



# On the Intersection of Language and Graph Models



Chuxu Zhang

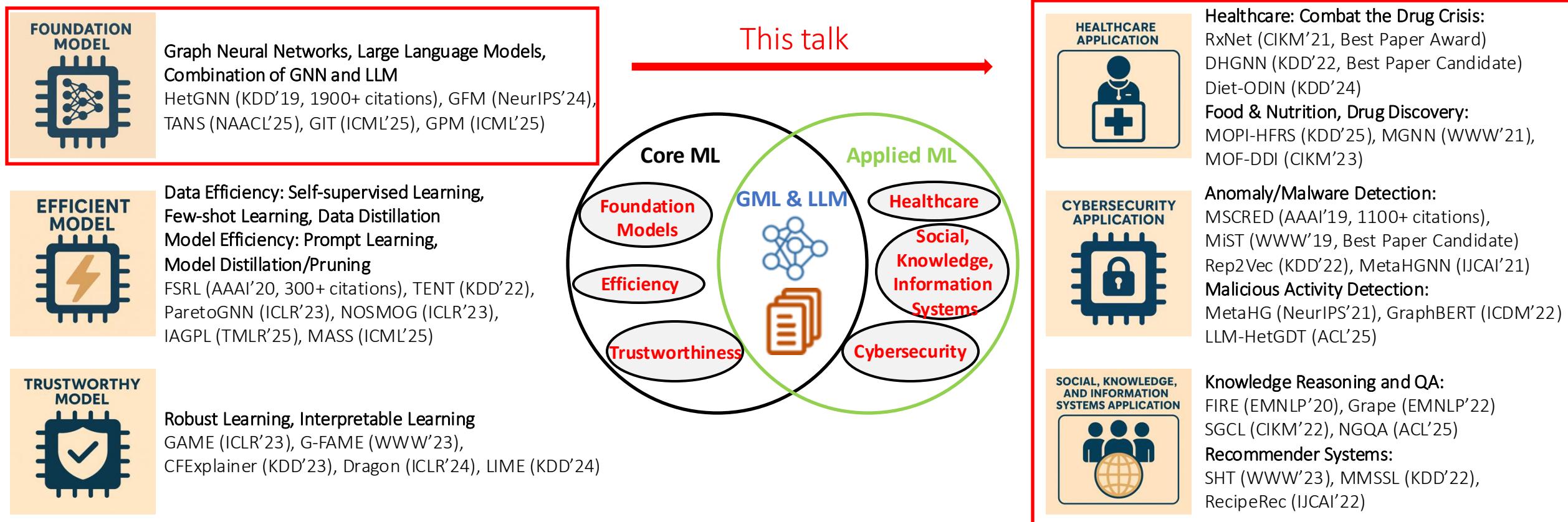
Associate Professor of Computer Science

University of Connecticut (UConn)

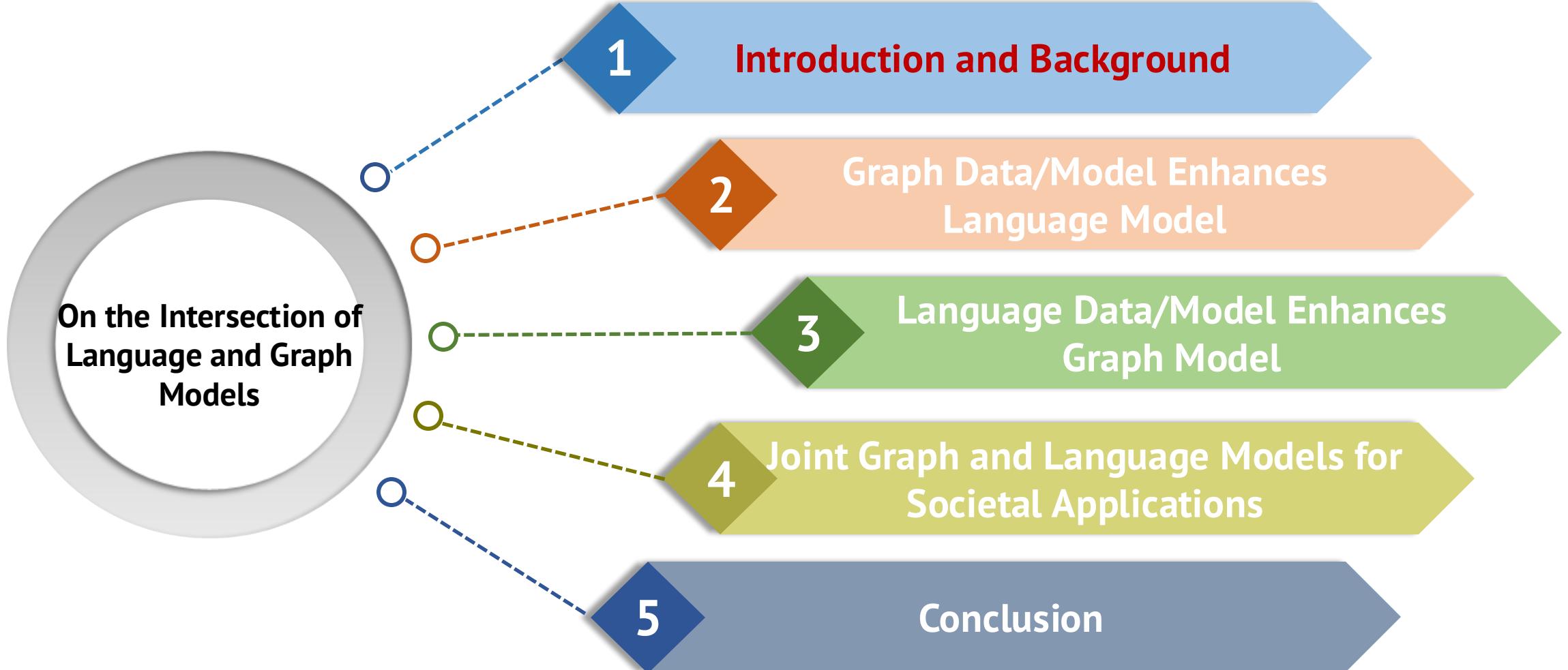
ICDM ISIR-eCom Workshop, Nov. 2025

# Research Overview: AI, ML, DS and Their Societal Applications

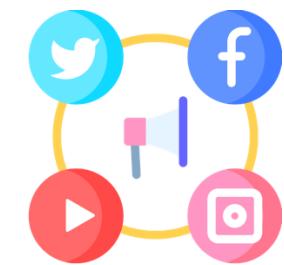
AI, Machine Learning, Data Science  
graph/network data (graph machine learning)  
text/language data (large language models)



# Outline



# Introduction: Various Data in Real-world Applications



Social Media



Cybersecurity/IoT



Knowledge System



Healthcare

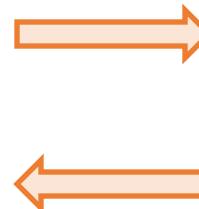
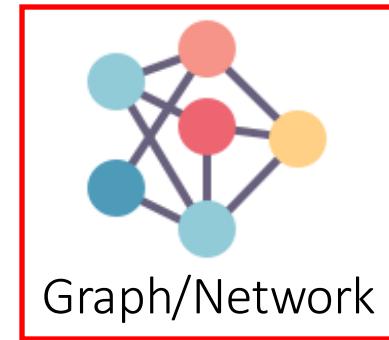
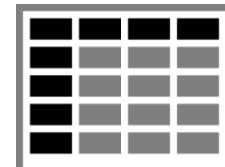


Science

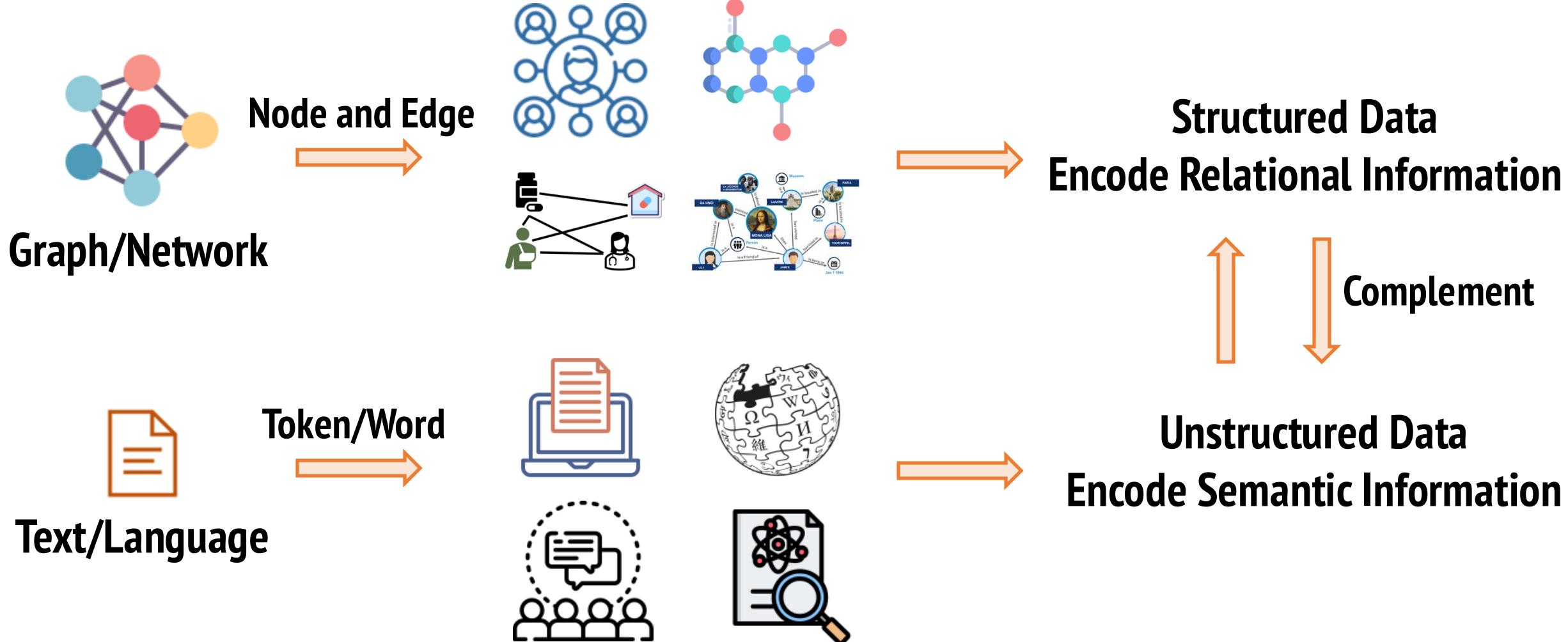


Web/Information

Spatial-temporal      Tabular

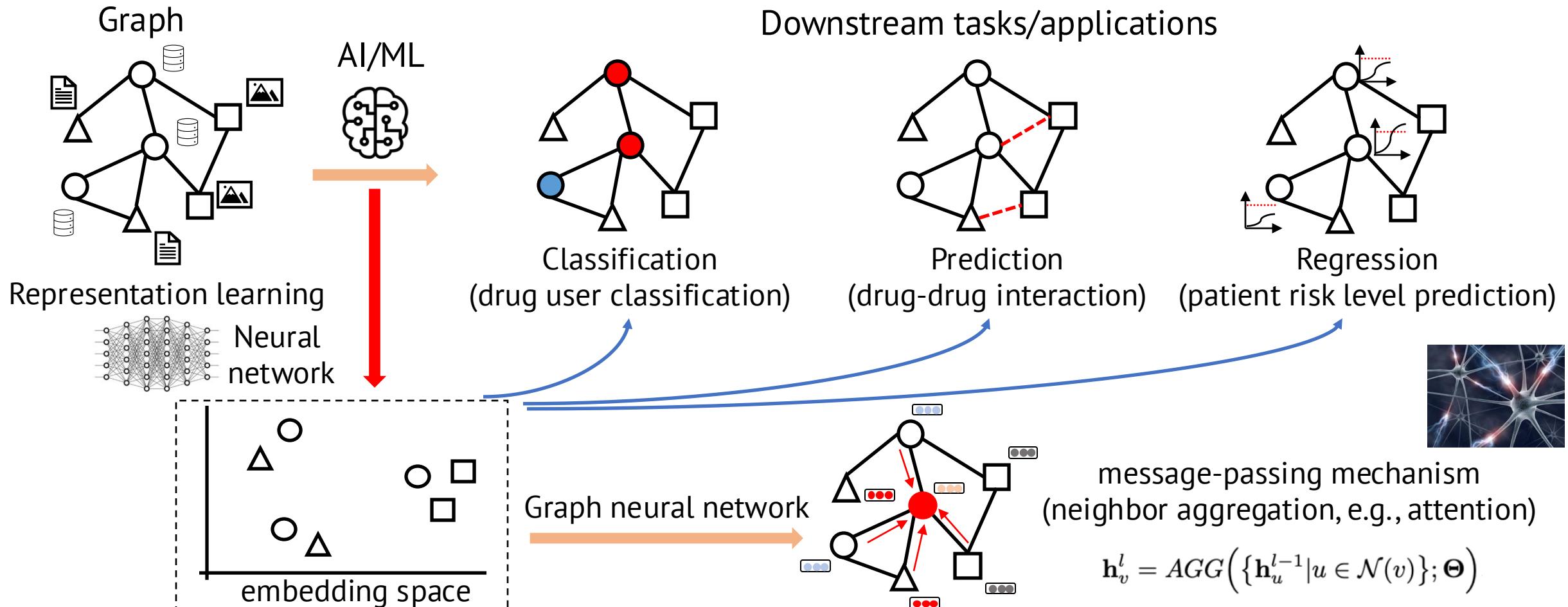


# Introduction: Graph (Network) and Text (Language) Data



# Introduction: Graph Neural Networks (GNNs)

- GNNs learn representations of nodes by iteratively **transforming** and **aggregating/propagating** the features from their neighborhoods



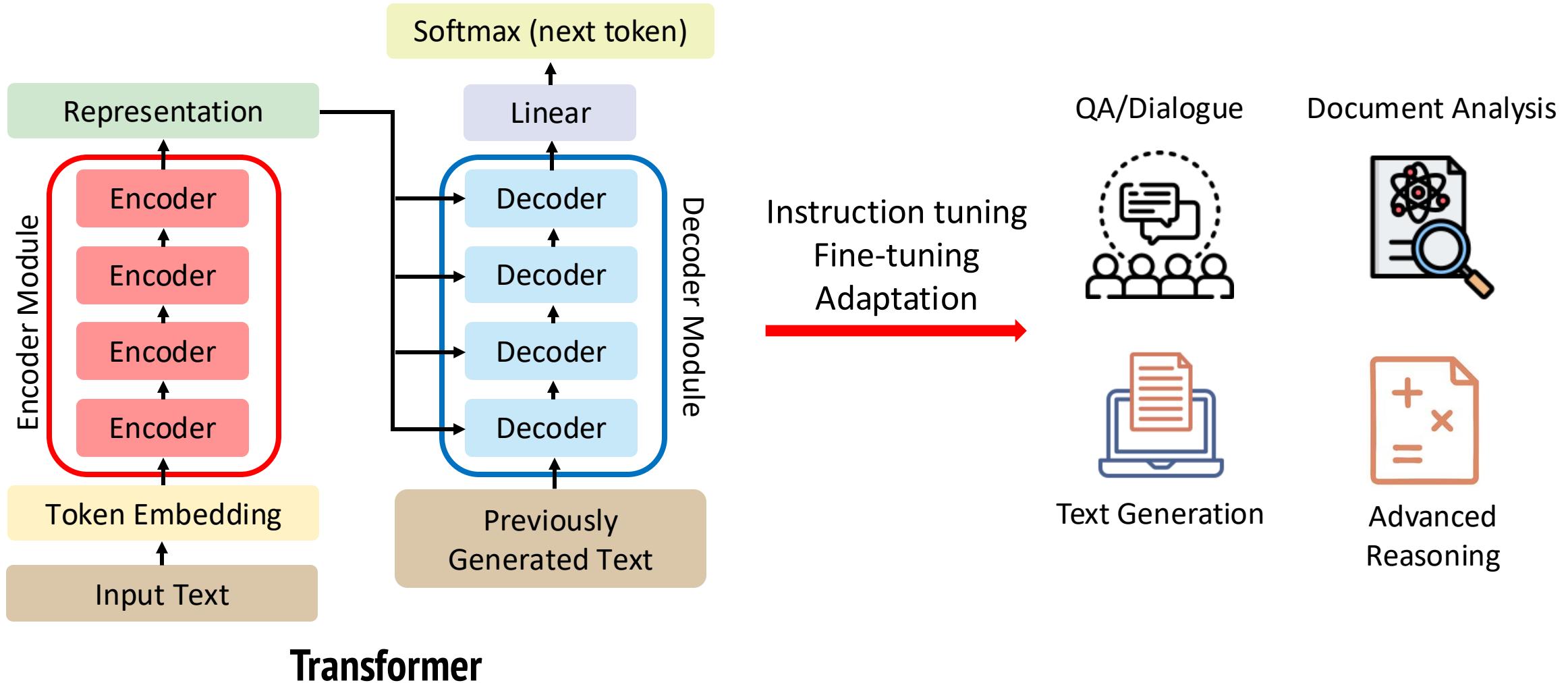
# Introduction: Graph Neural Networks (GNNs)

---

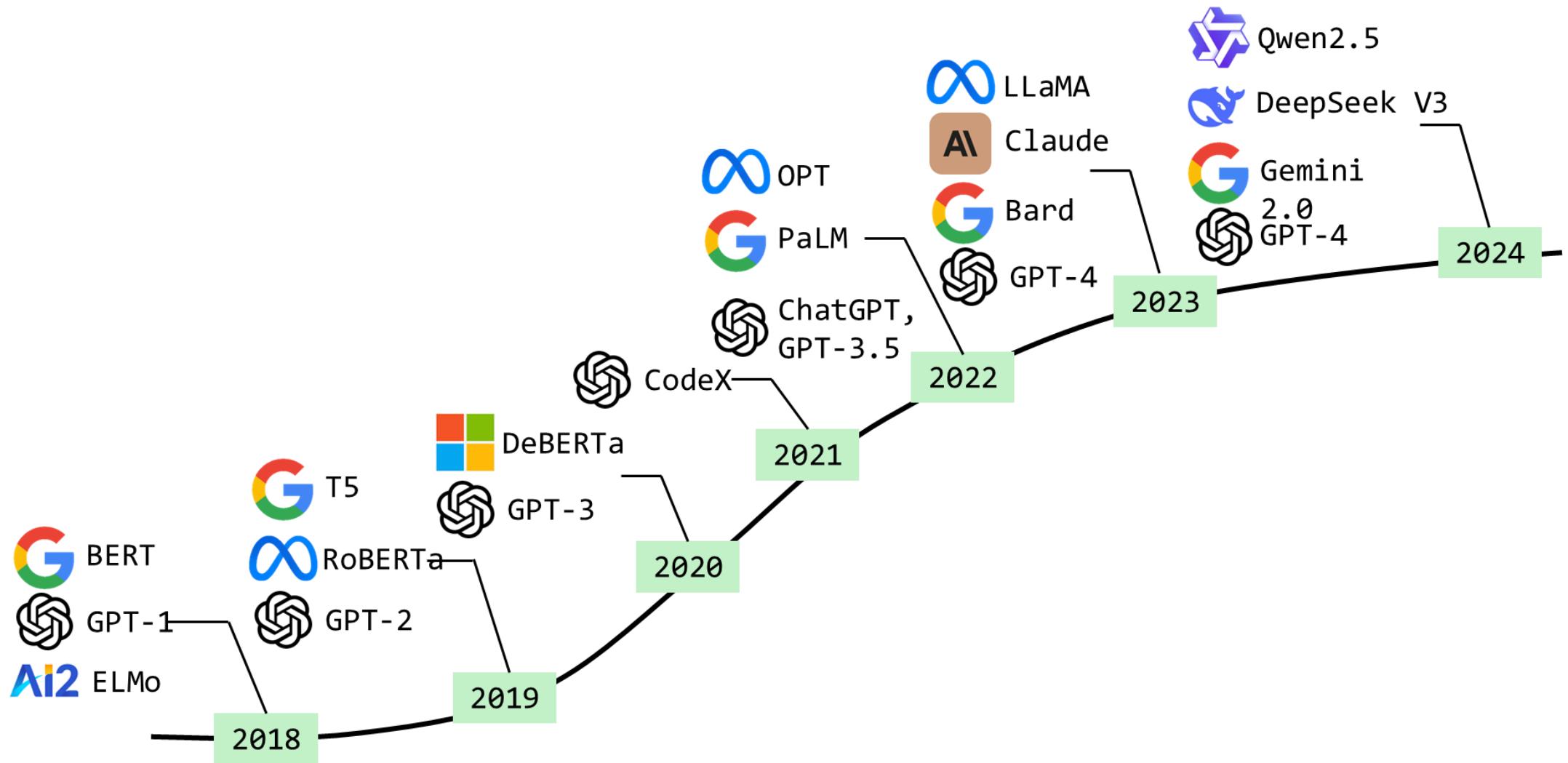
- **GCN (ICLR'16)**
  - Convolution Aggregation
- **GSAGE (NeurIPS'17)**
  - Pooling/Recurrent Aggregation
- **GAT (ICLR'18)**
  - Attention Aggregation
- **HetGNN (PhD work, KDD'19)**
  - The first GNN on Heterogeneous Graphs
  - 1900+ Citations, 2024 ICBS Frontiers of Science Award
- **Our recent works: GFT (NeurIPS'24), GIT (ICML'25), GPM (ICML'25), G<sup>2</sup>PM (NeurIPS'25)**
  - Graph Foundation Models Across Datasets/Tasks/Domains

# Introduction: Large Language Model (LLMs)

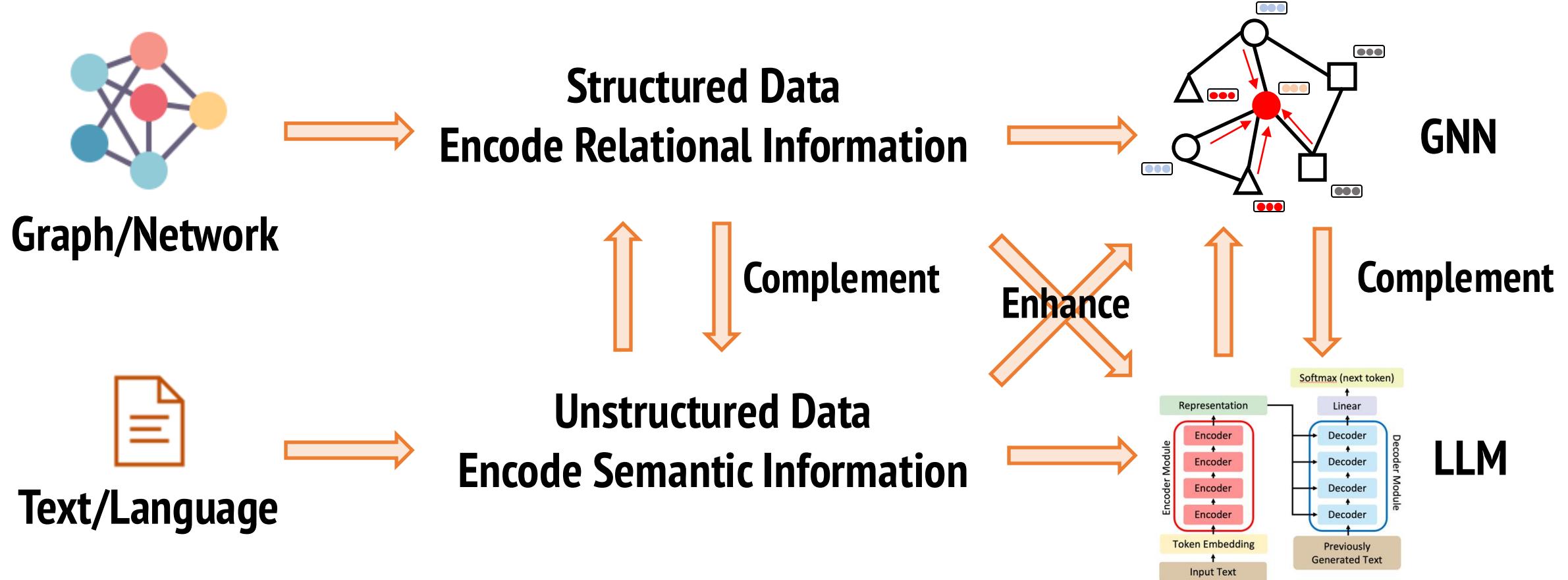
- Pretrained Architecture



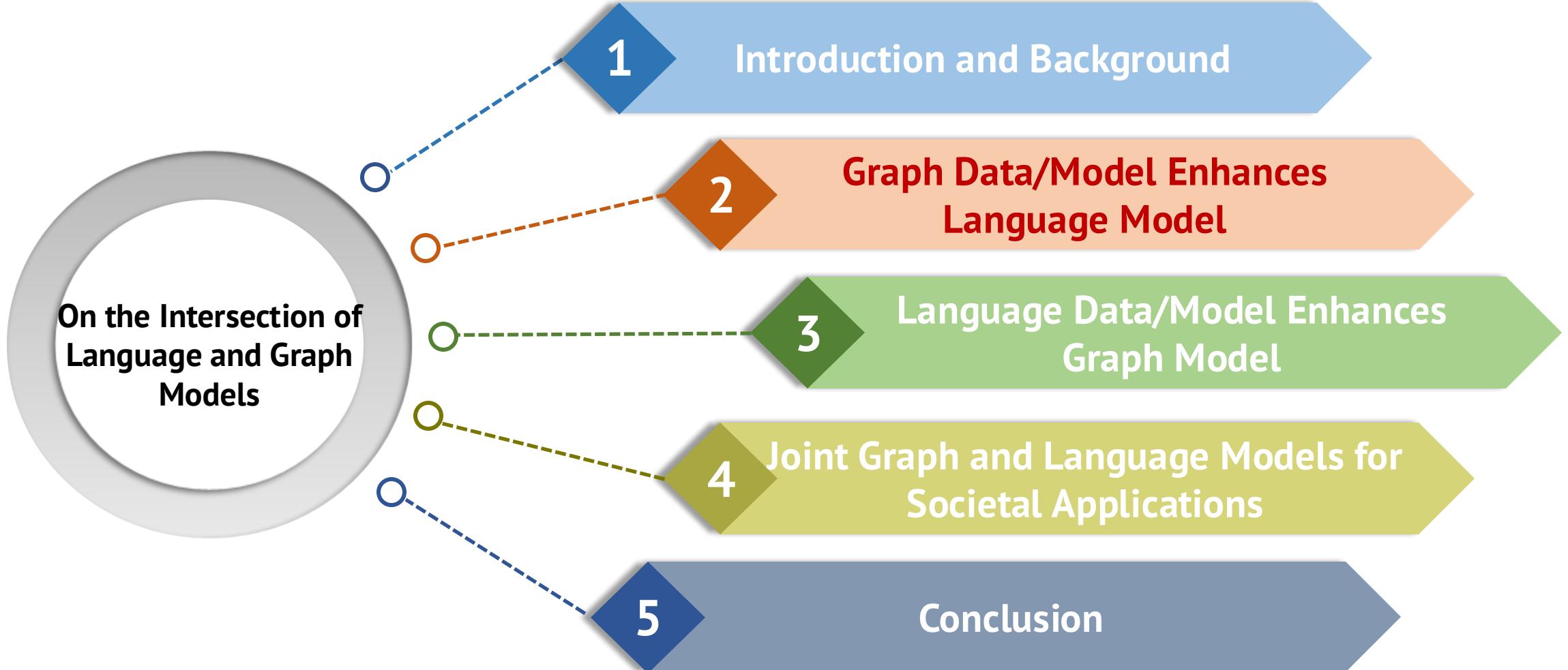
# Introduction: Large Language Model (LLMs)



# Introduction: Graph Model and Language Model

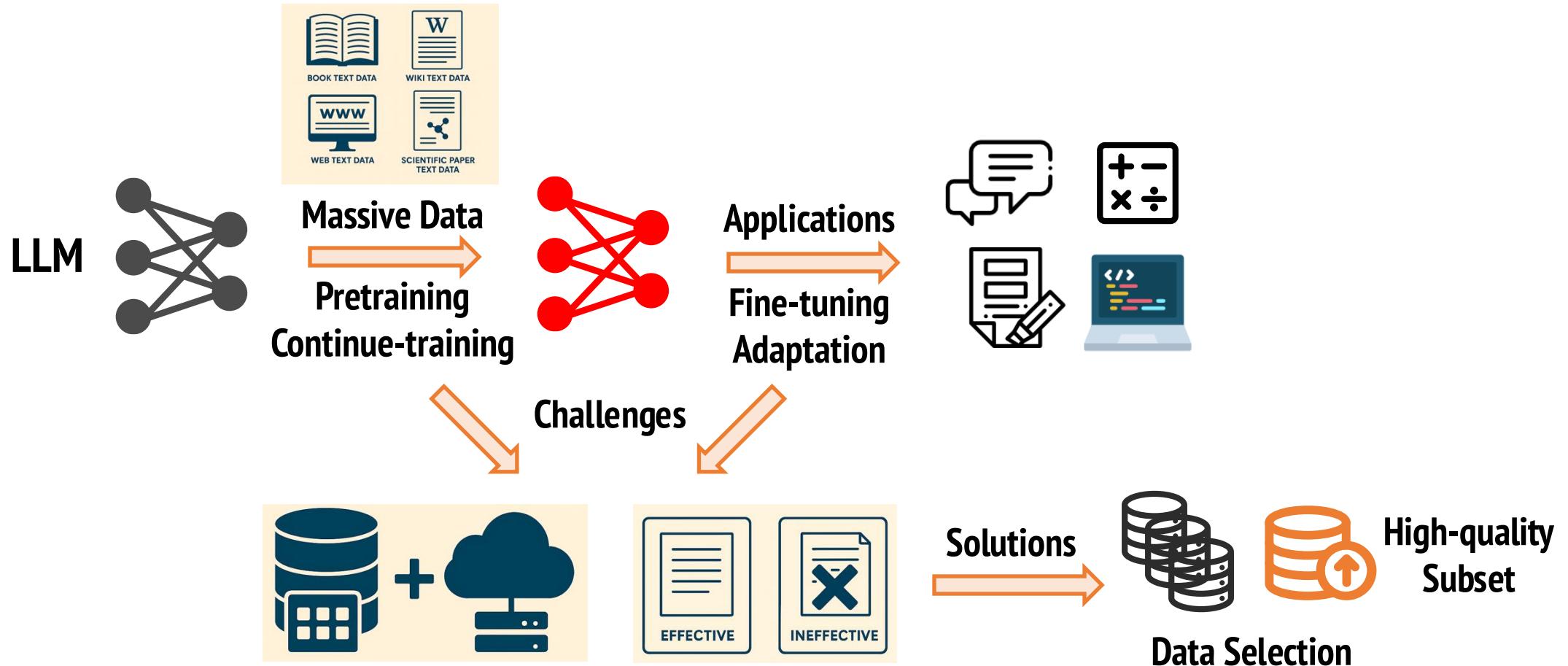


# Outline



# Graph Data/Model Enhances Language Model: MASS

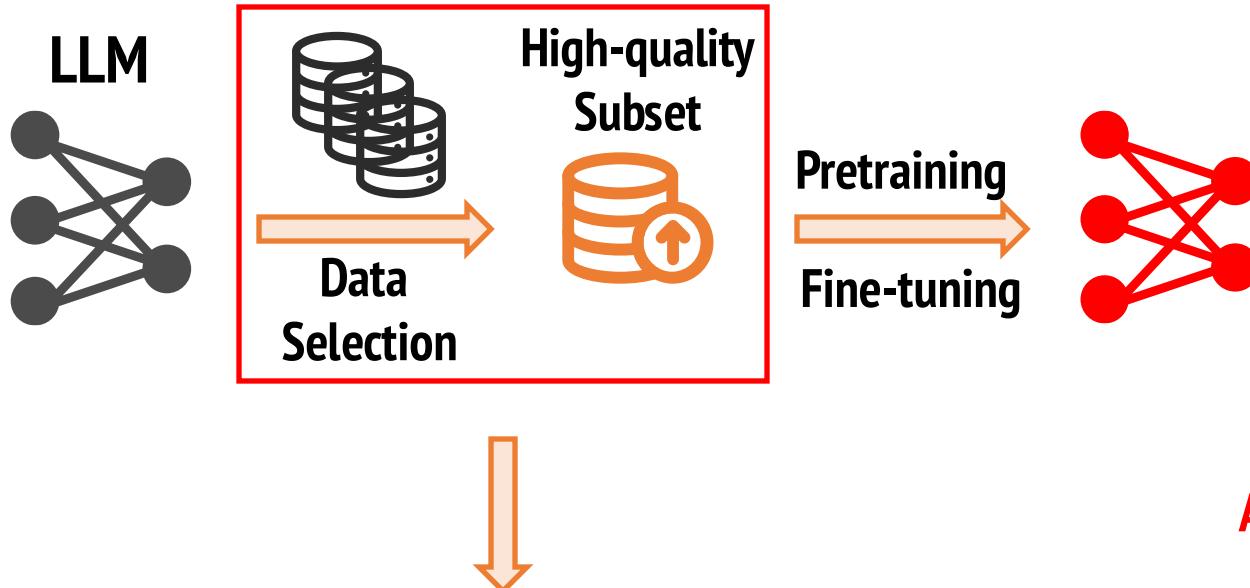
- Data Selection for Pretraining LLMs: Improve Training **Efficiency** and **Effectiveness**



MASS: MAthematical Data Selection via Skill Graphs for Pretraining Large Language Models, ICML'25

# Graph Data/Model Enhances Language Model: MASS

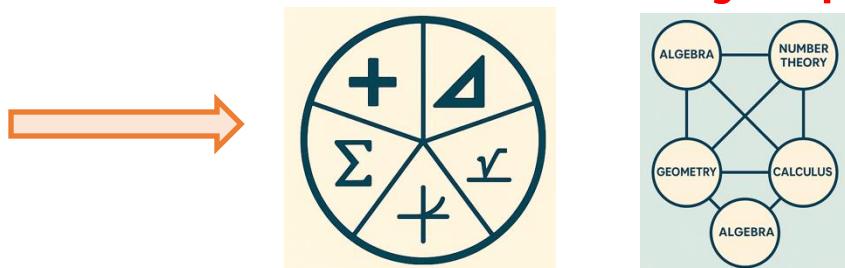
- Data Selection for Pretraining LLMs: Improve Training **Efficiency** and **Effectiveness**



Prior methods: e.g., RHO-1 [NeurIPS'24], AutoDS [ICLR'24]

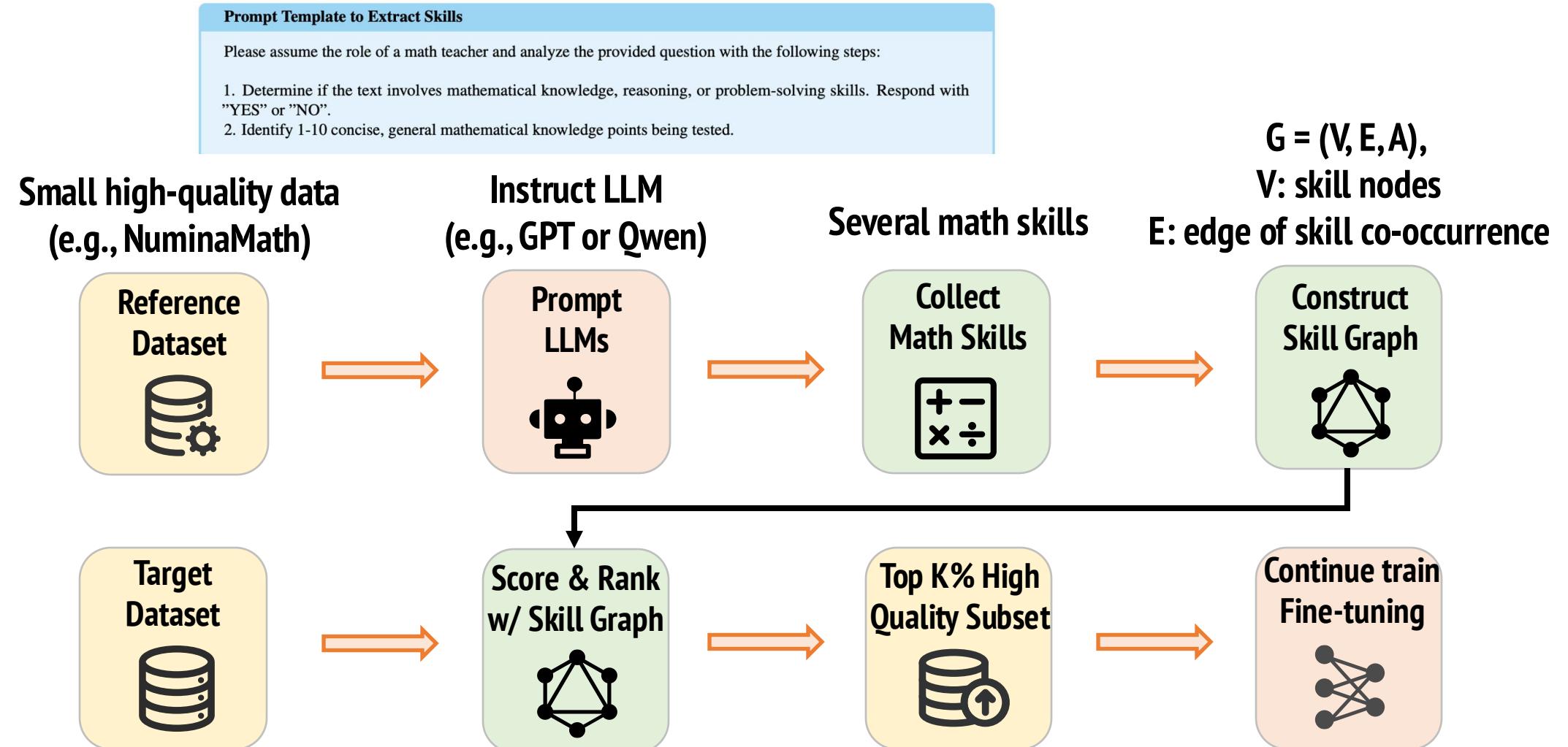
Limitation: concentrate on general domains while neglecting the underlying knowledge and their interrelations of specific domain data, such as mathematical skills for advanced reasoning capabilities.

**Assumption:** a data point reflecting (1) more important math skills or (2) more compositional information of math skills should receive a higher quality score.



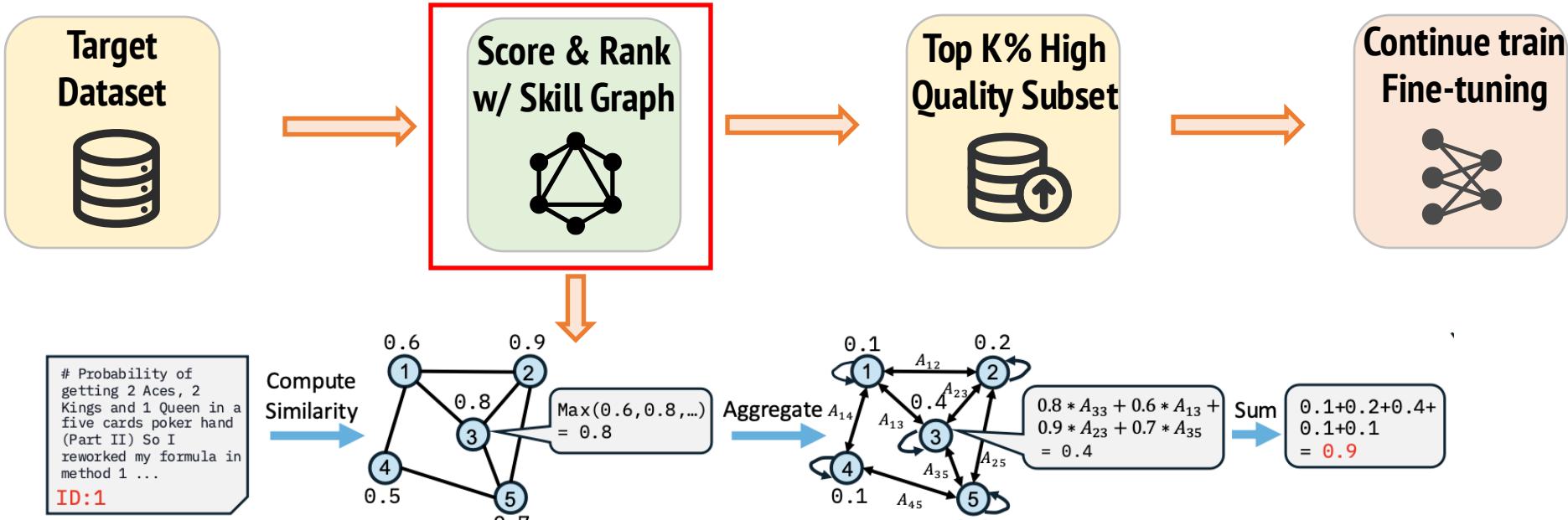
MASS: MAthematical Data Selection via Skill Graphs for Pretraining Large Language Models, ICML'25

# Graph Data/Model Enhances Language Model: MASS



MASS: MAthematical Data Selection via Skill Graphs for Pretraining Large Language Models, ICML'25

# Graph Data/Model Enhances Language Model: MASS



$$\mathbf{A}_{i,i} = \sigma(v_i^{cnt}, T) = \frac{\exp\left(\frac{v_i^{cnt}}{T}\right)}{\sum_{j=1}^{|V|} \exp\left(\frac{v_j^{cnt}}{T}\right)}$$

$$\mathbf{A}_{i,j} = \sigma(e_{ij}^{cnt}, T) = \frac{\exp\left(\frac{e_{ij}^{cnt}}{T}\right)}{\sum_{(e_{ij} \in E)} \exp\left(\frac{e_{ij}^{cnt}}{T}\right)}$$

$$\text{sim}(x_i, v_j) = \max_{k \in v_j^{ids}} \cos(\text{Emb}(x_i), \text{Emb}(d_k))$$

$$\text{score}(x) = \sum_{j=1}^{|V|} \text{sim}_{agg}(x, v_j)$$

$$= \sum_{j=1}^{|V|} \mathbf{A}_{j,j} \text{sim}(x, v_j) + \sum_{j=1}^{|V|} \sum_{v_k \in \mathcal{N}(v_j)} \mathbf{A}_{j,k} \text{sim}(x, v_k)$$

Why it works

- (1) more important math skills
- (2) more compositional information on math skills

# Graph Data/Model Enhances Language Model: MASS

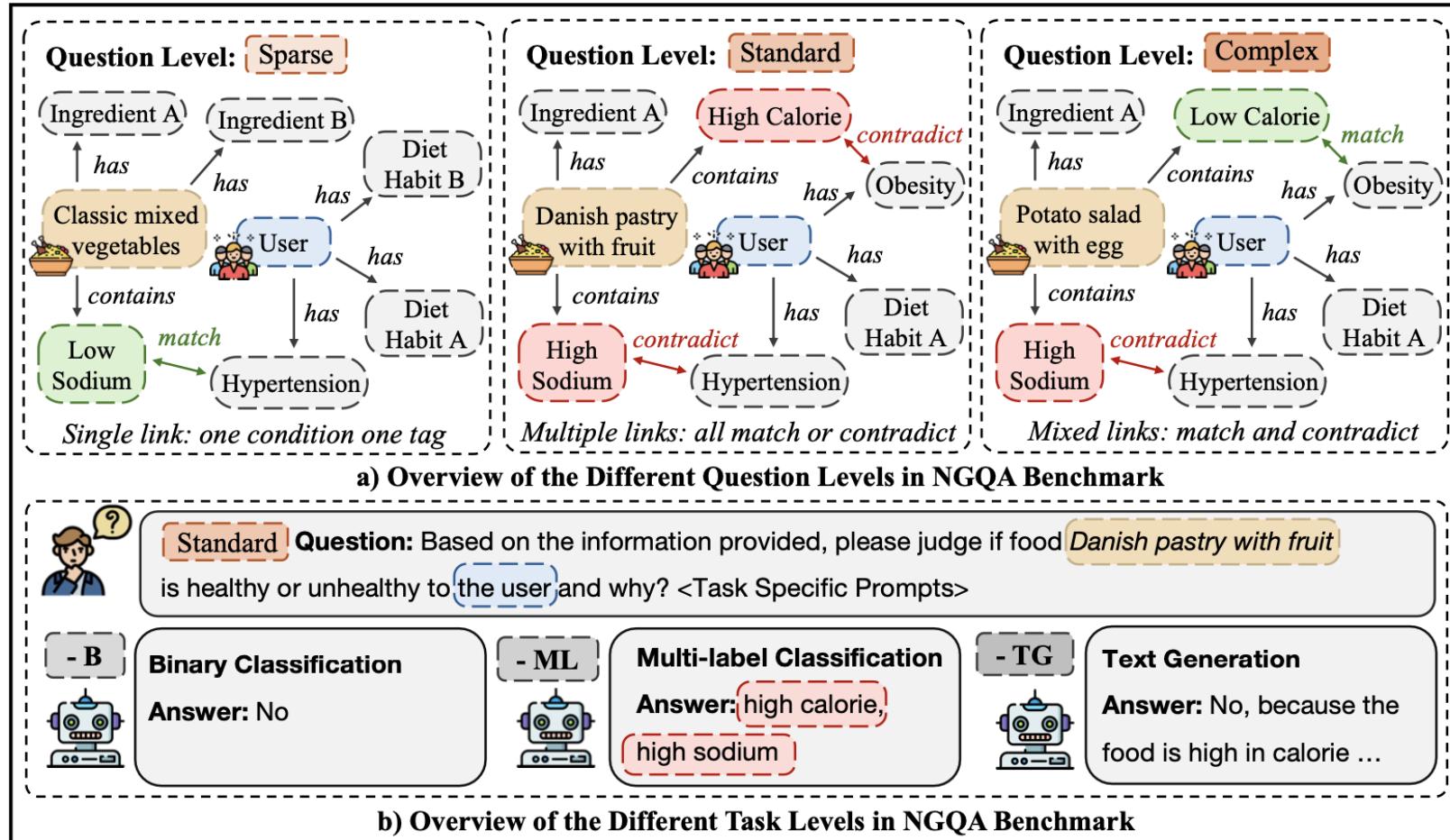
*Table 2.* The main experimental results. TinyLlama-1.1B and Mistral-7B are continuously pretrained using both the original and selected subsets of OpenWebMath, OpenWebMath-pro, and Jiuzhang3.0. The **bolded** entries indicate the best results within each setting.  
 \* indicates that results are from ProX (Zhou et al., 2024a)

Dataset	Method	Unique Tokens	Trained Tokens	GSM8K	MATH	SVAMP	ASDiv	MAWPS	TAB	MQA	MMLU STEM	SAT MATH	Avg.
TinyLlama-1.1B													
	w/o continual pretraining			2.7	2.8	10.9	17.9	20.5	12.5	14.0	16.3	21.9	13.3
OpenWebMath	-	14.6B	14.6B	5.2	3.0	20.7	31.4	41.0	14.6	10.1	19.5	<b>37.5</b>	20.3
	RULE*	6.5B	15B	4.5	2.8	17.5	29.4	39.3	15.1	12.4	19.4	25.0	18.4
	RHO-1*	14.6B	9B	7.1	<b>5.0</b>	23.5	41.2	53.8	-	<b>18.0</b>	-	-	-
	ProX	5.1B	14.6B	8.6	3.0	23.8	40.2	51.6	19.6	14.9	<b>26.1</b>	25.0	23.6
	DSIR	4.9B	14.6B	5.5	2.6	24.1	37.8	54.3	16.9	12.1	25.4	22.3	22.1
	AutoDS	4.9B	14.6B	7.3	2.4	22.9	39.2	52.7	18.4	13.8	23.2	24.1	22.7
	MASS	4.9B	14.6B	<b>9.0</b>	4.4	<b>24.9</b>	<b>41.4</b>	<b>54.8</b>	<b>21.5</b>	13.9	20.3	25.0	<b>23.9</b>
OpenWebMath -pro	-	5.1B	14.6B	8.6	3.0	23.8	40.2	51.6	19.6	14.9	26.1	25.0	23.6
	DSIR	3B	14.6B	8.8	3.2	<b>24.1</b>	41.5	53.1	18.9	14.3	<b>27.6</b>	27.5	24.4
	AutoDS	3B	14.6B	9.1	4.5	22.4	40.8	54.3	23.2	13.1	26.5	28.0	24.7
	MASS	3B	14.6B	<b>10.2</b>	<b>5.8</b>	23.8	<b>42.3</b>	<b>57.9</b>	<b>25.3</b>	<b>15.3</b>	27.0	<b>34.4</b>	<b>26.9</b>
Jiuzhang3.0	-	3.4B	6.8B	22.3	19.0	46.4	60.1	73.2	29.6	19.1	<b>24.0</b>	34.4	36.4
	DSIR	2.4B	6.8B	24.5	21.3	48.2	63.9	74.4	28.8	19.2	22.1	33.6	37.3
	AutoDS	2.4B	6.8B	26.7	20.8	51.3	66.7	73.5	31.1	19.3	22.4	32.8	38.3
	MASS	2.4B	6.8B	<b>30.1</b>	<b>24.8</b>	<b>52.5</b>	<b>69.1</b>	<b>80.7</b>	<b>32.9</b>	<b>20.4</b>	22.7	34.4	<b>40.8</b>
Mistral-7B													
	w/o continual pretraining			41.1	10.6	64.9	68.5	87.3	54.8	33.9	49.9	65.6	53.0
OpenWebMath	-	14.4B	9.6B	44.5	19.0	60.6	68.4	87.8	50.5	44.5	50.9	56.2	53.6
	MASS	4.8B	9.6B	<b>47.7</b>	<b>23.2</b>	<b>64.6</b>	<b>74.7</b>	<b>90.5</b>	<b>55.7</b>	<b>50.7</b>	<b>52.6</b>	<b>65.6</b>	<b>58.4</b>
OpenWebMath -pro	-	5.1B	5.1B	47.1	21.8	63.2	73.7	89.5	<b>58.2</b>	42.6	52.2	56.2	56.1
	MASS	3B	5.1B	<b>53.2</b>	<b>25.6</b>	<b>67.0</b>	<b>76.8</b>	<b>90.4</b>	57.6	<b>51.8</b>	<b>54.5</b>	<b>81.2</b>	<b>62.0</b>
Jiuzhang3.0	-	3.8B	3.8B	66.4	39.4	82.9	<b>85.9</b>	90.8	35.3	61.8	40.1	50.0	61.4
	MASS	2.7B	3.8B	<b>70.0</b>	<b>43.8</b>	<b>84.3</b>	85.7	<b>93.7</b>	<b>35.7</b>	<b>63.5</b>	<b>46.9</b>	<b>65.6</b>	<b>65.5</b>

MASS: MAthematical Data Selection via Skill Graphs for Pretraining Large Language Models, ICML'25

# Graph Data/Model Enhances Language Model: More Study

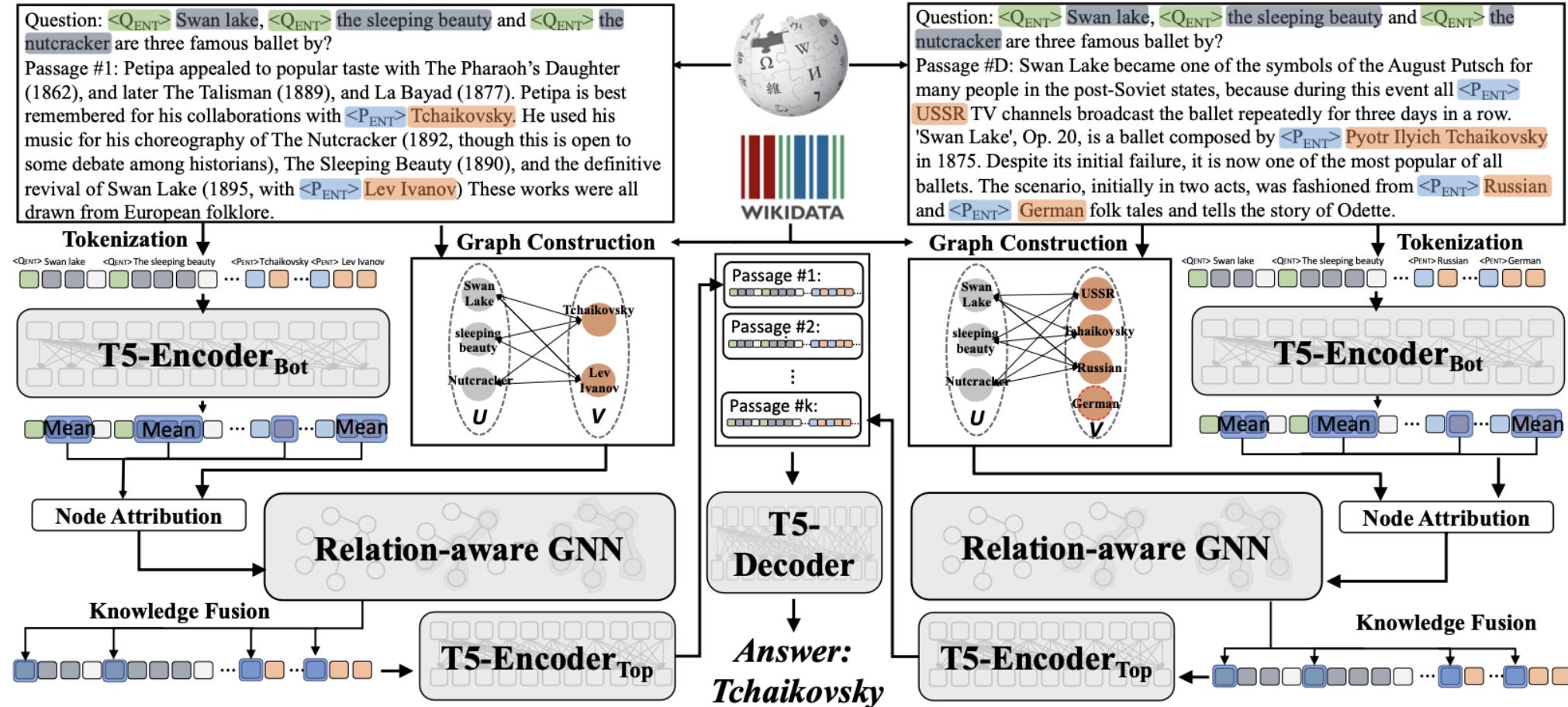
- QA: Knowledge Graph as Retrieval Augmentation for LLMs



NGQA: A Nutritional Graph Question Answering Benchmark for Personalized Health-aware Nutritional Reasoning, ACL 2025

# Graph Data/Model Enhances Language Model: More Study

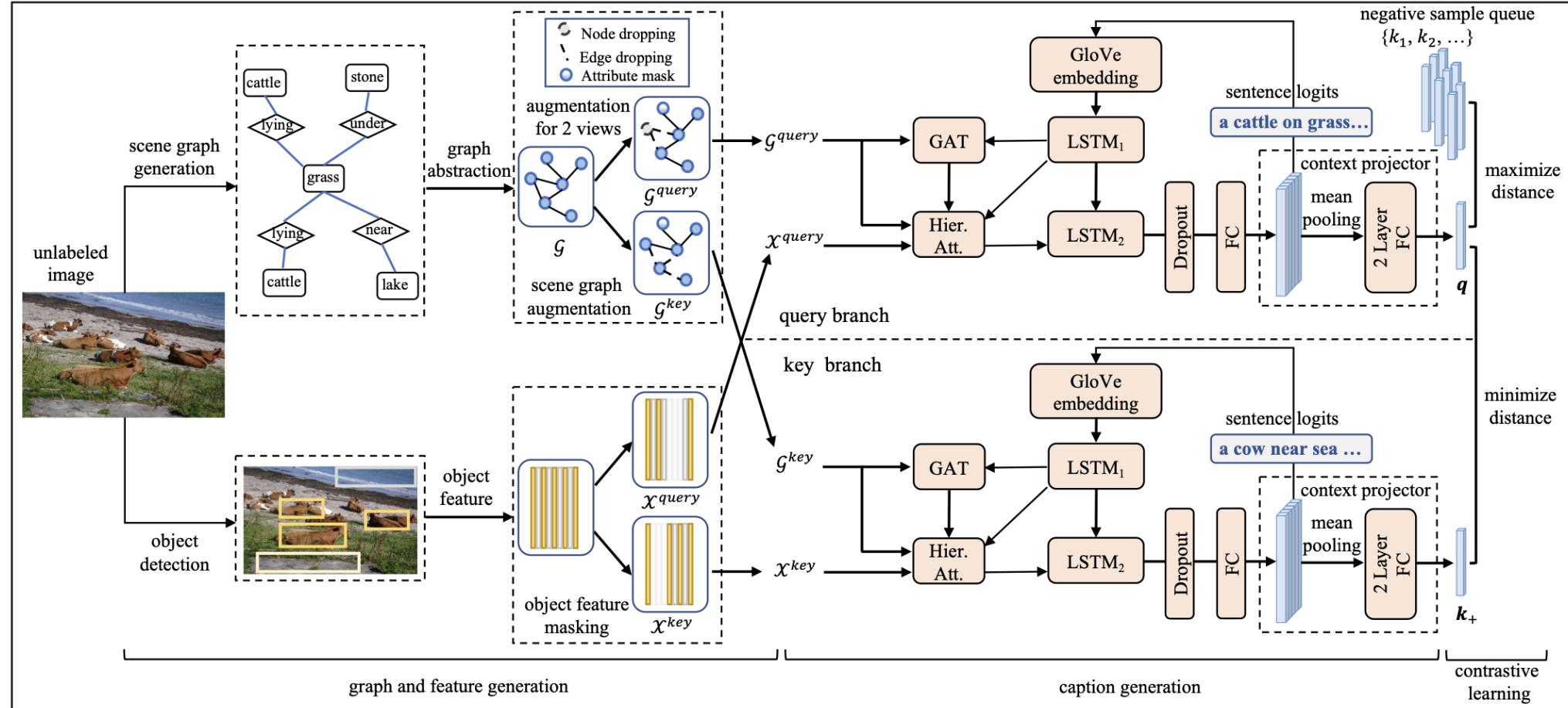
- QA: Graph as Augmented Information for LLMs



Knowledge Graph Enhanced Passage Reader for Open-domain Question Answering, EMNLP 2022

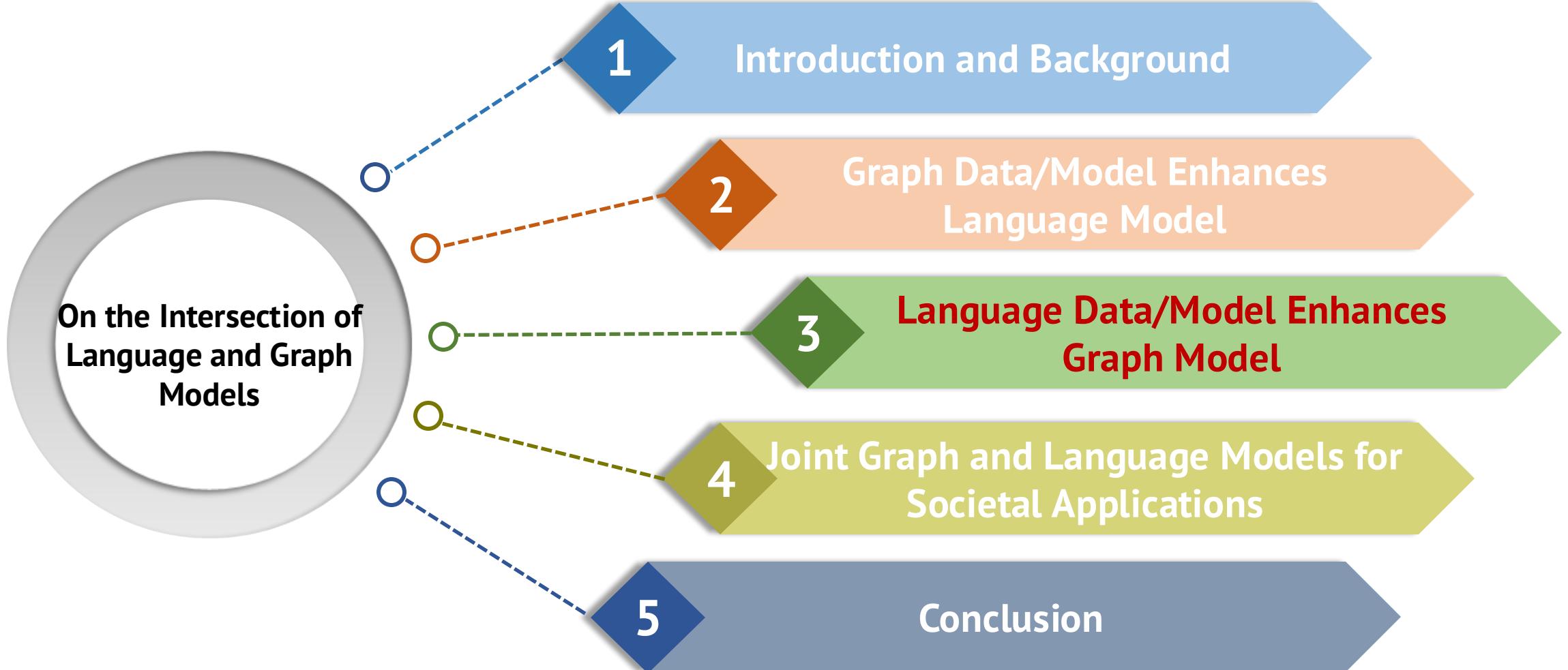
# Graph Data/Model Enhances Language Model: More Study

- Text Generation: Graph as Augmented Data for LLMs

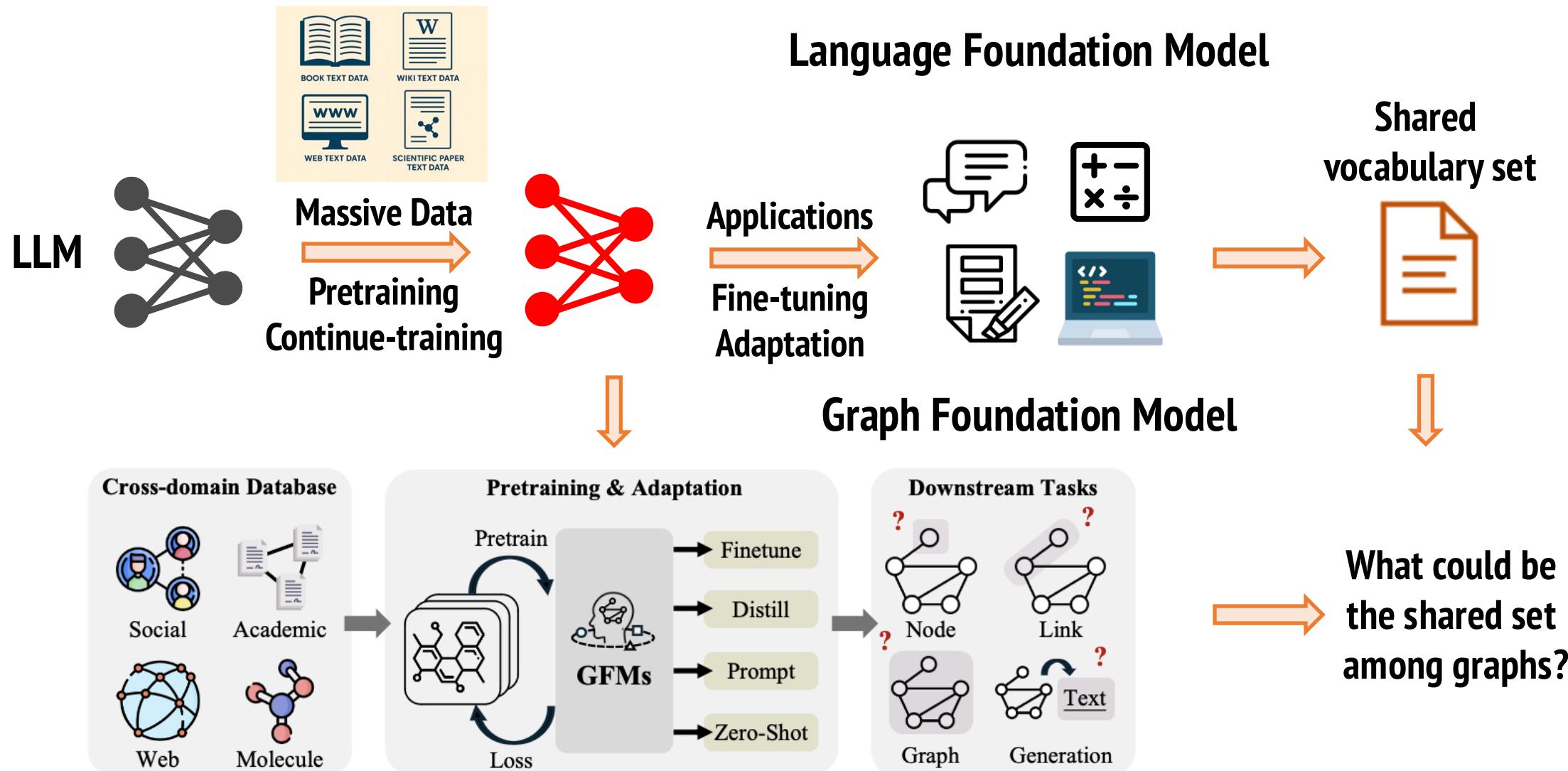


Look Twice as Much as You Say: Scene Graph Contrastive Learning for Self-Supervised Image Caption Generation, CIKM 2022

# Outline

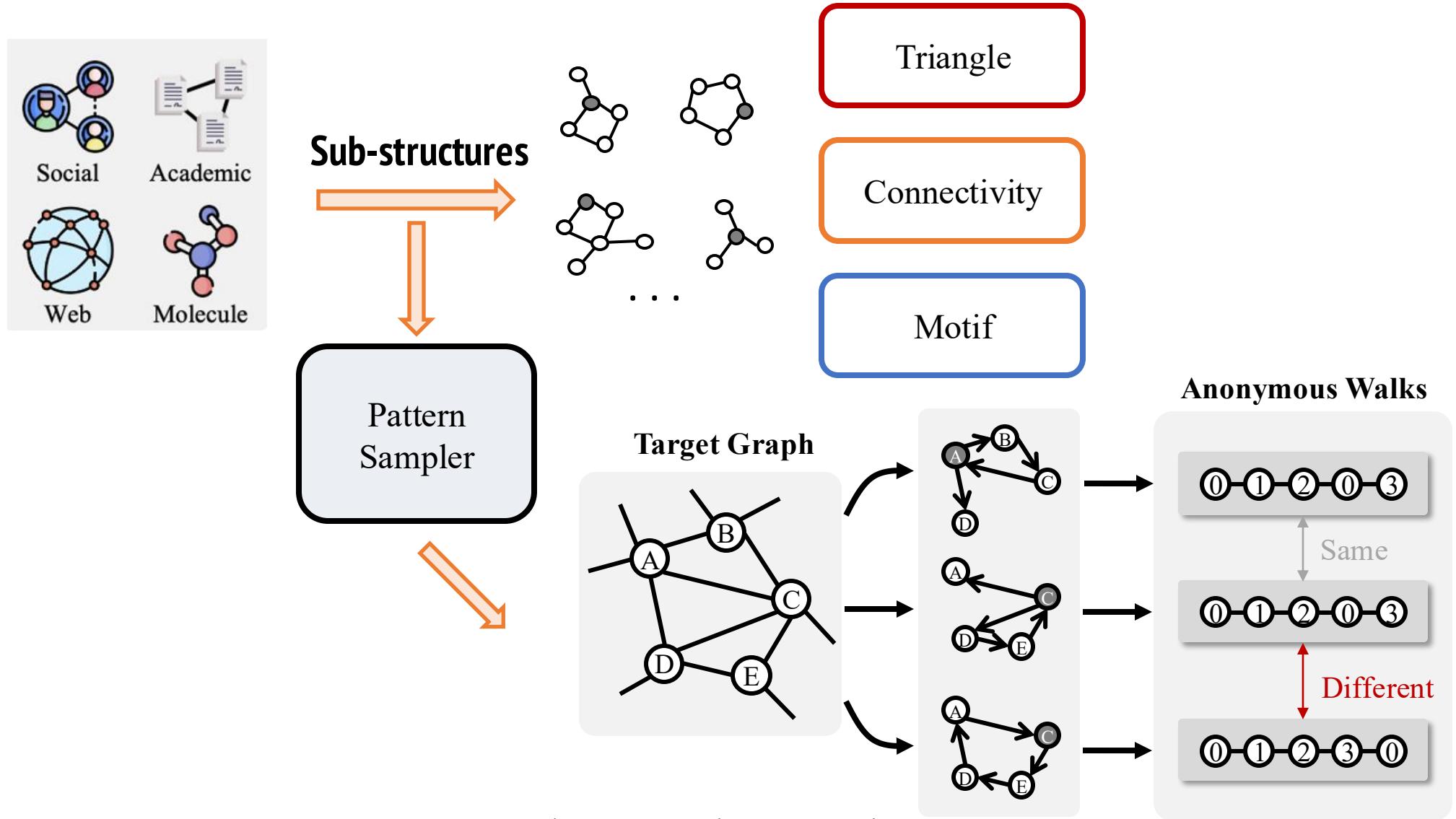


# Language Data/Model Enhances Graph Model: Foundation Model



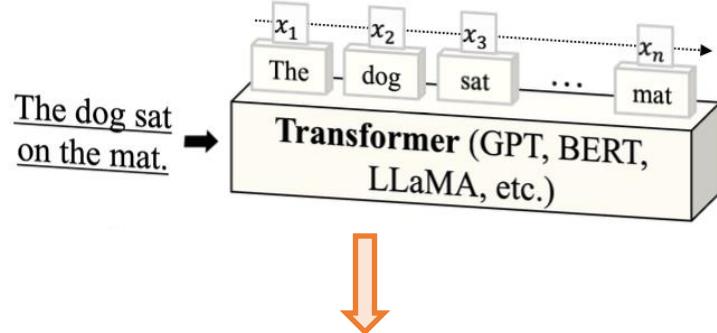
Generative Graph Pattern Machine, NeurIPS'25

# Language Data/Model Enhances Graph Model: Graph Foundation Model

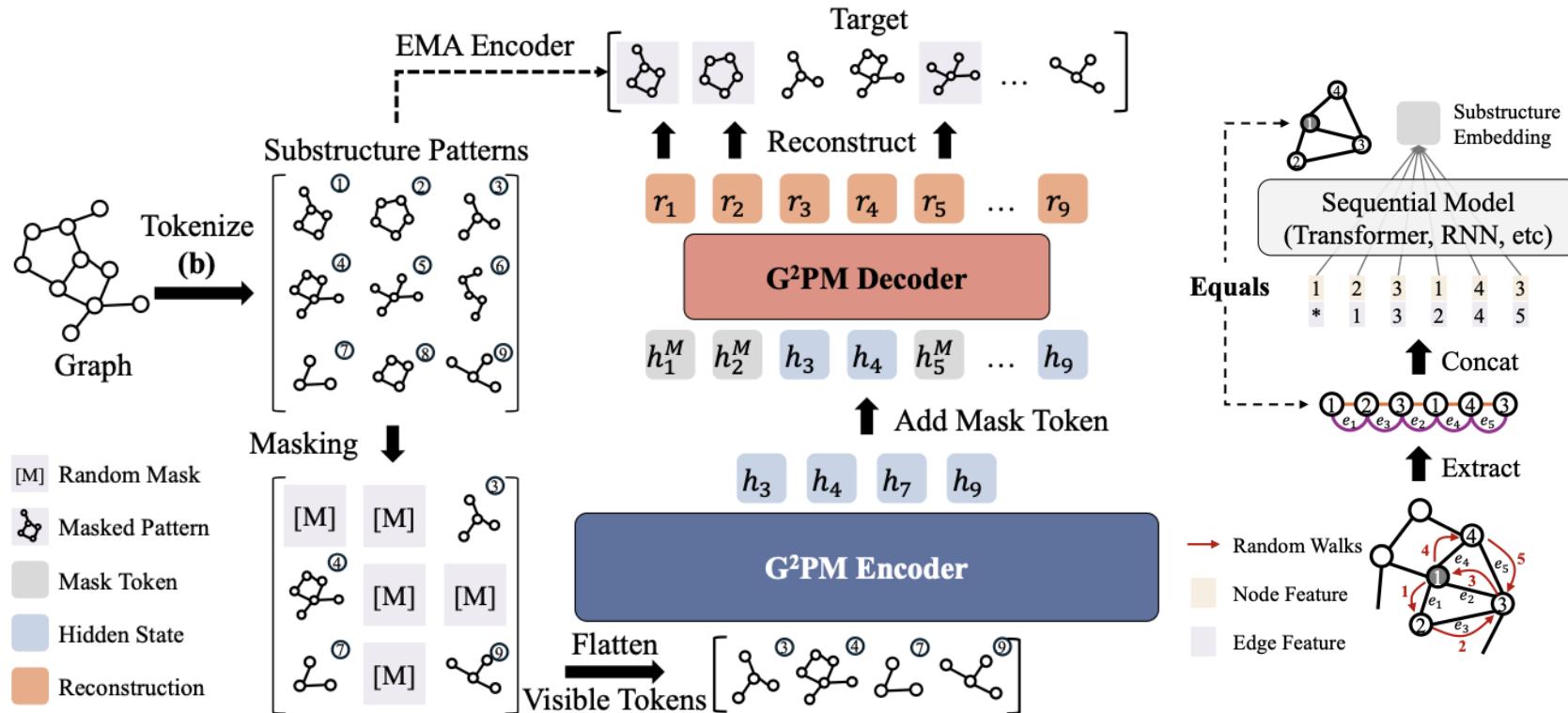
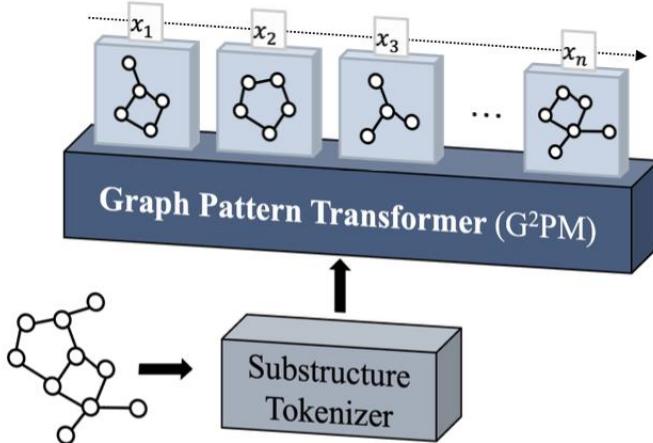


# Language Data/Model Enhances Graph Model: Graph Foundation Model

(a) Textual Generative Modeling via Word Sequence



(c) Graph Generative Modeling via Substructure Sequence (Ours)

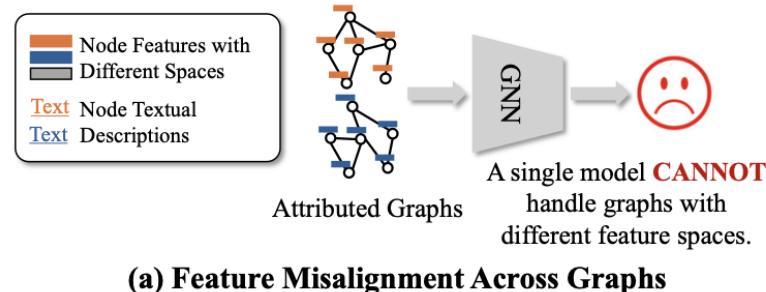


Generative Graph Pattern Machine, NeurIPS'25

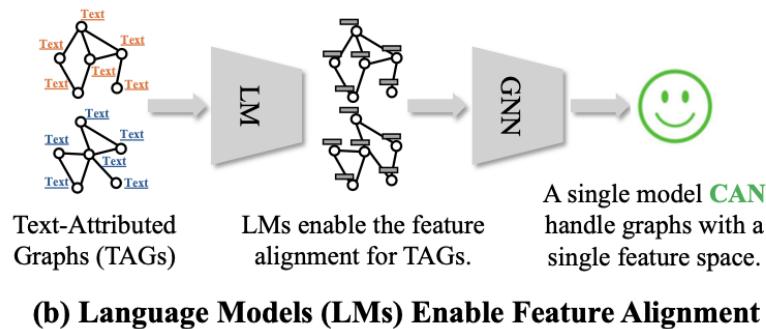


# Language Data/Model Enhances Graph Model

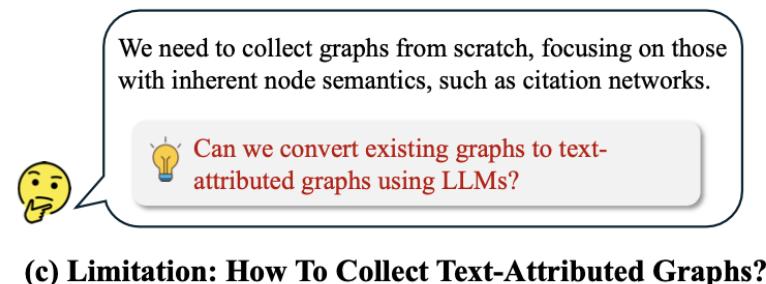
## • LLM as Text/Attribute Generation for GNNs



(a) Feature Misalignment Across Graphs



(b) Language Models (LMs) Enable Feature Alignment



(c) Limitation: How To Collect Text-Attributed Graphs?

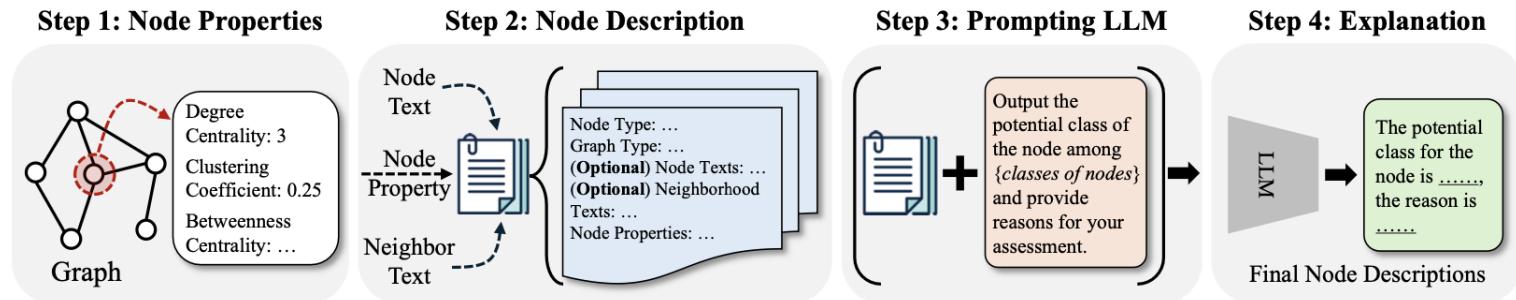


Figure 2: The framework of our topology-aware node description synthesis (TANS).

### Step 2: Generate Basic Node Descriptions

<b>Prefix</b>	Given a node from a {Graph Type} graph, where the node type is {Node Type} with {Node Number} nodes, and the edge type is {Edge Type} with {Edge Number} edges.
<b>Node Text (Optional)</b>	The original node description is {Original Textual Descriptions}.
<b>Neighbor Text (Optional)</b>	The following are the textual information of {k} connected nodes. The descriptions are: {Textual Descriptions of Selected Neighborhoods}.
<b>Node Property</b>	The value of {Node Property} is {Value of The Given Property}, ranked as {Rank of The Node}% among {Node Number} nodes.

### Step 3: Prompting LLMs

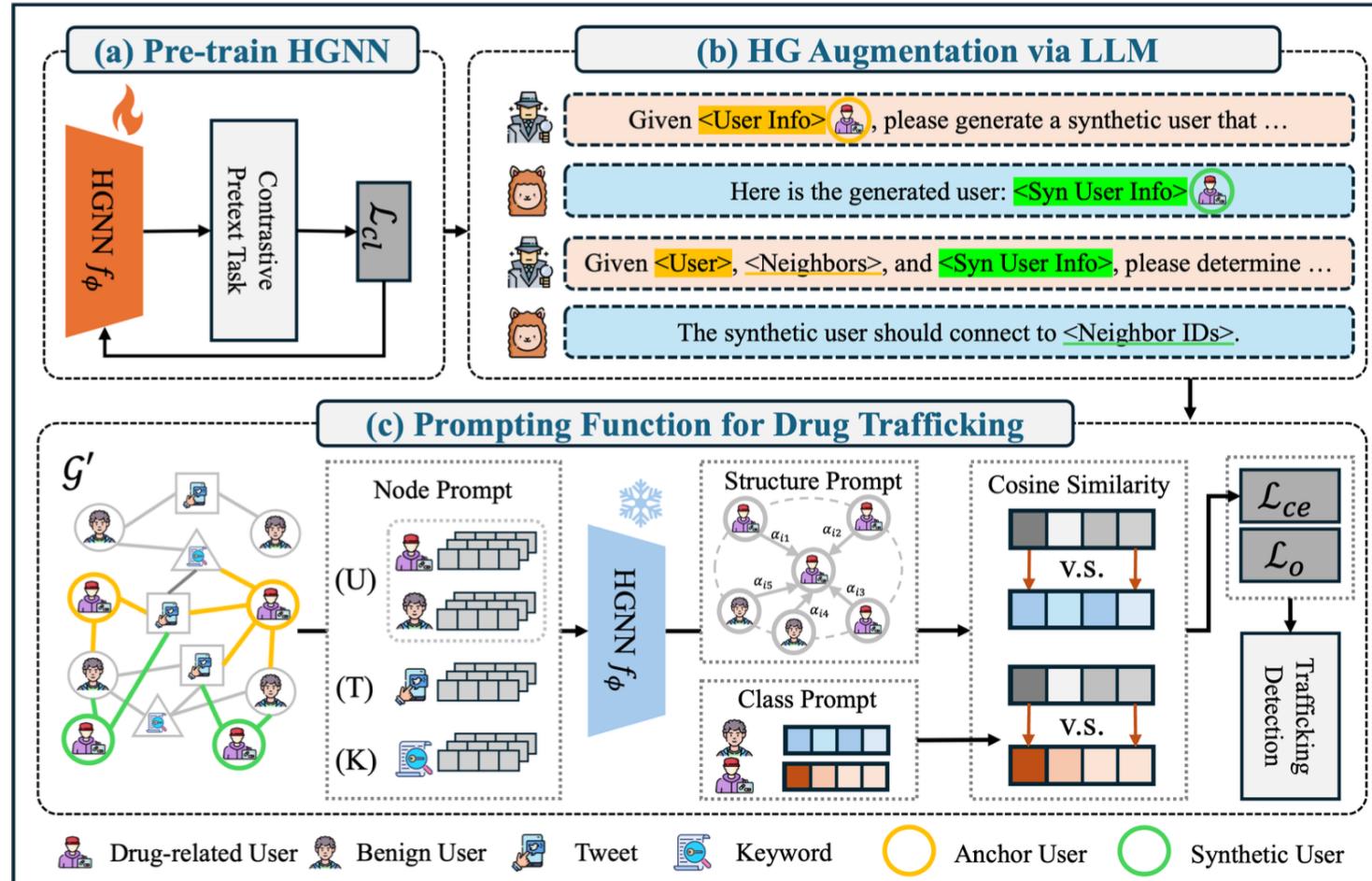
<b>Suffix</b>	Output the potential {k} classes of the node and provide reasons for your assessment. The classes include {Classes of Nodes}. Your answer should be less than 200 words.
---------------	--

Table 2: Prompt templates.

Can LLMs Convert Graphs to Text-Attributed Graphs? NAACL 2025

# Language Data/Model Enhances Graph Model

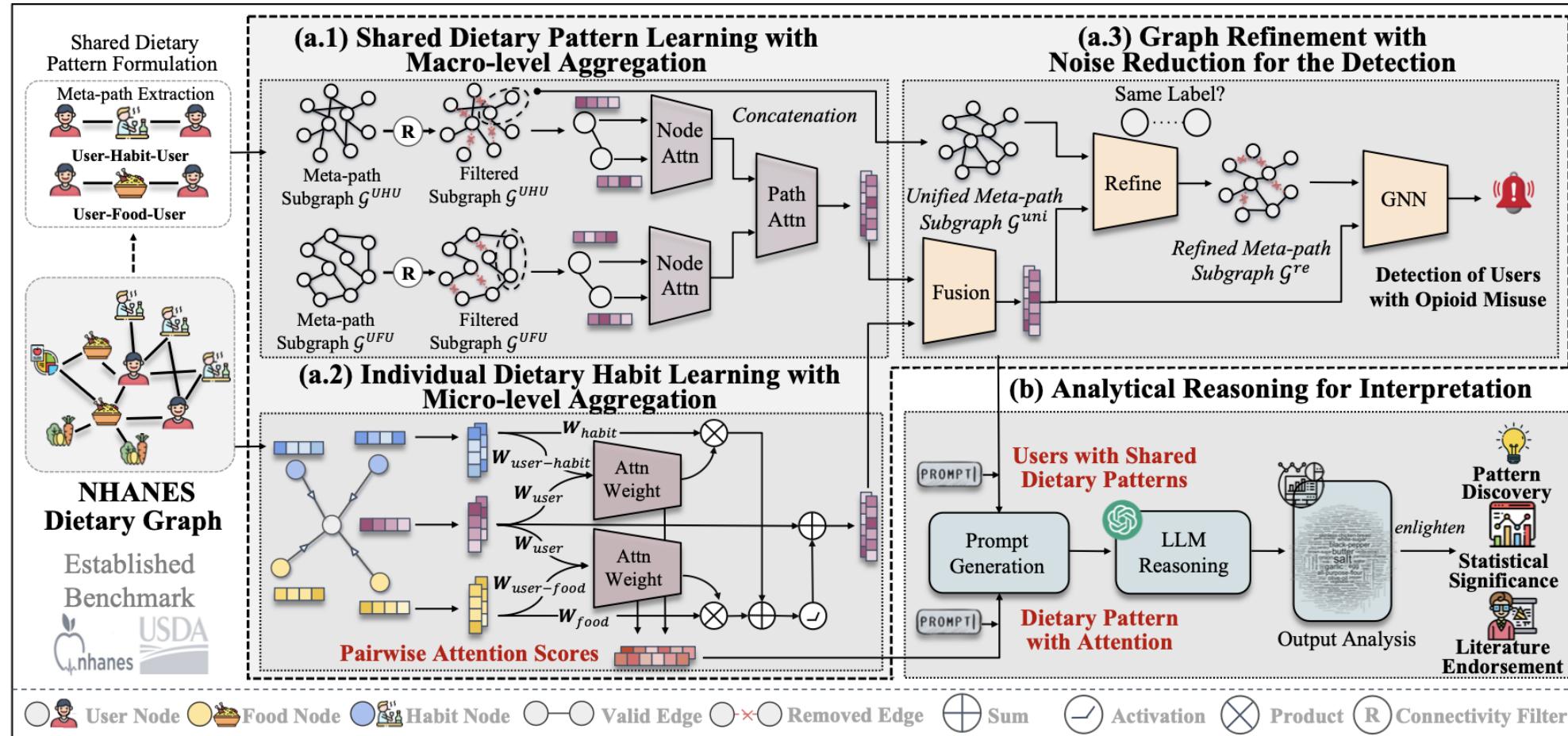
- LLM as Data Augmentation for GNNs



LLM-Empowered Class Imbalanced Graph Prompt Learning for Online Drug Trafficking Detection, ACL 2025

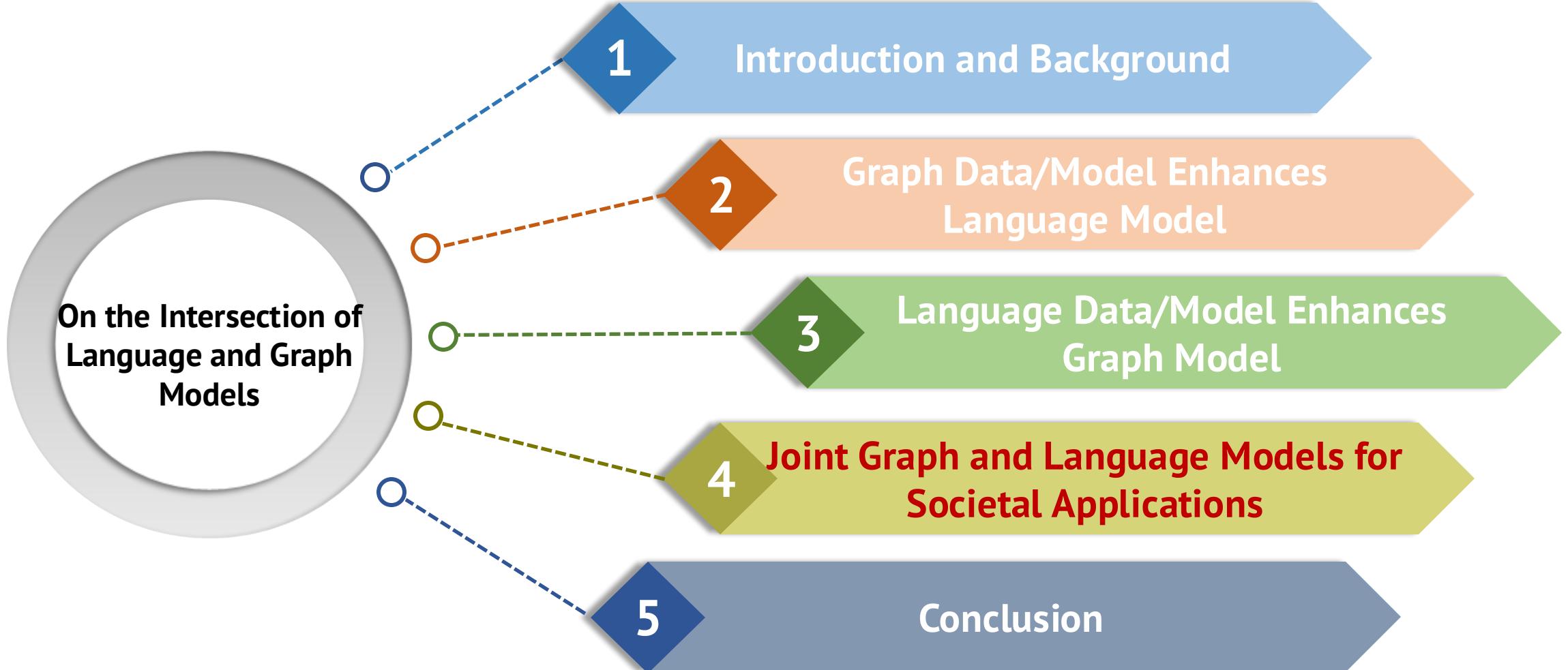
# Language Data/Model Enhances Graph Model

- Healthcare: LLM as Interpretator for GNN Output



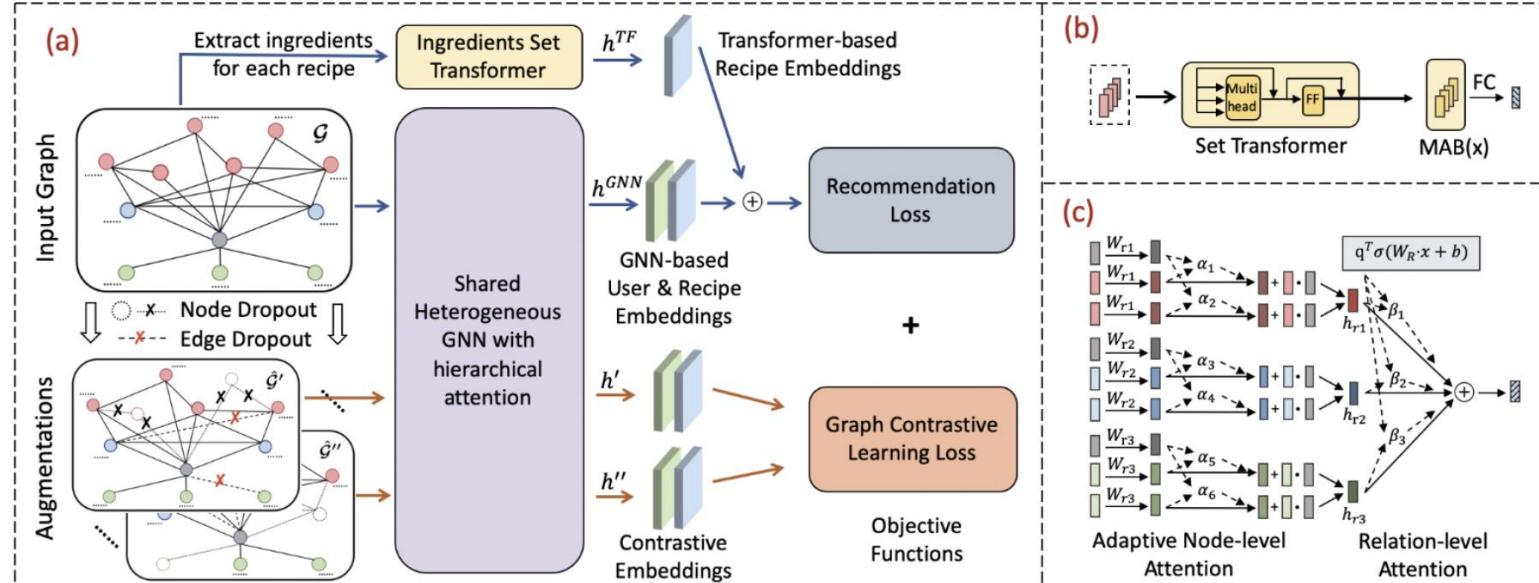
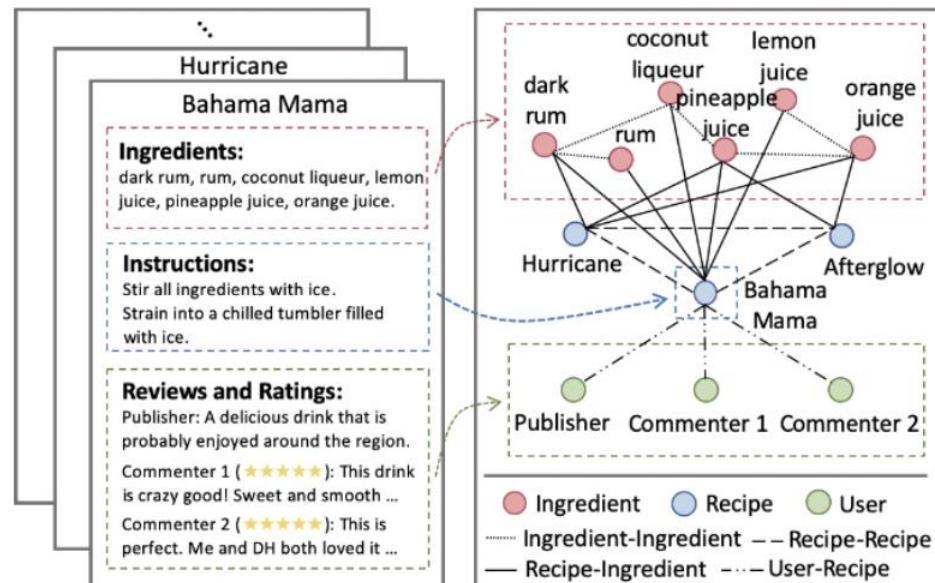
Diet-ODIN: A Novel Framework for Opioid Misuse Detection with Interpretable Dietary Patterns, KDD 2024

# Outline



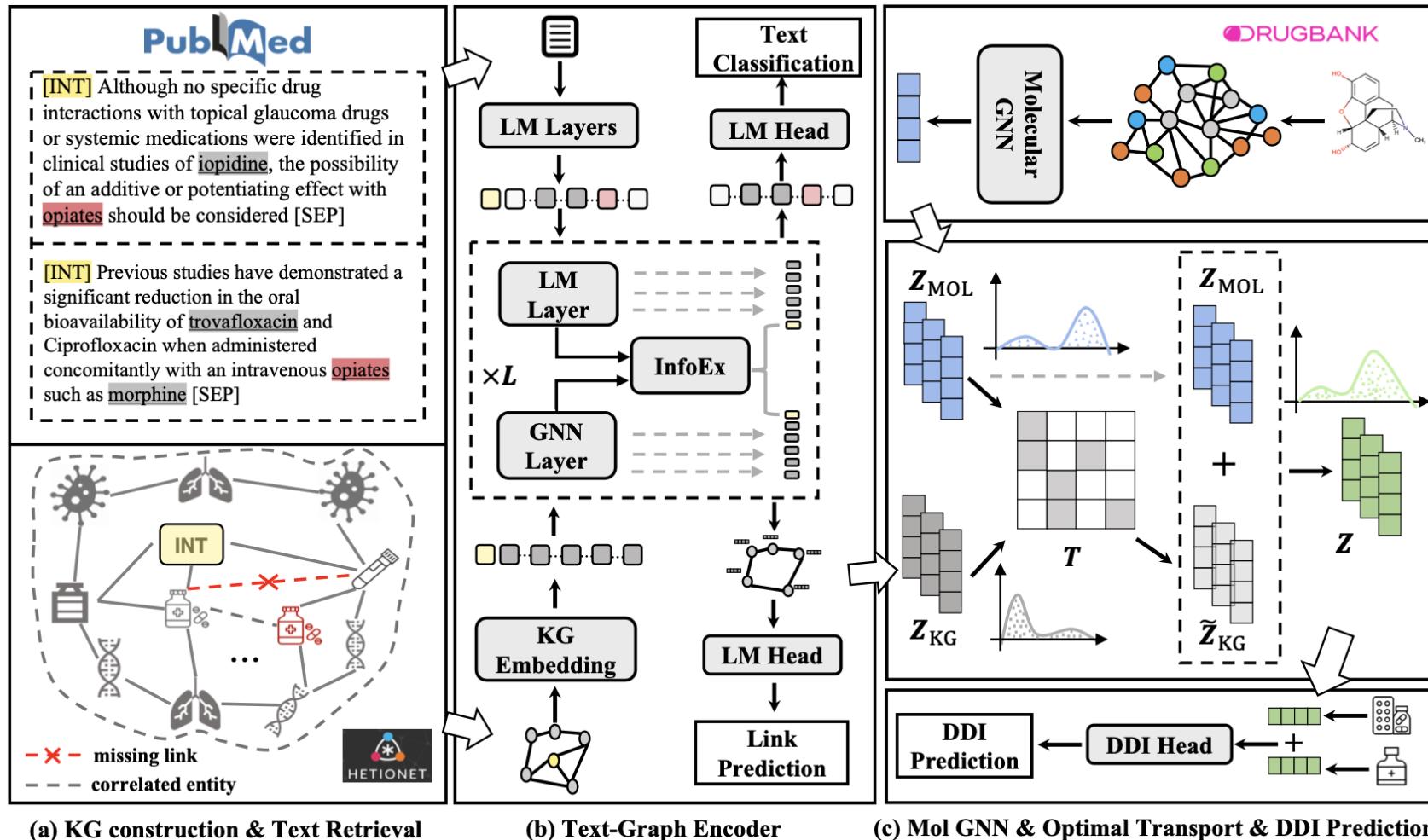
# Joint Graph and Language Models for Various Applications

- Recommender System: GNN and LLM as Multi-modal Data Encoder



# Joint Graph and Language Models for Various Applications

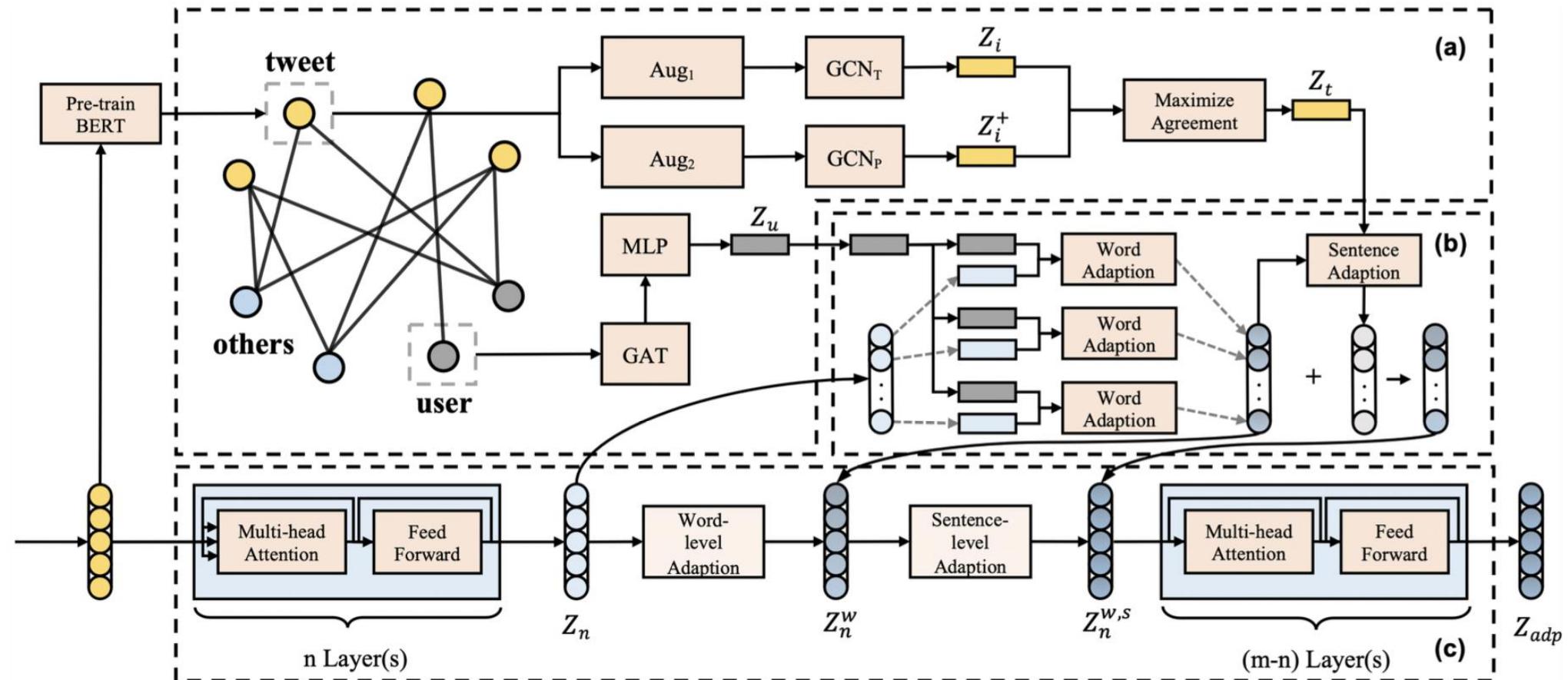
- Healthcare: GNN and LLM as Multi-modal Data Encoder



A Multi-Modality Framework for Drug-Drug Interaction Prediction by Harnessing Multi-source Data, CIKM 2023

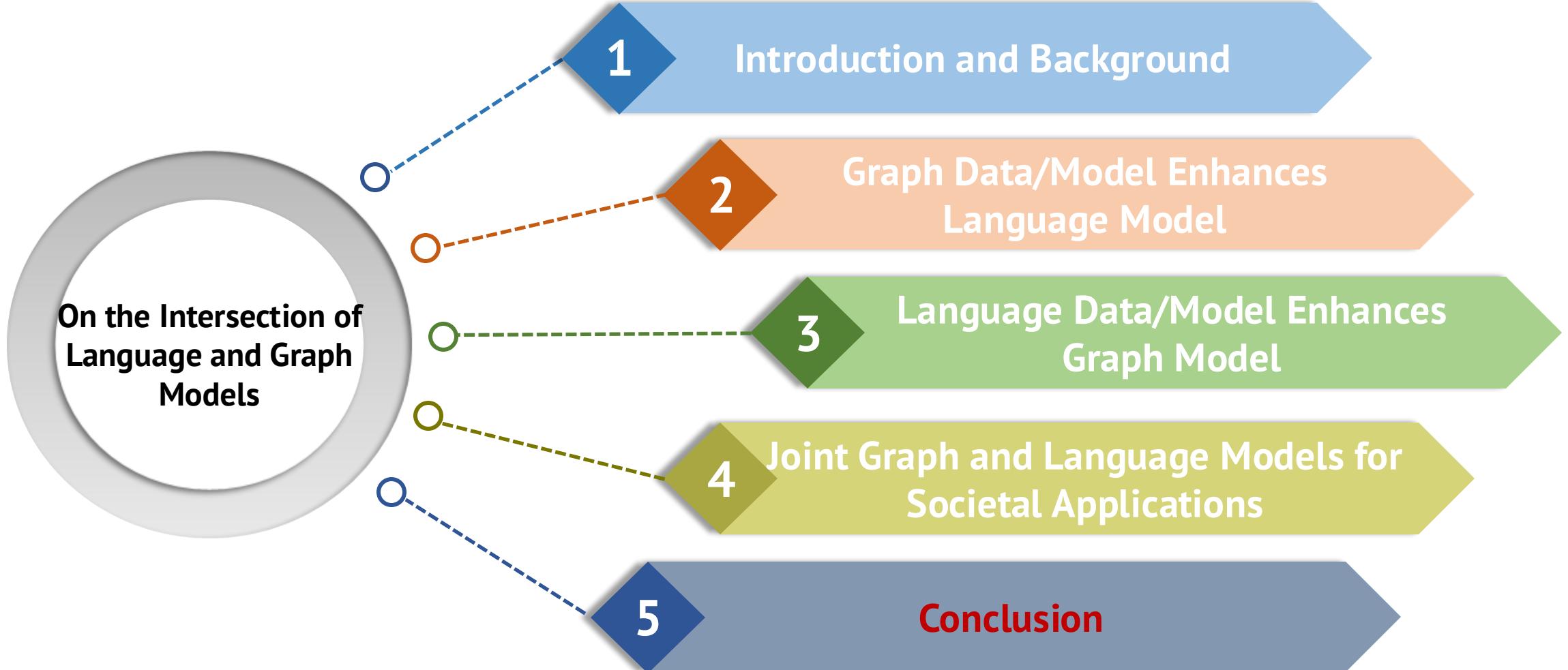
# Joint Graph and Language Models for Various Applications

- Social Network Analysis: GNN and LLM as Multi-modal Data Encoder



GraphBERT: Bridging Graph and Text for Malicious Behavior Detection on Social Media, ICDM 2022

# Outline



# Conclusion

---

- **Graph Data/Model Improves Language Model**
  - ❖ MASS: LLM for Advanced Reasoning
  - ❖ Graph as Augmented Information/Data for LLMs in QA, Text Generation, etc.
- **Language Data/Model Improves Graph Model**
  - ❖ GP2M: Graph Foundation Model
  - ❖ LLM as Text/Attribute Generation or Data Augmentation for GNNs, etc.
- **Joint Graph-Language Model for Societal Applications**
  - ❖ Recommender Systems, Social Network Analysis, Healthcare, etc.

# Q&A

- Feel free to contact me for any questions!

✉ Contact email: [chuxuzhang@gmail.com](mailto:chuxuzhang@gmail.com)

