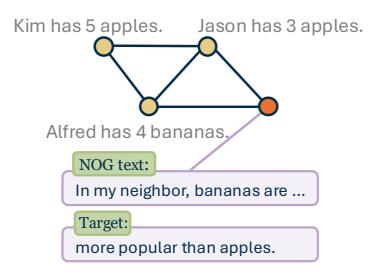


Figure 1: NOG Architecture.

0000000

Zhaokun Wang: GOFA

Methodology

Experiment



Conclusion

Introduction

Background

An Example

1. The Raw Graph

Nodes for "Nolan", "Inception", "Zimmer".

2. The TAG Version

Nodes now have text: "Director: Christopher Nolan", "Composer: Hans Zimmer", etc.

3. The User's Query

"What is the connection?"

4. The NOG is Created & Connected

A new NOG node with the query text appears and links to "Nolan" and "Zimmer" nodes.

5. The Model Generates the Answer at the NOG

"Christopher Nolan worked with composer Hans Zimmer on the film Inception."

Introduction

July 7, 2025

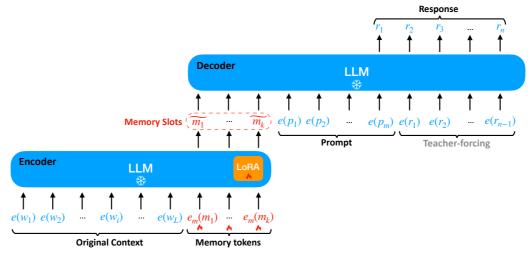
Background 0000000 Methodology

Experiment



In-Context Autoencoder (ICAE)

- What It Is: An autoencoder built with two LLMs (a compressor and a decoder) that compresses sentences into a fixed number of "memory tokens".
- **How It Works**: The decoder reconstructs the original text by only attending to these compressed memory tokens.
- **Purpose in GOFA**: To create dense, fixed-size text representations for graph nodes, which allows the GNN to perform message passing with minimal information loss.



Introduction

Background 0000000 Figure 2: ICAE Architecture.

Experiment



Transformer Convolutional GNN (TransConv)

- Core Design: GOFA uses a custom Transformer Convolutional GNN (TransConv), with its layers interleaved between the LLM's transformer layers.
- **Mechanism**: It functions like a standard Transformer attention layer, where a node's representation is updated from an attention-weighted sum of its neighbors' features.
- **Key Details**: The implementation is stabilized with pre-layer normalization and residual connections.

Introduction

Background 000000

Methodology

Experiment 000000000



The GOFA Architecture

Key Design

A hybrid design that interleaves GNN and frozen LLM layers.

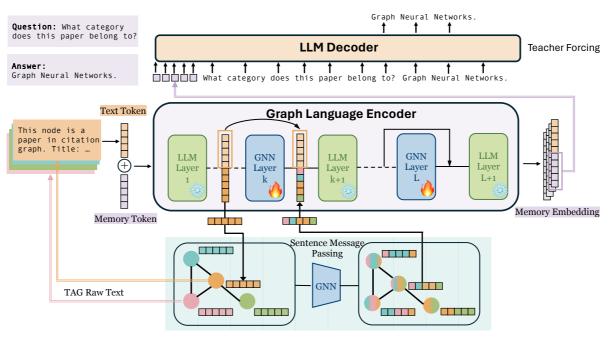


Figure 3: GOFA Architecture.

Introduction

Background

Methodology

• 0 0 0 0 0 0 0

Experiment 000000000



GOFA Encoder = GNN + LLM (Compressor)

How to encode a TAG?

- Sentence compressor (ICAE): converts node text → token-level embeddings.
- GNN layers: exchange messages between token embeddings, not just whole nodes.
- **LLM layers**: perform attention across tokens within a single node.

Information Flow

- Tokens within a node attend to each other
- Tokens from neighboring nodes message-pass via GNN

Introduction

Background

Methodology •••••• Experiment 000000000



Why Token-Level Message Passing?

Naive GNN: Treats node as one big vector

Pooling loses sentence detail.

GOFA's Token-GNN: Passes messages between token positions

- Preserves word-level meaning.
- Enables finer-grained structure reasoning.

Like passing phrases, not just labels, between neighbors.

Introduction

Background

Methodology 0000000 Experiment 000000000



Step 1: Smart Compression with Memory Tokens

Solution: Distill, Don't Discard

GOFA uses a pre-trained sentence compressor (ICAE) to distill a sentence into a fixed set of K Memory Tokens.

■ These are not a single summary; they are *rich*, *multi-faceted "bullet points"* that preserve the original information.

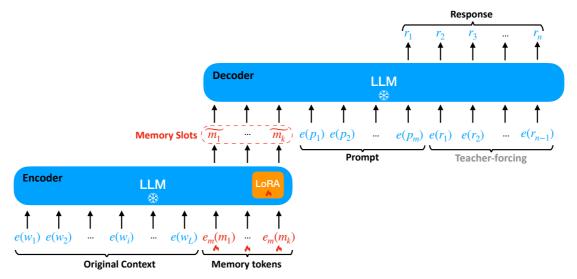


Figure 4: ICAE Architecture.

Introduction

Background

Methodology 00000000 Experiment 000000000



Step 2: Semantic Processing

The Memory Tokens for each node first pass through a standard **LLM (Transformer) layer**.

- At this point, the model processes each node's text in isolation.
- The self-attention mechanism deepens semantic understanding of the node's own content.

Goal

Understand what the node is about *before* considering its connections.

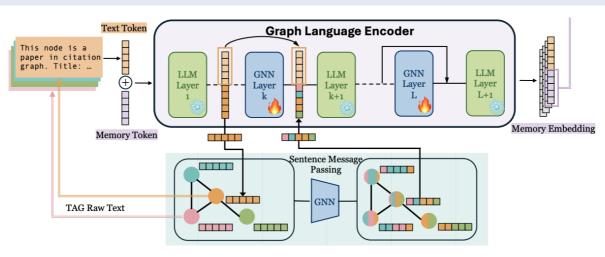


Figure 5: The Encoder.

Introduction

July 7, 2025

Background

Methodology 00000000 Experiment 000000000



Step 3: The GNN Layer

The output immediately flows into a GNN Layer which performs message passing at the token level.

■ The k-th memory token of a node only receives messages from the k-th memory token of its neighbors.

Analogy: Parallel Phone Calls

Line 1 of Node A only talks to Line 1 of its neighbors. Line 2 only talks to Line 2, and so on. No crosstalk... yet.

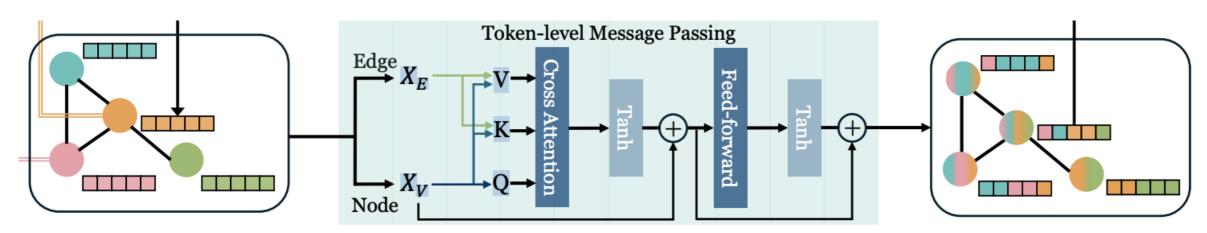


Figure 6: The GNN Layer.

Introduction

Background 0000000

Methodology 00000000 Experiment 000000000



Step 4: The Fusion

The now structurally-aware tokens flow into the **next LLM layer**.

- The LLM's powerful self-attention mechanism now fuses everything.
- The 1st token (which heard from its neighbors) can now interact with the 2nd, 3rd, and all other tokens of its own node.

Analogy: The Conference Room

After the separate phone calls (GNN), everyone gets into a conference room (the LLM layer) to share notes and build a complete, unified picture.

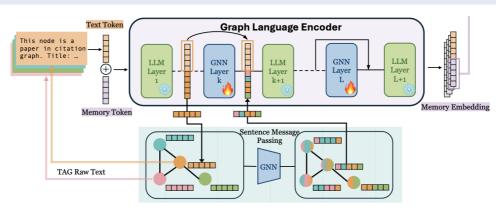


Figure 7: The Encoder.

Introduction

July 7, 2025

Background

Methodology 00000000 Experiment 000000000



Decoding the Answer

- 1. The LLM \rightarrow GNN cycle repeats, progressively deepening the fused representation.
- 2. A tanh gate on the GNN output acts as a safety switch, allowing the model to ignore graph structure if it's not useful.
- 3. Finally, the enriched memory tokens from the **NOG** (the query node) are fed to the LLM Decoder.
- 4. The Decoder, having received a perfectly briefed and context-rich summary, generates the final answer.

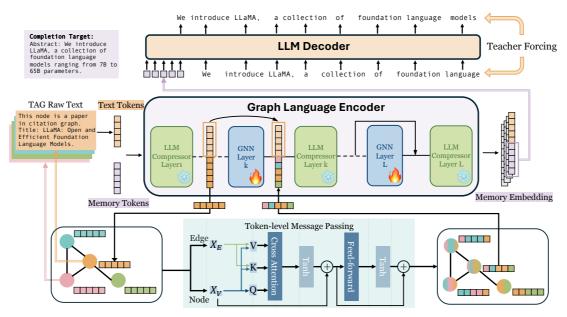


Figure 8: GOFA Architecture.

Introduction 00000

July 7, 2025

Background

Methodology 0000000 Experiment 000000000



Summary of the Architecture

Key Takeaways

- GOFA is modular: can plug into any decoder-only LLM.
- GNN layers are interleaved for structural awareness.
- Handles any graph task via a generative interface.
- A truly graph-aware, instruction-following LLM.

Introduction

Background 0000000 Methodology 0000000 Experiment 000000000



Experimental Overview

We will answer four fundamental questions:

- Q1: Effectiveness Does the training work?
- **Q2:** Generalization Can it handle unseen tasks?
- **Q3: Efficiency** Is it better than just using an LLM?
- **Q4: Fluency** Is it a truly flexible reasoner?

Introduction 000000

July 7, 2025

Background

Methodology

Experiment •000000000



The Setup

Core Method

■ All tasks are converted into subgraph problems, using k-hop subgraphs around target nodes as model input.

Pre-training

- Objective: GOFA is pre-trained using four self-supervised tasks.
- Data: Training utilizes large-scale, diverse graph datasets like MAG240M, Arxiv, Pubmed, and Wikikg90m.
- Setup: Initially, only the GNN layers are tuned. Training was performed on 8 NVIDIA A100 GPUs.

Fine-tuning

- **Zero-Shot**: The model is instruction-tuned on a small set of tasks (e.g., from Arxiv/Pubmed) to learn task formats. Prompts for this stage include detailed instructions and options.
- **Supervised**: The model is simultaneously fine-tuned on the training sets of all evaluation datasets, with prompts designed for direct answer generation.

Introduction

July 7, 2025

Background

Methodology 00000000 Experiment

O O O O O O O O O



The Setup

Datasets	Baselines
 Citation networks: Cora, Arxiv Knowledge graphs: FB15K237 Product graph: Amazon 	LLMs: LLaMA2, MistralGraph LLMs: GraphGPT, UniGraph, OFA
■ Commonsense QA: ExplaGraphs	

Introduction 000000

Background 00000000 Methodology 00000000 Experiment 000000000



4 Types of Pretraining Tasks

A powerful architecture needs powerful training. GOFA is pre-trained on four novel tasks designed to build the three key properties of a GFM.

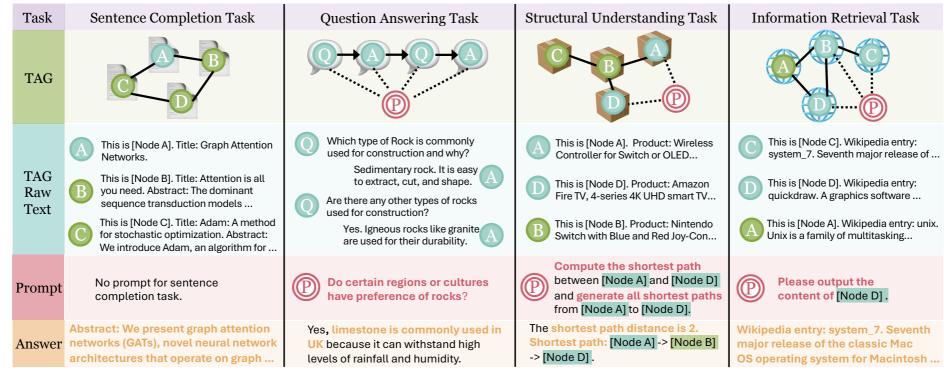


Figure 9: Examples of our pre-training tasks.

Introduction

Background

Methodology

Experiment 000000000



Q1: Does the training work?

Setup

- Perplexity ↓: How "surprised" the model is by new text. Lower is better.
- SPD & CN RMSE ↓: Error in predicting graph structure. Lower is better.
- **Test**: Evaluate on unseen datasets vs. base LLM.

Results

■ Language: GOFA perplexity is much lower (21.3 vs 30.1).

■ Structure: GOFA error (RMSE) is significantly lower.

Table 1: Evaluation for pre-trained GOFA.

	Perplexity ↓	SPD ↓	CN↓
Mistral-7B	30.12	1.254	1.035
GOFA-SN	26.20	_	_
GOFA	21.34	0.634	0.326

Introduction

Background 0000000 Methodology 00000000 Experiment 0000 00000



Q2: Can it handle unseen tasks?

Zero-Shot Classification Results

Setup: Zero-Shot Evaluation

- Process: Instruction-tuned on a small set of tasks.
- **Test**: Evaluated on completely new, unseen datasets.
- Baselines: Compared against SOTA Graph-LLMs (UniGraph, LLaGA).

Table 2: Node Classification Task Results (Accuracy).

Task	Cora-	Node	Wik	iCS		Products	3	ExplaGraphs	Cora-Link
Way	7	2	10	5	47	10	5	2	2
LLama2-7B Mistral-7B	47.92 60.54	73.45 88.39	40.10 63.63	58.77 71.90	27.65 43.99	58.71 70.16	64.33 74.94	57.76 68.77	48.15 49.43
OFA-Llama2 GraphGPT UniGraph ZeroG LLaGA	28.65 44.65 69.53 64.21 51.85	56.92 - 89.74 87.83 62.73	21.20 - 43.45 31.26	35.15 - 60.23 48.25	19.37 18.84 38.45 31.24 23.10	30.43 - 66.07 51.24 34.15	39.31 - 75.73 71.29 39.72	51.36	52.22 50.74 - - 88.09
GOFA-T GOFA-F	70.81 69.41	85.73 87.52	71.17 68.84	80.93 80.62	54.60 56.13	79.33 80.03	87.13 88.34	79.49 71.34	85.10 86.31

Introduction 000000

July 7, 2025

Background

Methodology

Experiment 00000000000



Q2: Can it handle unseen tasks?

Zero-Shot Classification Results

Setup: Zero-Shot Evaluation

- Process: Instruction-tuned on a small set of tasks.
- **Test**: Evaluated on completely new, unseen datasets.
- Baselines: Compared against SOTA Graph-LLMs (UniGraph, LLaGA).

Table 3: Link Classification Task Results (Accuracy).

Task	FB15K237	SceneGraphs
Format	10-Way	QA
Llama2-7B	48.32	38.62
Mistral-7B	62.48	45.95
GOFA-T	73.59	34.06
GOFA-F	80.69	31.36

Introduction 00000

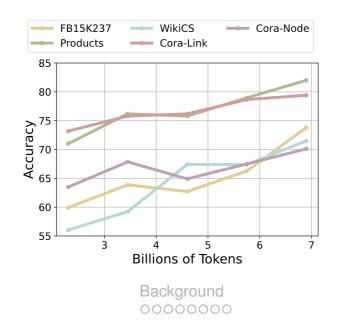
Background 0000000 Methodology 00000000 

Intuition Behind Zero-Shot Ability

1. Scaling Works

More pre-training data leads to better performance. This is a hallmark of a successful foundation model.

Figure 10: Performance vs. pre-training sample size.



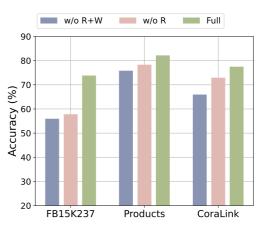
Introduction 00000

Methodology

2. Every Pre-training Task Matters

An ablation study shows that removing specific pre-training tasks hurts downstream performance. Removing the **Information Retrieval** task is particularly harmful for KGs.

Figure 11: Ablation study results.



 Conclusion



33/42 July 7, 2025 Zhaokun Wang: GOFA

Q3: Is it better than just using an LLM?

The Experiment

■ **GOFA:** Gets the graph as input.

LLM-N: Gets all node texts *plus* textual descriptions of every connection (e.g., "Node A connects to Node B").

Table 4: Comparison between GOFA and LLM with the same input.

Task Metric	ExplaGraphs Acc ↑	Time sec/sample ↓	WikiCS Acc ↑	Time sec/sample ↓	Cora-Link Acc ↑	Time sec/sample ↓	FB15k237 Acc ↑	Time sec/sample ↓
LLM-N	74.13	1.50	OOM	OOM	50.36	3.84	51.25	3.92
GOFA-F	79.49	0.48	71.17	2.43	85.10	1.67	73.49	3.37
Improvement	7.23%	68.00%	NA	NA	68.98%	56.51%	43.40%	14.03%

Accuracy		Efficienc	Efficiency			
GOFA is far mor improvement).	re accurate on all tasks (e.g	., 69% GOFA is mon large gr	nuch faster and avoids "Out of raphs.	f Memory" errors		
Introduction	Background	Methodology	Experiment	Conclusion		



000000

0000000000

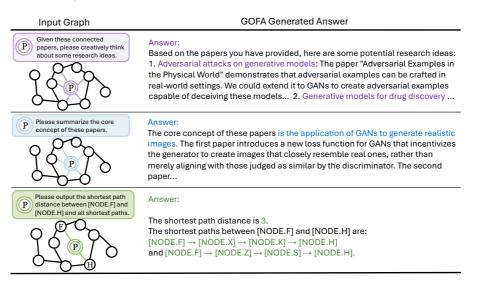
000000

Q4: Can GOFA "Think" Flexibly?

Setup

The same graph is given to GOFA with three different open-ended prompts.

Figure 12: Comparison between GOFA and LLM with the same input.



Introduction

July 7, 2025

Background

Methodology

Experiment 0000000000



Summary of Findings

- GOFA matches LLMs in text tasks
- Outperforms on graph-related tasks
- Generalizes to unseen datasets and instructions

One model. All tasks. Structure + Language.

Introduction 00000

Background

Methodology

Experiment 00000000



The Bigger Picture: Towards Unified Al Reasoning

What GOFA enables:

- Complex structural QA: "List all paths from X to Y"
- Conceptual synthesis: "Summarize the main idea of a citation cluster"
- Graph querying via prompts

The Pieces Coming Together:

Text + Graph + Instruction = Universal Reasoning Interface

Introduction

Background

Methodology

Experiment 000000000

Conclusion ●○○○○○



Limitations

■ Frozen LLM Compressor

The compressor can't adapt to the specific semantics of the graph data.

Format Dependence

Instruction tuning assumes known task templates; completely new formats are still a challenge.

No Truly End-to-End Training

The encoder and decoder components are still trained or utilized in separate stages.

Introduction 000000

Background

Methodology

Experiment

00000



Looking Ahead: What Comes Next?

Unified Training

Train the compressor, GNN, and decoder components end-to-end for deeper integration.

Universal Graph Agents

Combine GOFA with retrieval systems or external tools to act on real-world knowledge graphs dynamically.

Cross-Modal Graph Reasoning

Extend the architecture to reason over graphs where nodes are images, audio, or other complex data types (e.g., visual scene graphs).

Introduction

Background

Methodology

Experiment 000000000



A Thought to Leave With...

"Graphs are how the world connects. Language is how we explain it. What if a model could do both?"

GOFA is a bold step. But it's only the beginning.

Introduction

July 7, 2025

Background

Methodology

Experiment 000000000





References

- [1] L. Kong, J. Feng, H. Liu, C. Huang, J. Huang, Y. Chen, M. Zhang (2024). Gofa: A Generative One-for-All Model for Joint Graph Language Modeling. arXiv:2407.09709.
- [2] Z. Wang, Z. Liu, T. Ma, J. Li, Z. Zhang, X. Fu, Y. Li, Z. Yuan, W. Song, Y. Ma, et al. (2025). *Graph Foundation Models: A Comprehensive Survey*. arXiv:2505.15116.
- [3] Y. Shi, Z. Huang, S. Feng, H. Zhong, W. Wang, Y. Sun (2021).

 Masked Label Prediction: Unified Message Passing Model for Semi-Supervised Classification.

 arXiv:2009.03509.
- [4] T. Ge, J. Hu, L. Wang, X. Wang, S.-Q. Chen, F. Wei (2024).

 In-context Autoencoder for Context Compression in a Large Language Model.

 arXiv:2307.06945.

Introduction

July 7, 2025

Background 0000000

Methodology

Experiment 000000000





Q&A

Thank You! Questions?

Introduction 00000 Background 0000000 Methodology 00000000 Experiment 000000000

