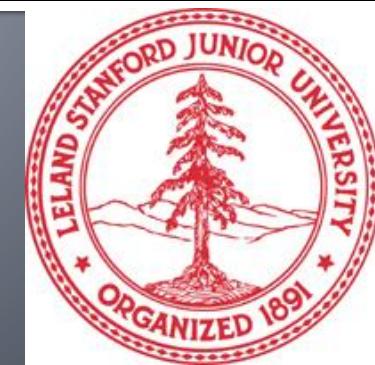


Note to other teachers and users of these slides: We would be delighted if you found our material useful for giving your own lectures. Feel free to use these slides verbatim, or to modify them to fit your own needs. If you make use of a significant portion of these slides in your own lecture, please include this message, or a link to our web site: <http://cs224w.Stanford.edu>

Stanford CS224W: Advanced architectures in Relational Deep Learning

CS224W: Machine Learning with Graphs
Charilaos Kanatsoulis, Stanford University
<http://cs224w.stanford.edu>

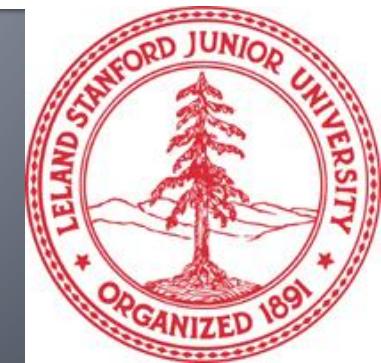


Announcements

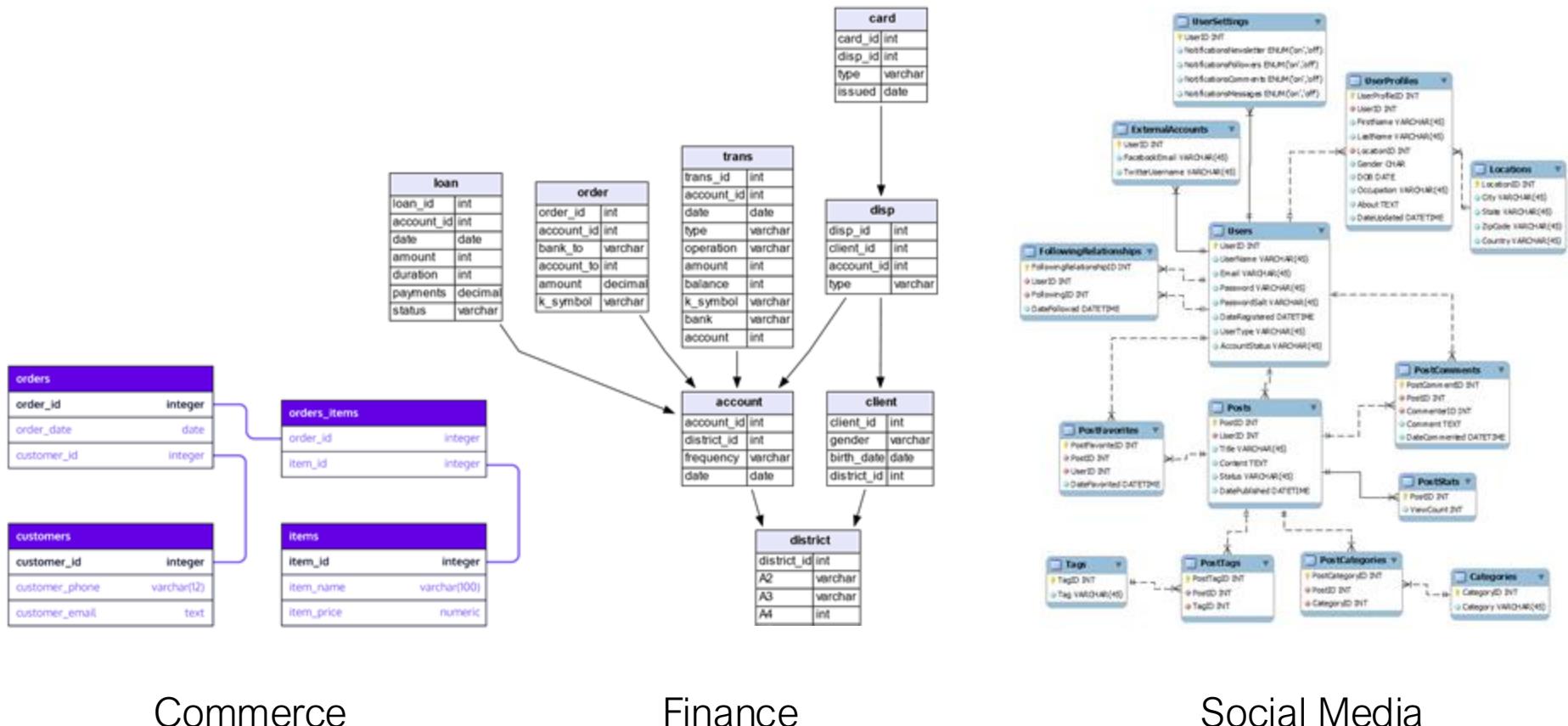
- **Colab 3** is due today
- **Homework 3** is due next Thursday 11/13
 - Recitation session recording on Ed
 - **Project Milestone** deadline has been pushed back: now due Tuesday 11/11
 - **Colab 4** will be released today
 - **Homework 2 grades will be released today!**
 - Regrade requests open until 11/13
 - **Practice exam will be released today (on Ed)**

Stanford CS224W: **Relational Deep Learning**

CS224W: Machine Learning with Graphs
Charilaos Kanatsoulis, Stanford University
<http://cs224w.stanford.edu>



Data stored in Relational Databases



Commerce

Finance

Social Media

Predictions on Relational Data

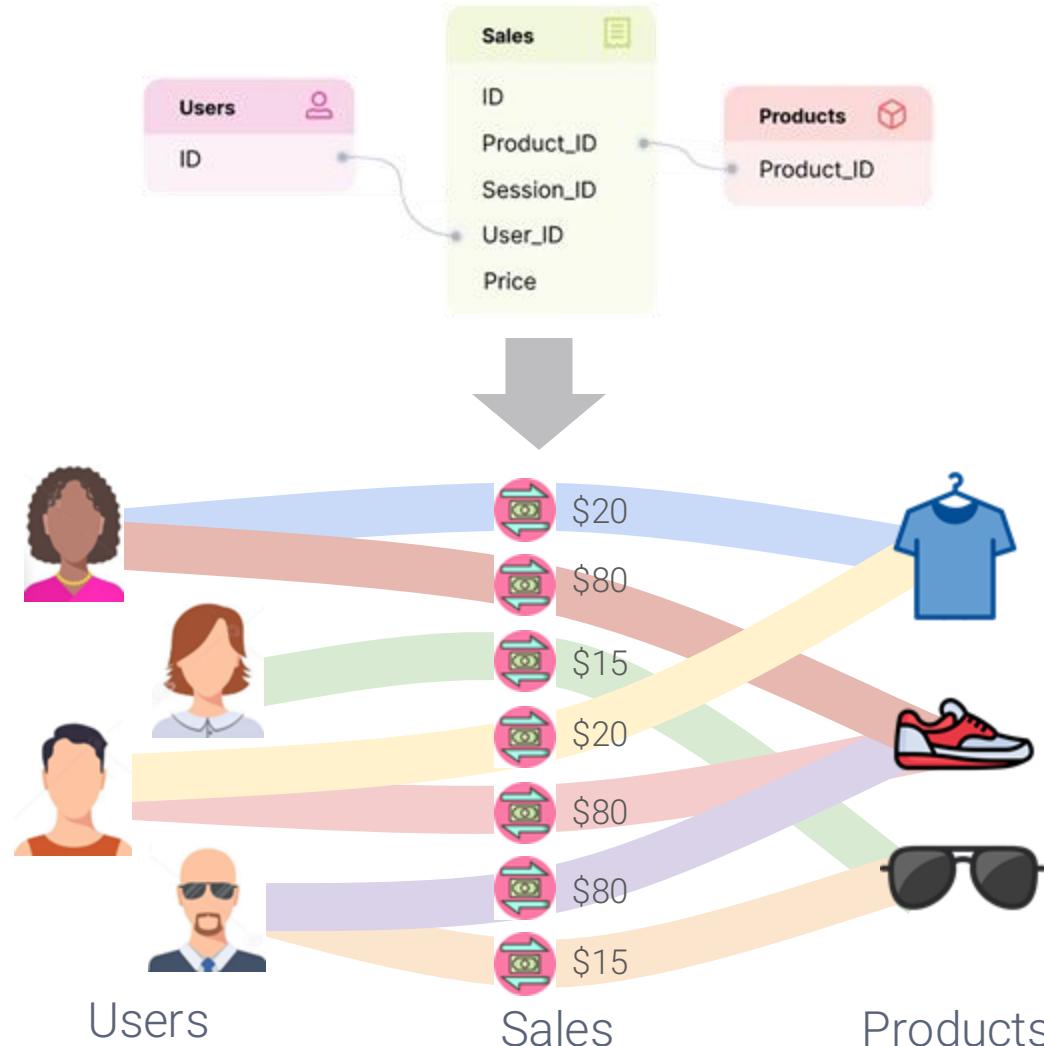
- + Which products will a user purchase in the next 7 days?
- + Will an active user churn in the next 90 days?
- + What will be the total sales for each product in the next 30 days?



Insight: A Data is a graph!



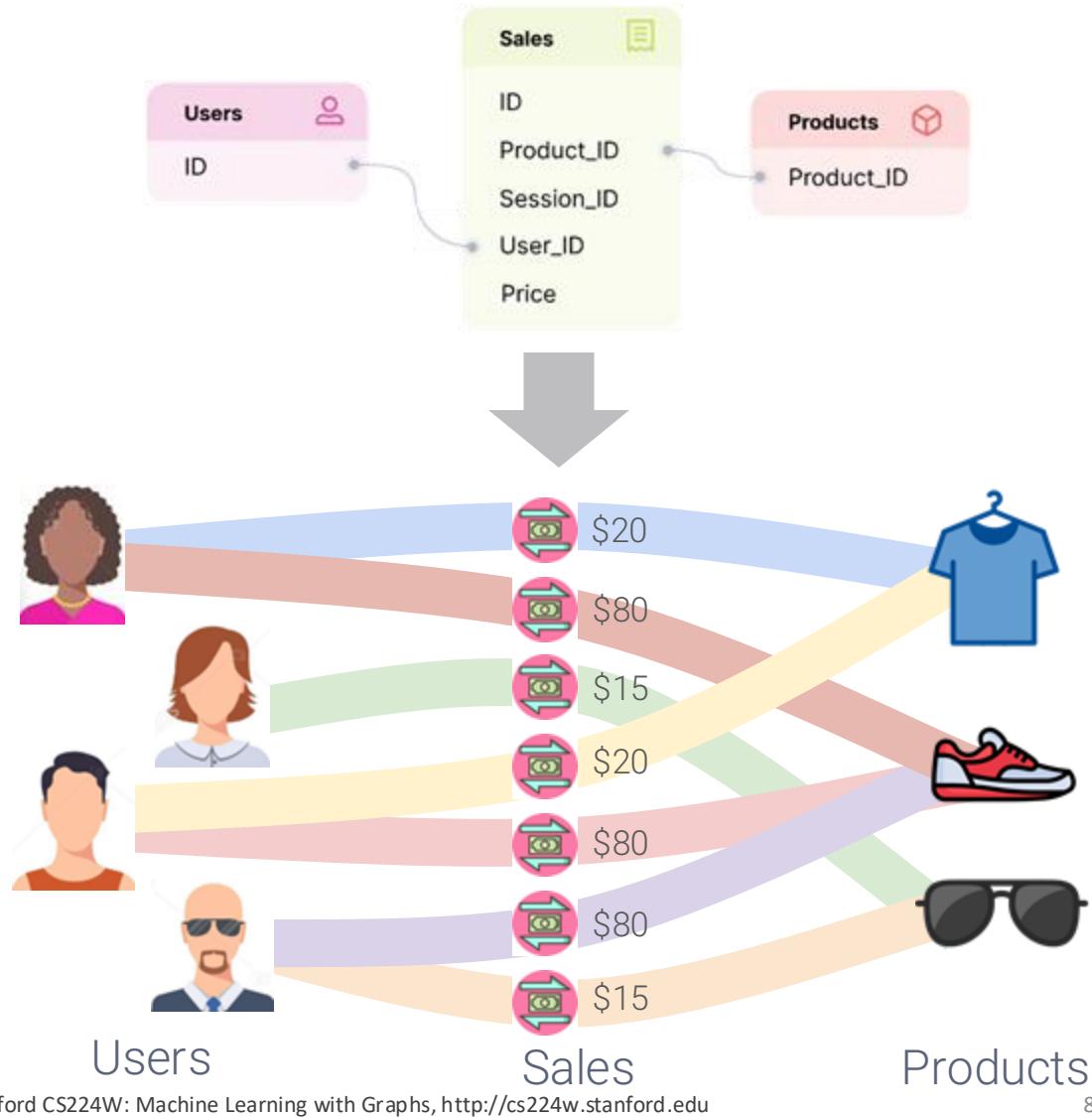
A Database is a graph!



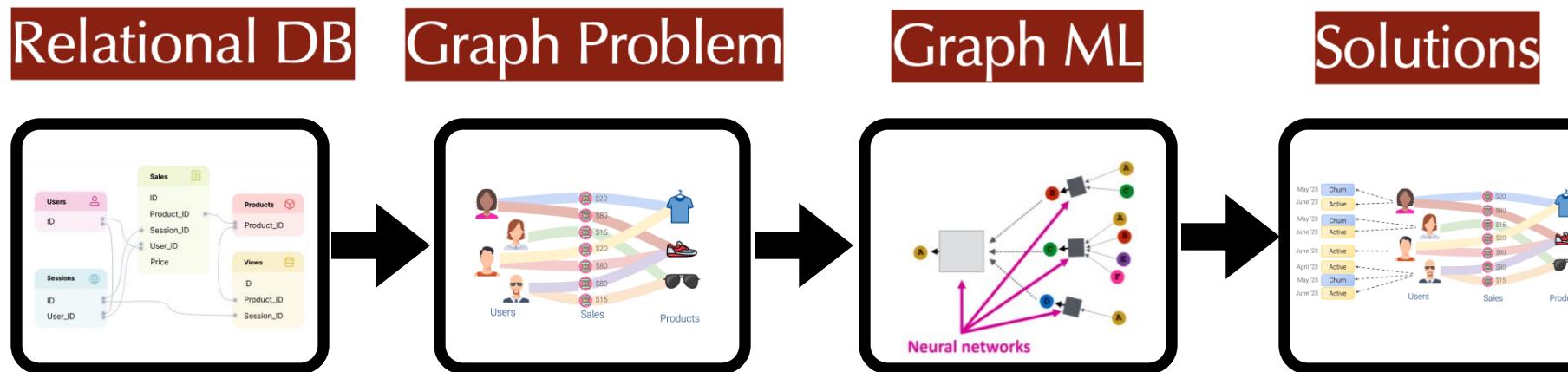
Just do ML on a Graph!

ML in the language of graphs:

- **Node-level:**
 - Churn
 - Life-time value
 - Next best action
- **Link-level:**
 - Product affinity
 - Recommendations
- **Graph-level:**
 - Fraud, money laundering



Graph ML Problem Solving Pipeline



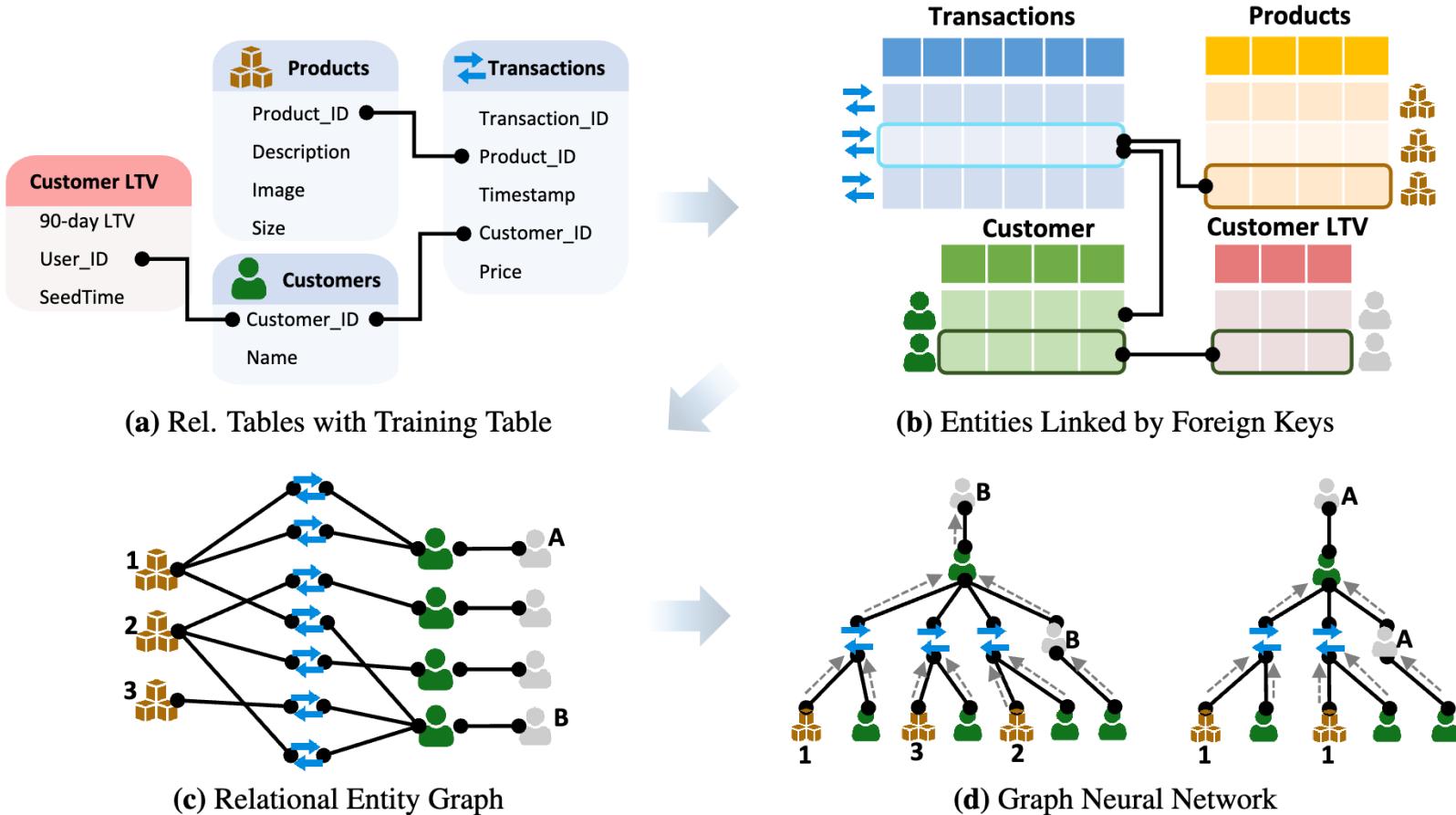
Defining a Task: Training Table

- **Training Table:** A special table containing training labels
 - (entity ID, time, labels)
 - Classification, Regression, Multi-class

Training Table

Entity ID	Timestamp	Label
99	10172024	1
99	10182024	1
...
100	10172024	1
100	10182024	0
...

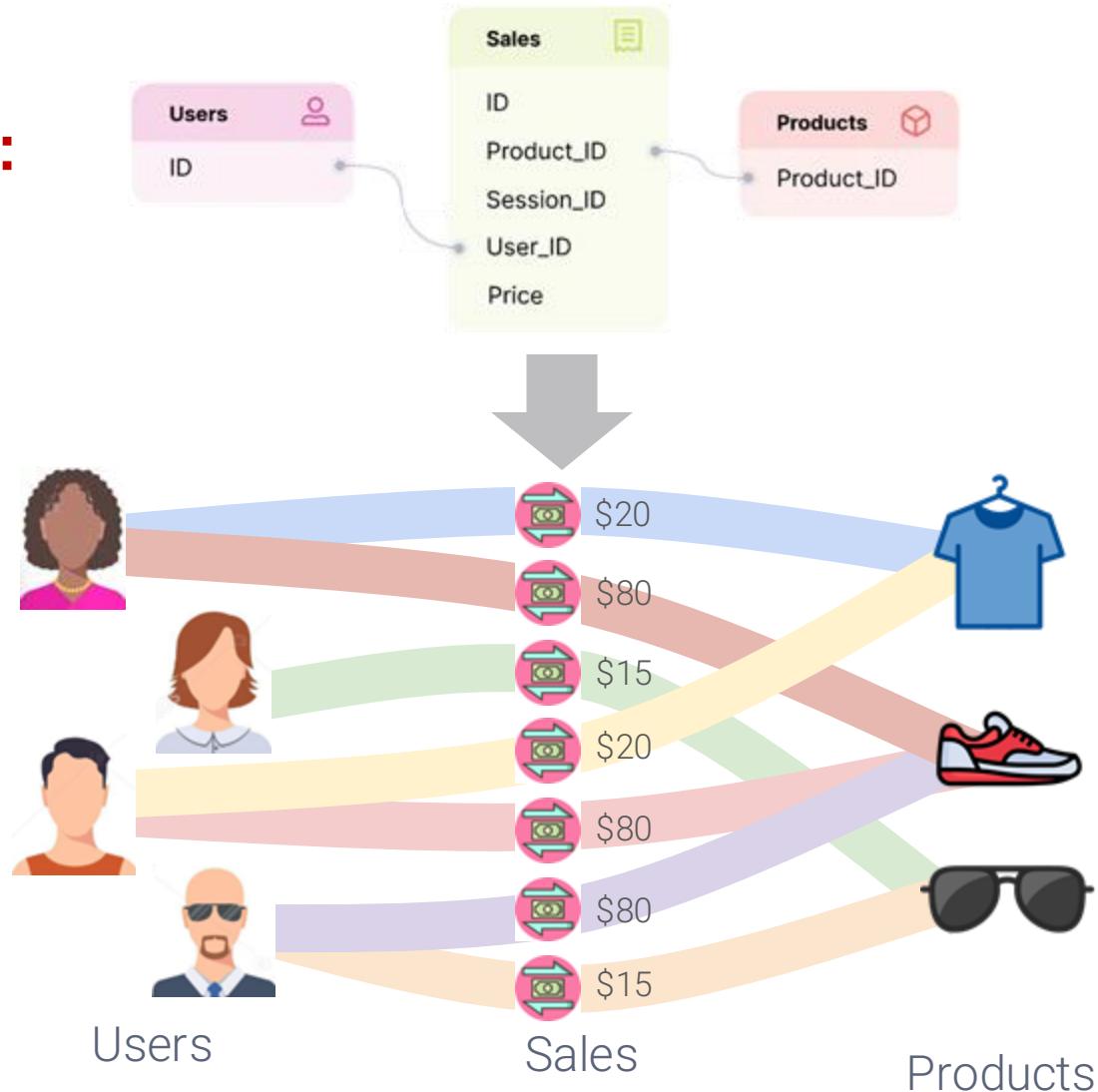
Relational Deep Learning



Relational Entity Graph

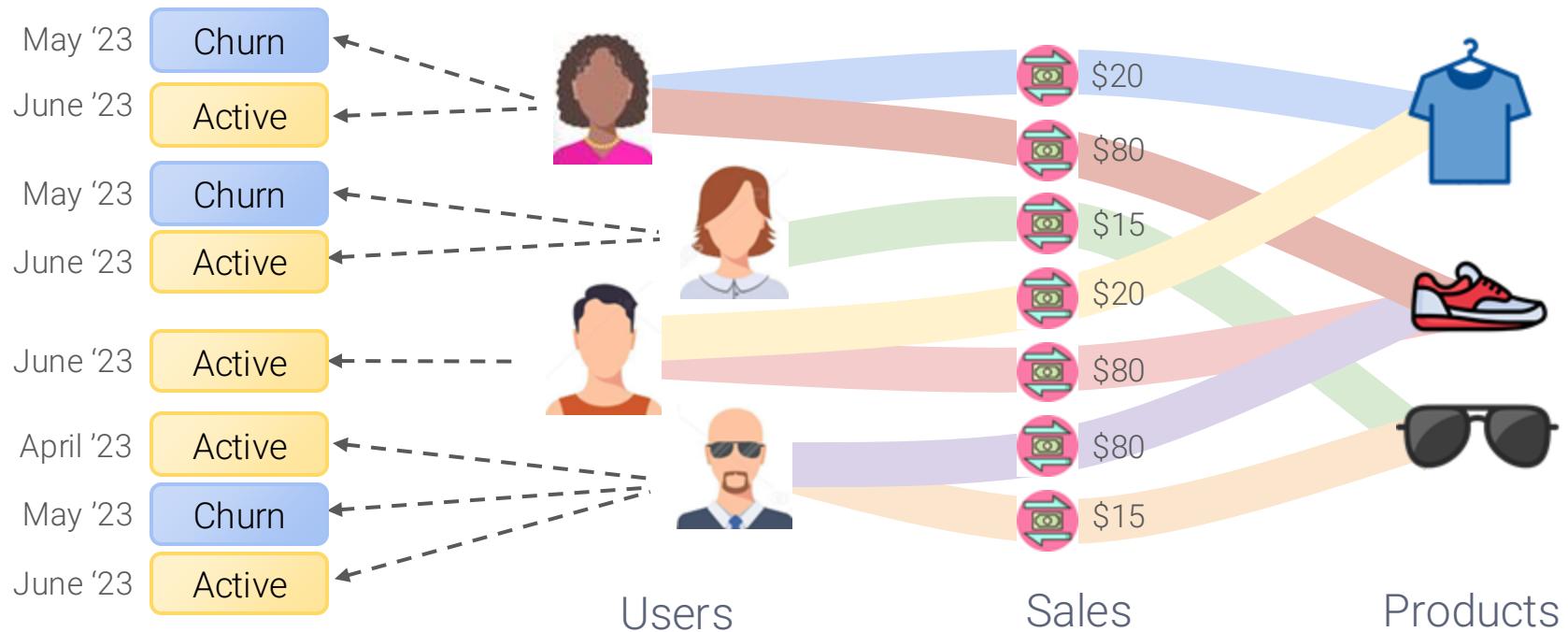
Relational Entity Graph:

Create connections via primary-foreign keys



Connect the Training Table

Training labels together with **timestamps** are attached to the graph



GNN on the Entity Graph

Node's neighborhood defines a computation graph

Nodes learn how to *optimally* use information from neighbors to obtain enhanced node representations



Why not use an LLM?

Some have tried

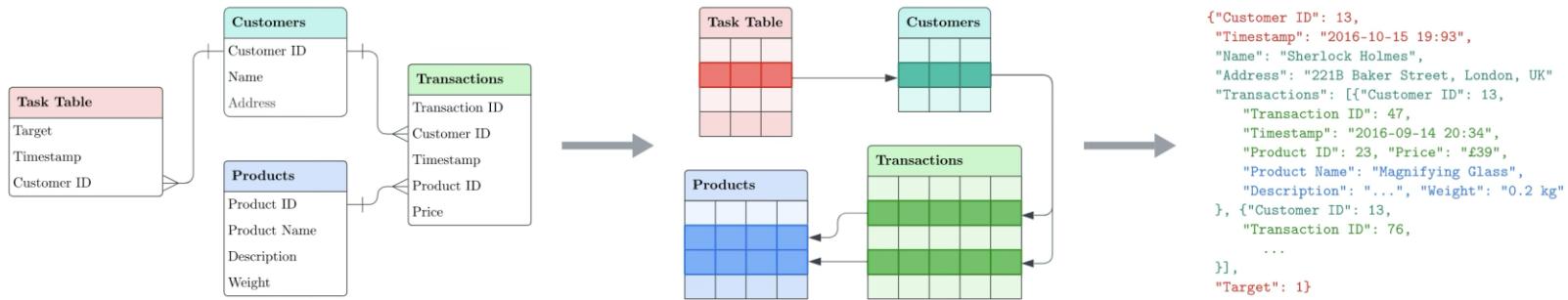


Figure 1: Process of constructing a single example for LLM-based inference.

Tackling prediction tasks in relational databases with LLMs

Marek Wydmuch^{*,†}

Łukasz Borchmann*

Filip Graliński^{*,‡}

* Snowflake AI Research

{first-name}.{last-name}@snowflake.com

† Poznań University of Technology / Poznań, Poland

‡ Adam Mickiewicz University / Poznań, Poland

Prompt

Why not LLMs?

- LLMs are not trained to understand and effectively learn from the **relational nature** of Databases
- Database prediction is **NOT** sequence modeling
- LLMs are **NOT** trained to predict the **future**
- **Context size** is **limited**

RDL gains over LLM

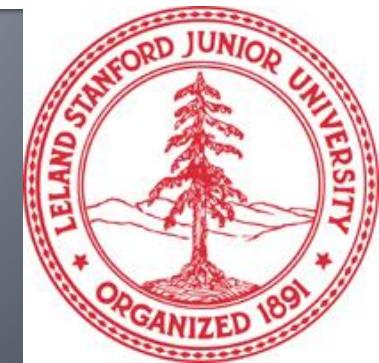
Entity classification results (ROC-AUC(%), higher is better) on RELBENCH test set. Best values are in bold.

Dataset	Task	ICL Llama 3.2 3B	RELGNN(ours)	Relative Gain
rel-amazon	user-churn	62.55	70.99	13%
	item-churn	73.41	82.64	13%
rel-avito	user-visits	53.36	66.18	24%
	user-clicks	54.07	68.23	26%
rel-event	user-repeat	70.11	79.61	14%
	user-ignore	68.65	86.18	26%
rel-f1	driver-dnf	80.03	75.29	-6%
	driver-top3	87.11	85.69	-2%
rel-hm	user-churn	63.81	70.93	11%
rel-stack	user-engagement	81.23	90.75	12%
	user-badge	79.99	88.98	11%
rel-trial	study-outcome	59.17	71.24	20%

▣ Stay tuned for more discussion on LLMs later!

Stanford CS224W: Next-Generation Architectures for RDL

CS224W: Machine Learning with Graphs
Charilaos Kanatsoulis, Stanford University
<http://cs224w.stanford.edu>



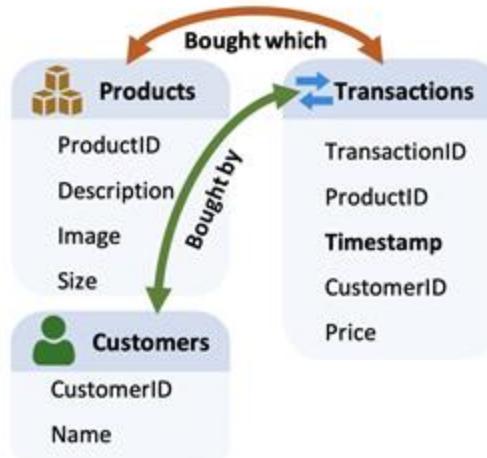
The need for specialized models

Relational Databases: **multi-modal, multi-scale**
information

The need for specialized models

Relational Databases: **multi-modal, multi-scale** information

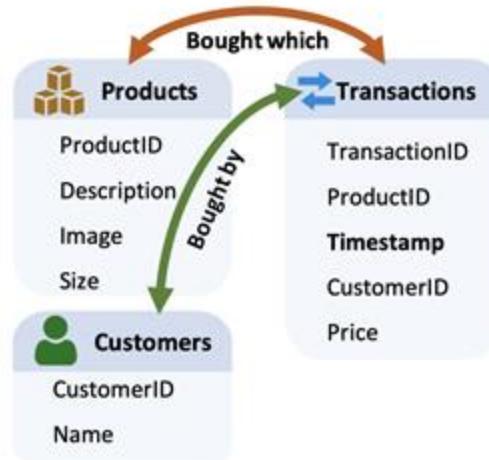
Entity level: Text, image, numerical, categorical, timestamp



The need for specialized models

Relational Databases: **multi-modal, multi-scale** information

Entity level: Text, image, numerical, categorical, timestamp



Database level: Relational structure, Temporal Structure

Relational DBs structure

Tripartite patterns



Cycle patterns



Structure changes over time!



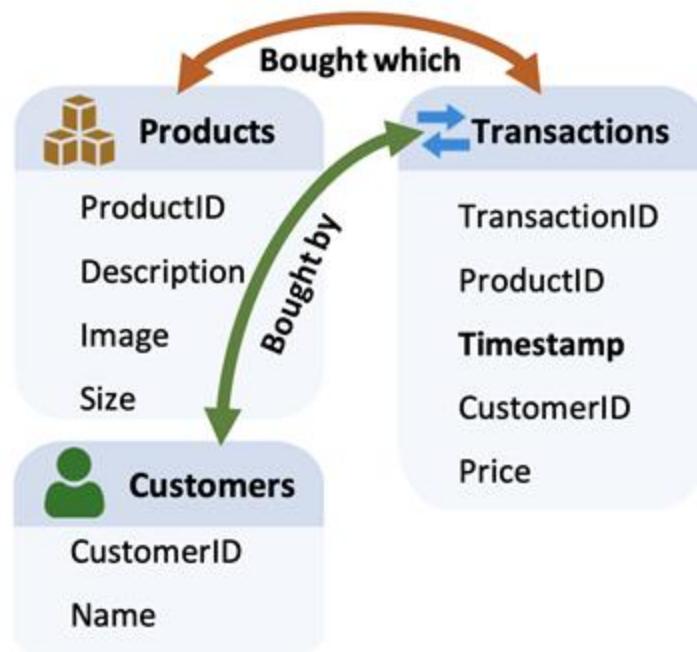
Stanford CS224W: RelGNN: Composite Message-Passing for RDL

CS224W: Machine Learning with Graphs
Charilaos Kanatsoulis, Stanford University

<http://cs224w.stanford.edu>

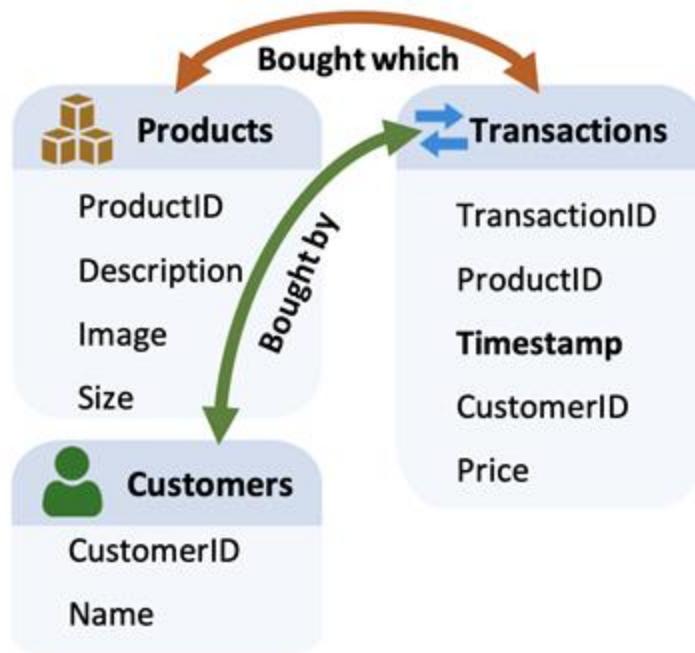


Relational DBs structure



standard heterogeneous graph

Relational DBs structure



standard heterogeneous graph

relational entity graph

Observation

- Edges in relational entity graphs are **defined by primary–foreign key links** that merely record table connectivity **without semantics**.



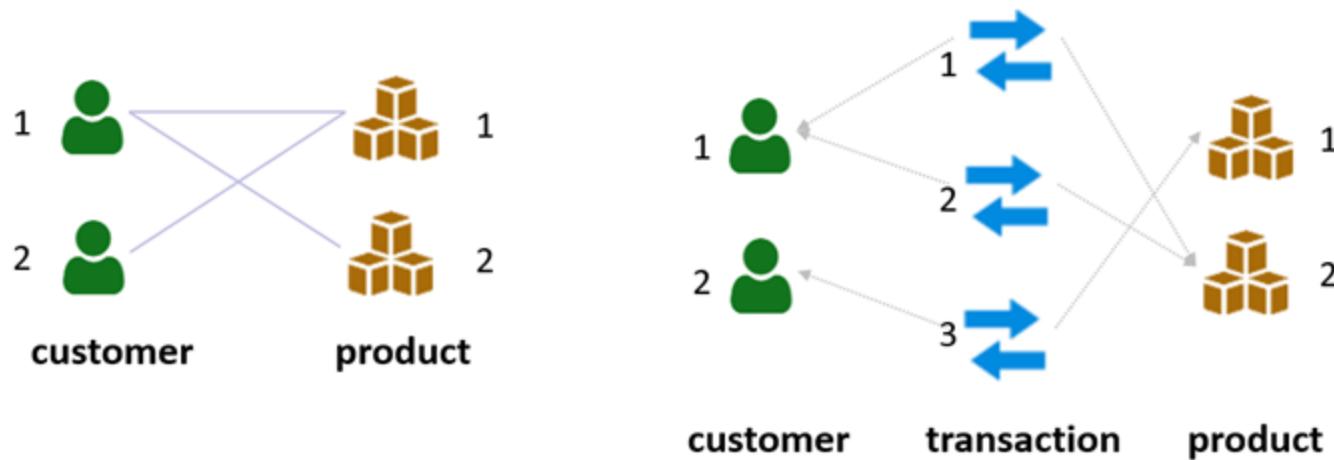
standard heterogeneous graph



relational entity graph

Observation

- Edges in relational entity graphs are **defined by primary–foreign key links** that merely record table connectivity **without semantics**.
- Junction tables are introduced, decomposing each many-to-many association into **a pair of one-to-many links**.





Stanford CS224W: RELBENCH

CS224W: Machine Learning with Graphs

Jure Leskovec, Stanford University

<http://cs224w.stanford.edu>



RelBench Datasets

7 Diverse Datasets

E-Commerce



- rel-amazon
- rel-avito
- rel-hm



Social

- rel-event
- rel-stack

Sports



- rel-f1



Medical

- rel-trial

Junction Tables in Schema Graph

rel-amazon



rel-hm



Junction Tables in Schema Graph

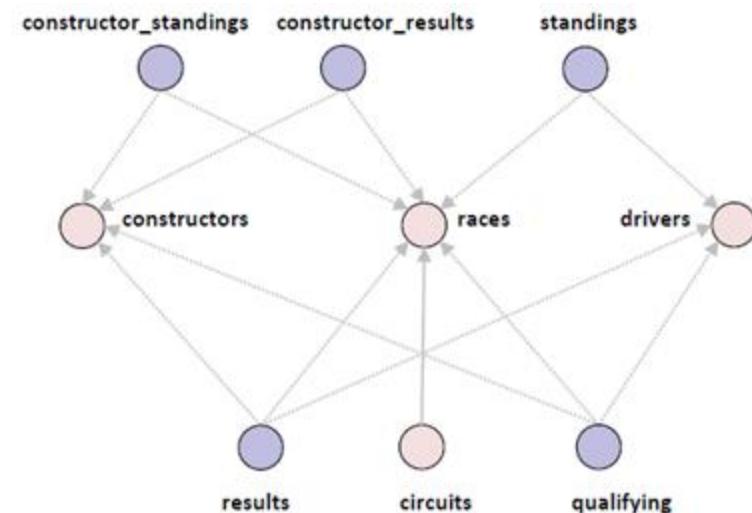
rel-amazon



rel-hm



rel-F1

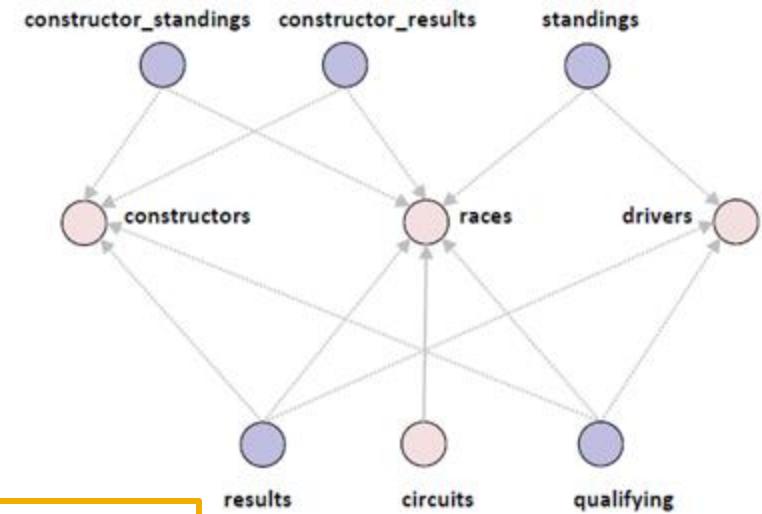


Junction Tables in Schema Graph

rel-amazon



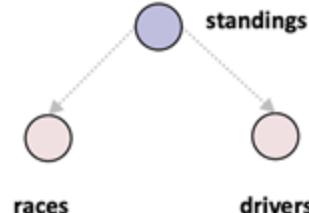
rel-F1



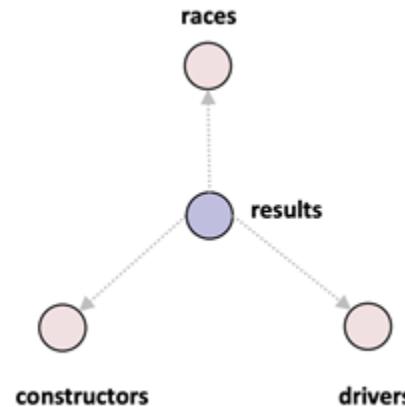
rel-hm



bridge structure

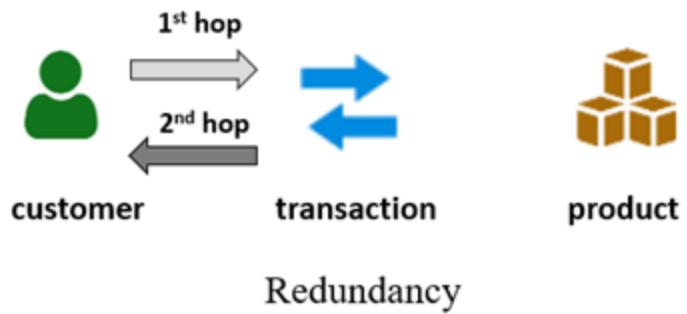
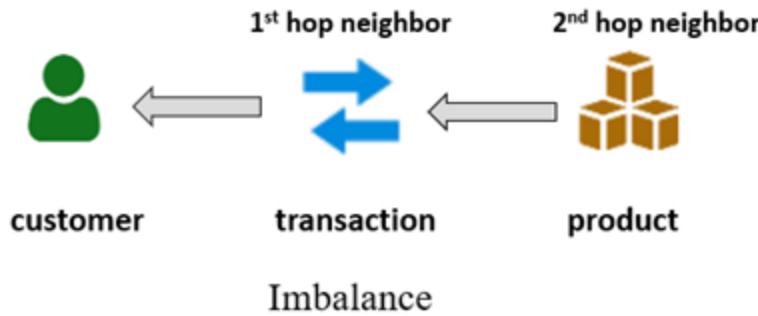


hub structure



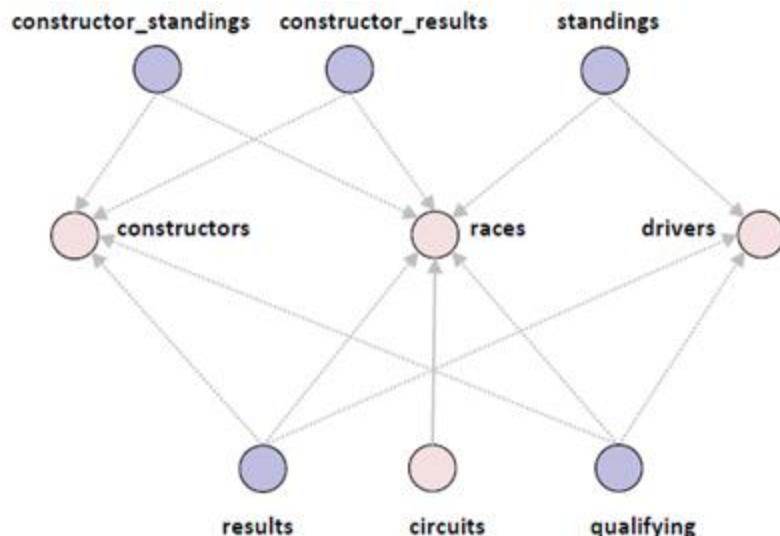
Limitation

- ▣ (1) *Imbalance*: junction node are 1-hop neighbors and **more informative node** become 2-hop.
- ▣ (2) *Redundancy*: information from the source node is passed to the junction node in the first hop and routed back to itself in the second, aggregated with the junction and 2-hop node information.

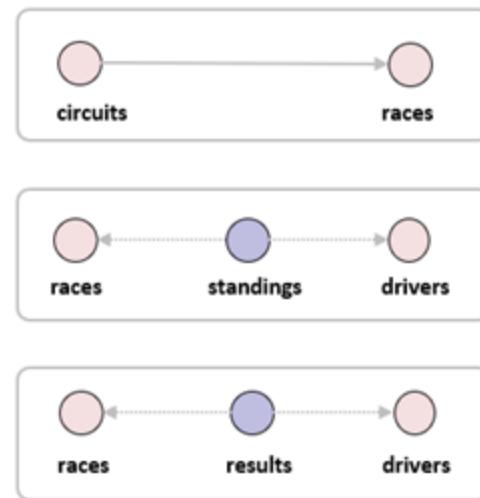


RelGNN: Atomic Routes

- Informal: An atomic route is a **simple path between pink node-types**.
- Atomic routes can be **derived automatically without manual intervention or domain-expert knowledge.**



Primary-foreign Key Relation of rel-fl Dataset



Example of Atomic Routes

RelGNN: Composite Message Passing

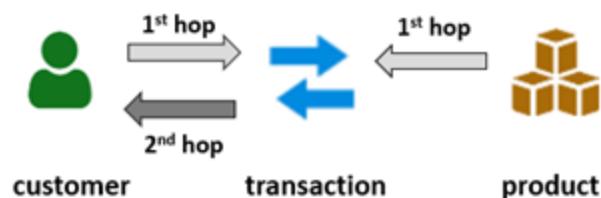
- **Composite Message Passing:** aggregates messages along atomic routes in a single step.

$$\mathbf{m}_{(\text{dst}, \text{mid}, \text{src})}^{(l+1)} = \text{AGGR}(\mathbf{h}_{\text{dst}}^{(l)}, \{\{\text{FUSE}(\mathbf{h}_{\text{mid}}^{(l)}, \mathbf{h}_{\text{src}}^{(l)})\}\})$$

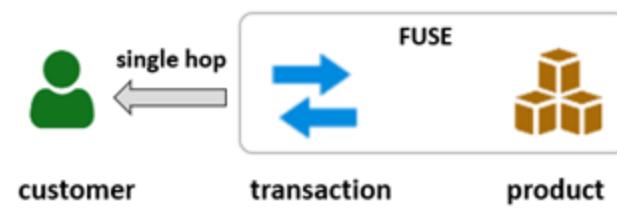
- Instantiation example:

$$\text{FUSE}(\mathbf{h}_{\text{mid}}^{(l)}, \mathbf{h}_{\text{src}}^{(l)}) = \mathbf{W}_1 \mathbf{h}_{\text{mid}}^{(l)} + \mathbf{W}_2 \mathbf{h}_{\text{src}}^{(l)}.$$

$$\text{AGGR}(\mathbf{h}_{\text{dst}}^{(l)}, \{\{\mathbf{h}_{\text{fuse}}^{(l)}\}\}) = \mathbf{W}_{\text{proj}} \mathbf{h}_{\text{dst}}^{(l)} + \sum_{\text{fuse} \in \mathcal{N}(\text{dst})} \alpha_{\text{dst}, \text{fuse}} \mathbf{W}_V \mathbf{h}_{\text{fuse}}^{(l)}$$



Standard Message Passing



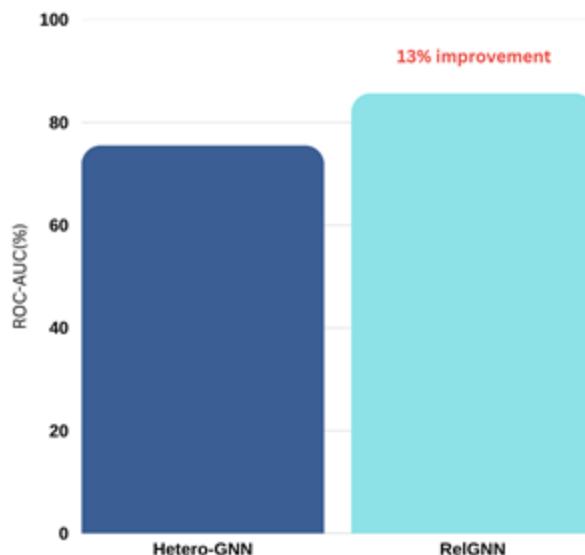
Composite Message Passing

ReIGNN results



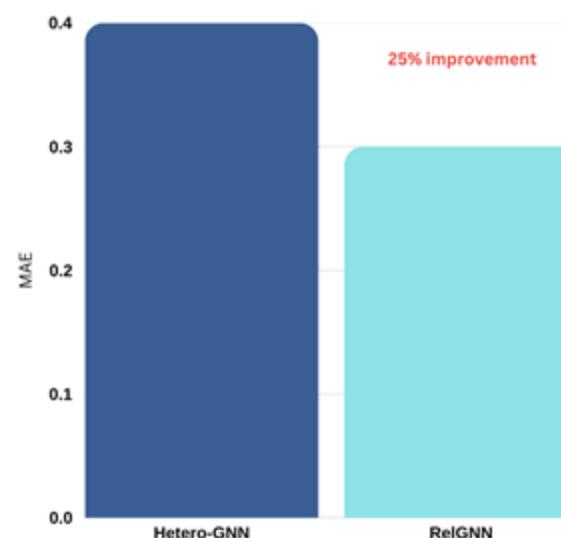
rel-F1

Predict if the driver will qualify in the top-3 for a race in the next 1 month.



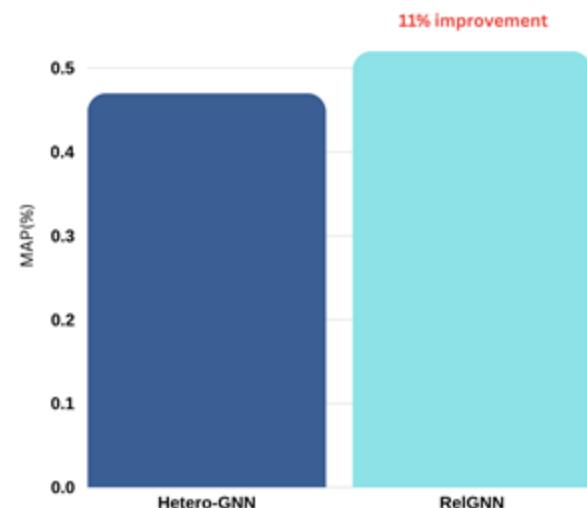
rel-trial

Predict the success rate of a trial site in the next 1 year.



rel-amazon

Predict the list of distinct items each customer will purchase and give a detailed review in the next 3 months.



■ ReIGNN achieves **SOTA performance** on the vast majority of tasks from RelBench, with improvements of up to 25%.

Stanford CS224W: RelGT: Relational Graph Transformer

CS224W: Machine Learning with Graphs

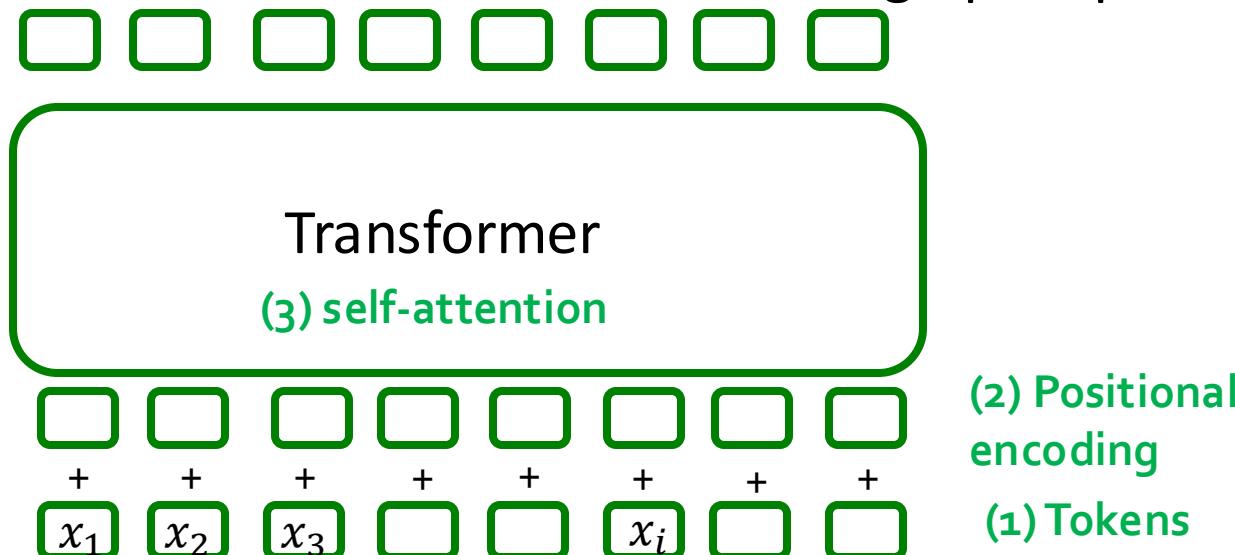
Jure Leskovec, Stanford University

<http://cs224w.stanford.edu>



Components of a Transformer

- Key components of Transformer
 - (1) tokenizing
 - (2) positional encoding
 - (3) self-attention
- **Key question:** What should these be for a graph input?



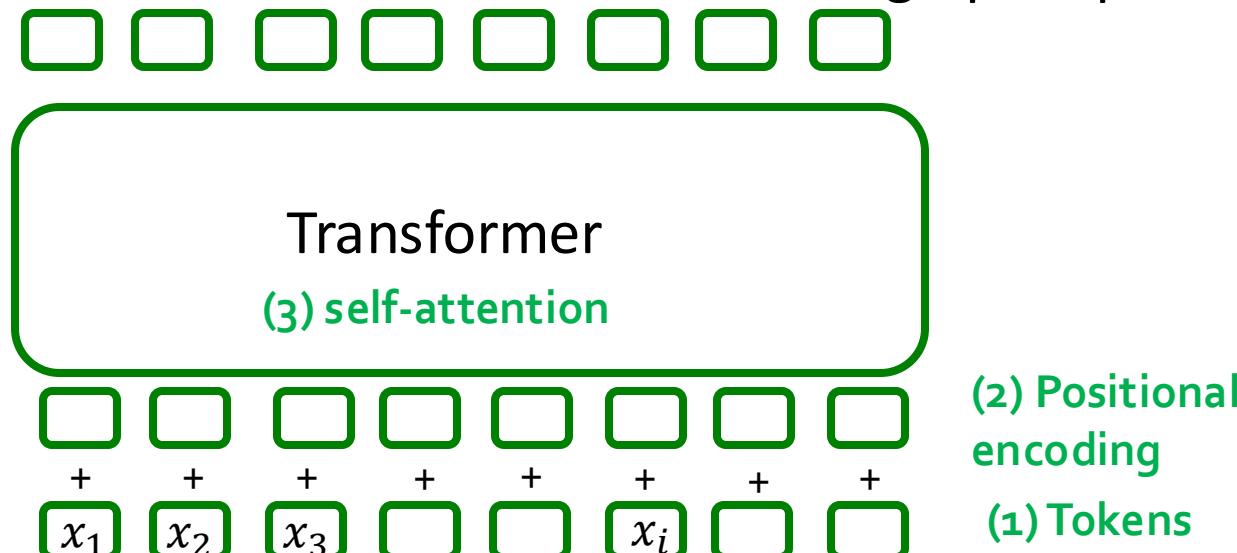
Components of a Transformer

- Key components of Transformer

- (1) tokenizing
- (2) positional encoding
- (3) self-attention

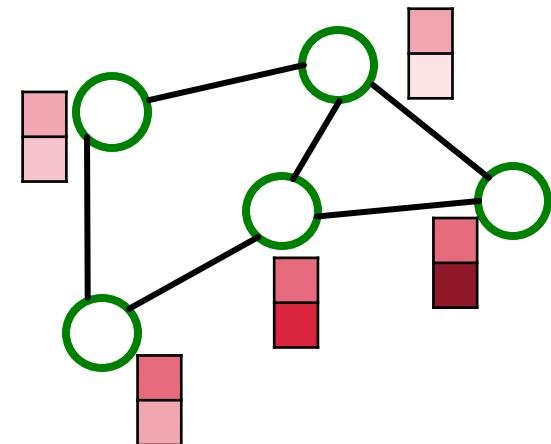
How to chose these
for graph data?

- Key question:** What should these be for a graph input?



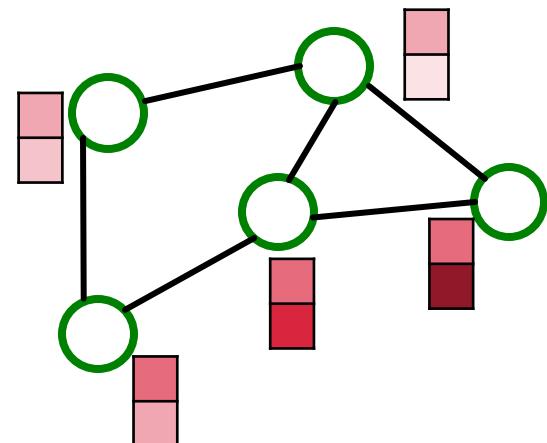
Processing Graphs with Transformers

- A graph Transformer must take the following inputs:
 - (1) Node features?
 - (2) Adjacency information?
 - (3) Edge features?
- Key components of Transformer
 - (1) tokenizing
 - (2) positional encoding
 - (3) self-attention



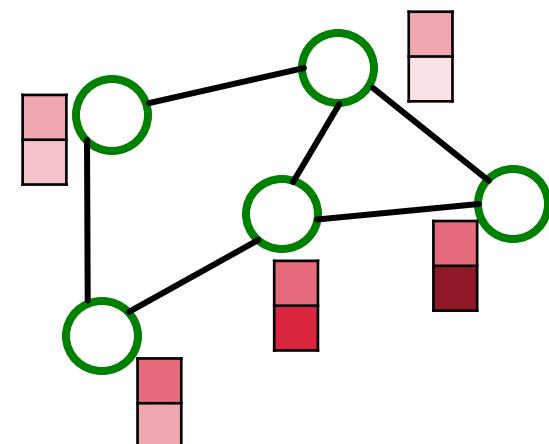
Processing Graphs with Transformers

- A graph Transformer must take the following inputs:
 - (1) Node features?
 - (2) Adjacency information?
 - (3) Edge features?
- Key components of Transformer
 - (1) tokenizing
 - (2) positional encoding
 - (3) self-attention
- There are many ways to do this
- Different approaches correspond to different “matchings” between graph inputs (1), (2), (3) transformer components (1), (2), (3)



Processing Graphs with Transformers

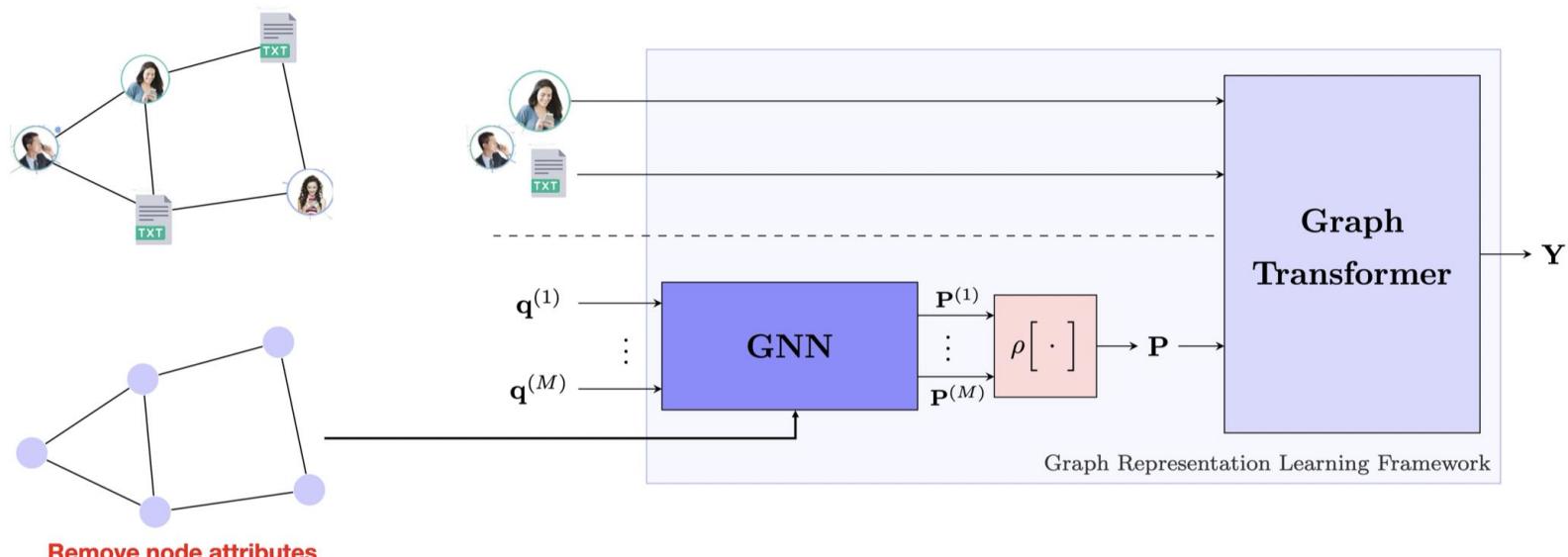
- A graph Transformer must take the following inputs:
 - (1) Node features?
 - (2) Adjacency information?
 - (3) Edge features?
- Key components of Transformer
 - (1) tokenizing
 - (2) positional encoding
 - (3) self-attention
- There are many ways to do this
- Different approaches correspond to different “matchings” between graph inputs (1), (2), (3) transformer components (1), (2), (3)



PEARL + Transformer

■ How to use PEARL: in practice?

- Step 1: Sample node ids from a probability distribution.
- Step 2: Process each set of node samples independently via a GNN.
- Step 3: Summarize the outputs via empirical expectation.
- Step 4: concatenate PEARL embeddings with node features X.
- Step 5: pass through main GNN/Transformer as usual.
- Step 6: Backpropagate gradients to train PEARL + Prediction model jointly.

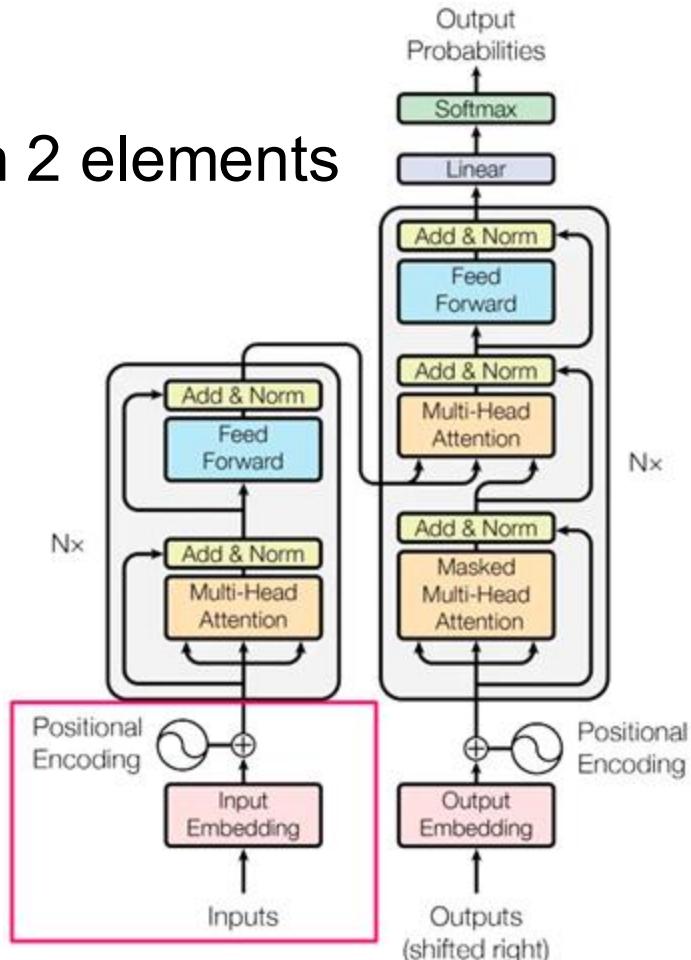


Observation

(Text) Transformers

- An input token is represented with 2 elements
 - o Token id, position

feature structure



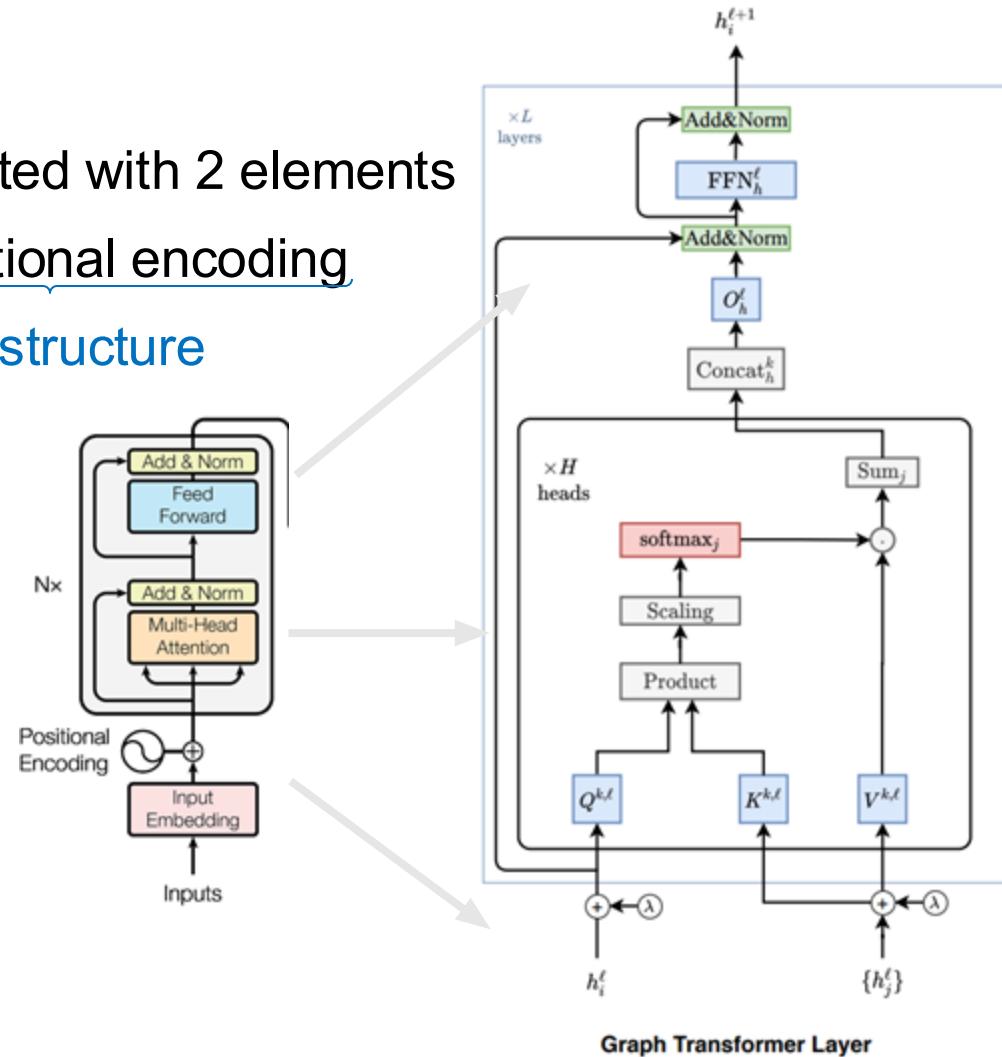
Observation

Graph Transformers

- An input token is represented with 2 elements
 - o Node feature, graph positional encoding

feature

structure



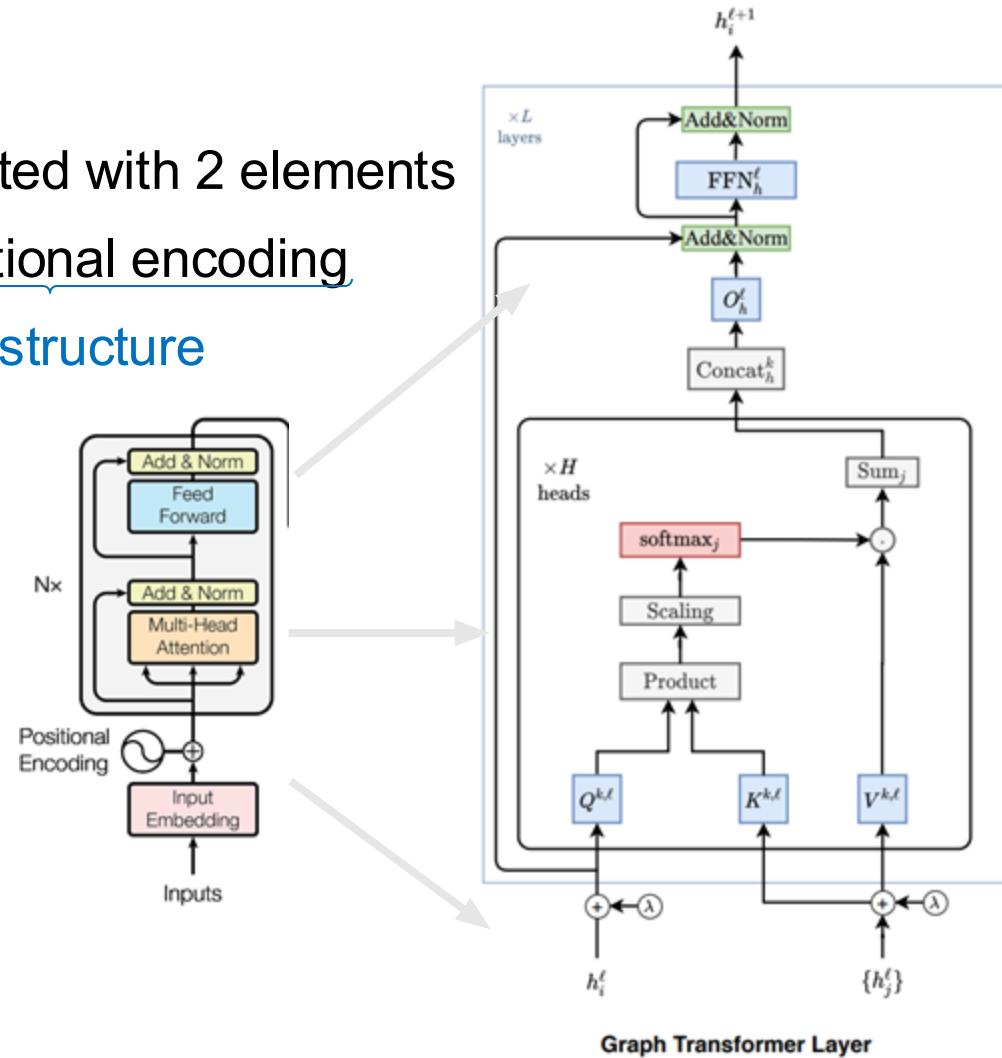
Observation

Graph Transformers

- An input token is represented with 2 elements
 - o Node feature, graph positional encoding

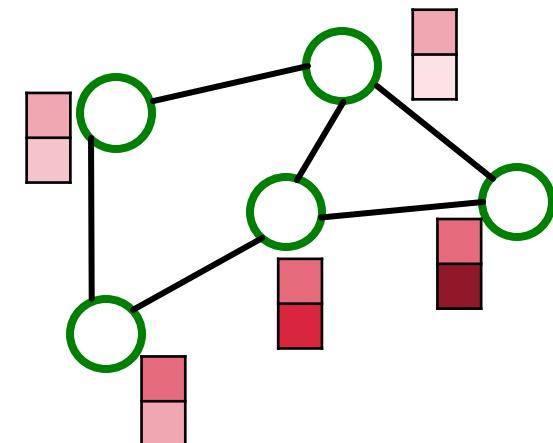
feature structure

- Most GTs are designed for homogeneous, static graphs
 - However, for additional RDL complexity such as **heterogeneity** and **temporality**, this design needs to be improved

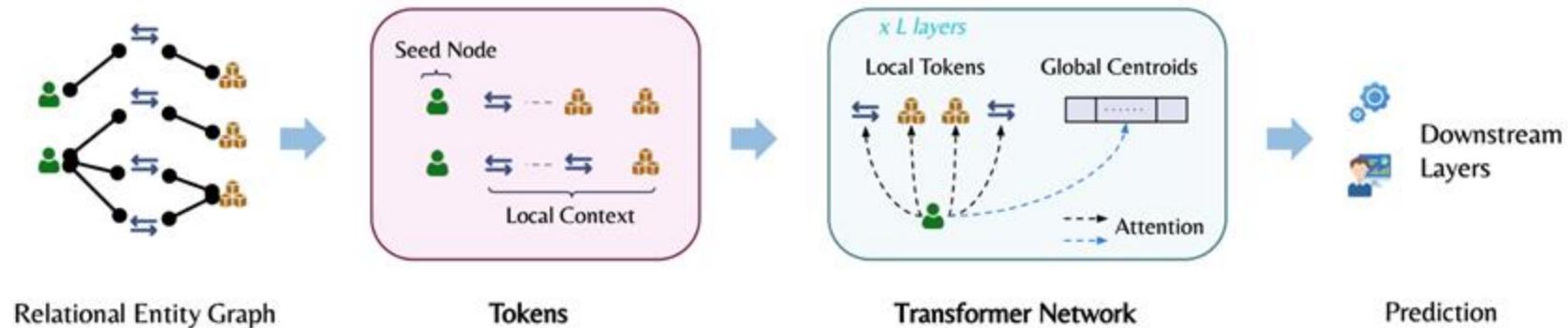


Processing Graphs with Transformers

- A graph Transformer must take the following inputs:
 - (1) Node features?
 - (2) Adjacency information?
 - ~~(3) Edge features?~~
 - (3) Node types
 - (4) Timestamp
- Key components of Transformer
 - (1) tokenizing
 - (2) positional encoding
 - (3) self-attention



Relational Graph Transformer



In Relational Graph Transformer

- We use multiple elements to represent the graph structure
 - Node feature, <multiple information elements>

Limitations of existing PEs for RDL

- Designed for homogeneous graphs
- No temporal dynamics
- Computational overhead
- Graph changes over time

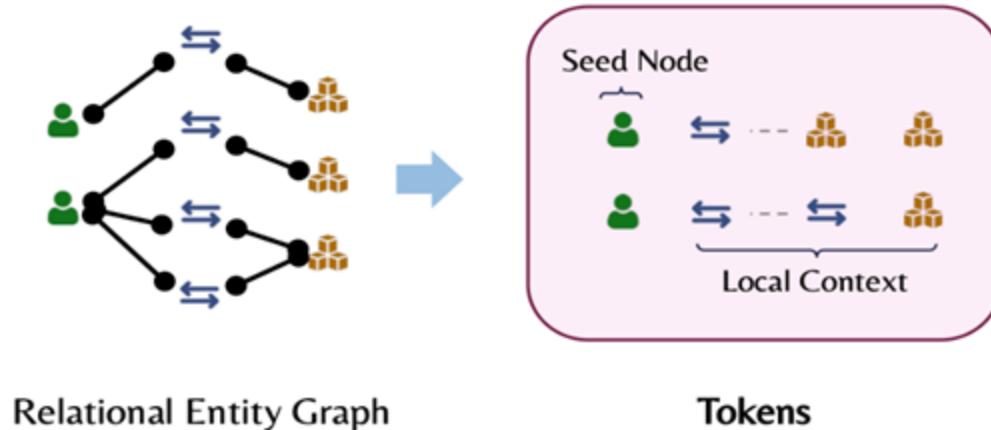
Unique properties of Relational Graphs

- The relational entity graph is:
 - Large-scale
 - Heterogeneous
 - Temporal
- Add a categorical encoding to capture the heterogeneity
- Add a time encoding to capture the temporal structure

How to encode structure

- **Challenge:** The relational entity graph is
 - Large-scale
- **Solution:** We need to sample a subgraph around the seed node

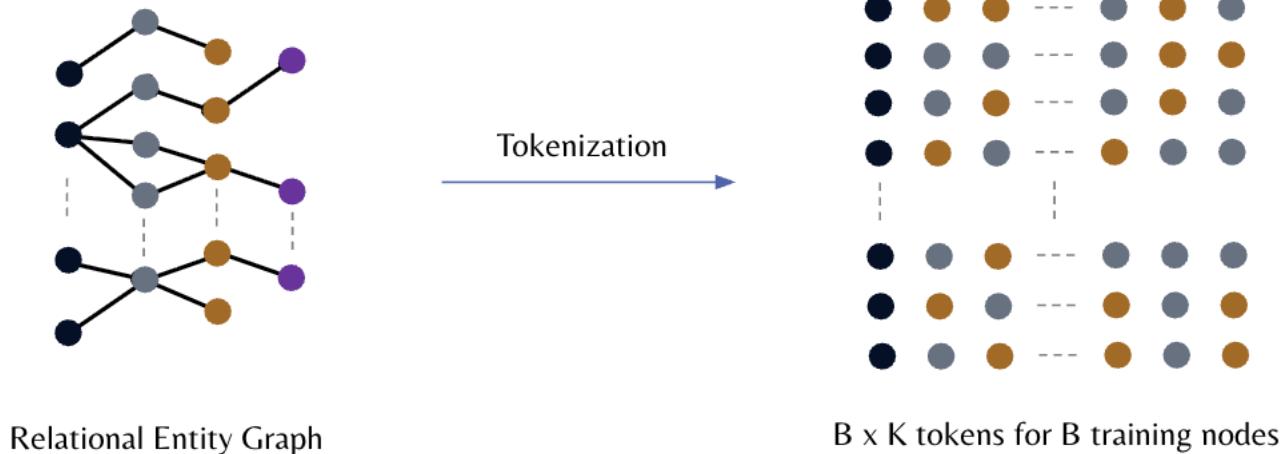
Relational Graph Transformer



Relational Entity Graphs

- **Training (seed) nodes** which correspond to node-types with respect to which tasks are defined.
- We are going to sample a context (subgraph) around each seed node.

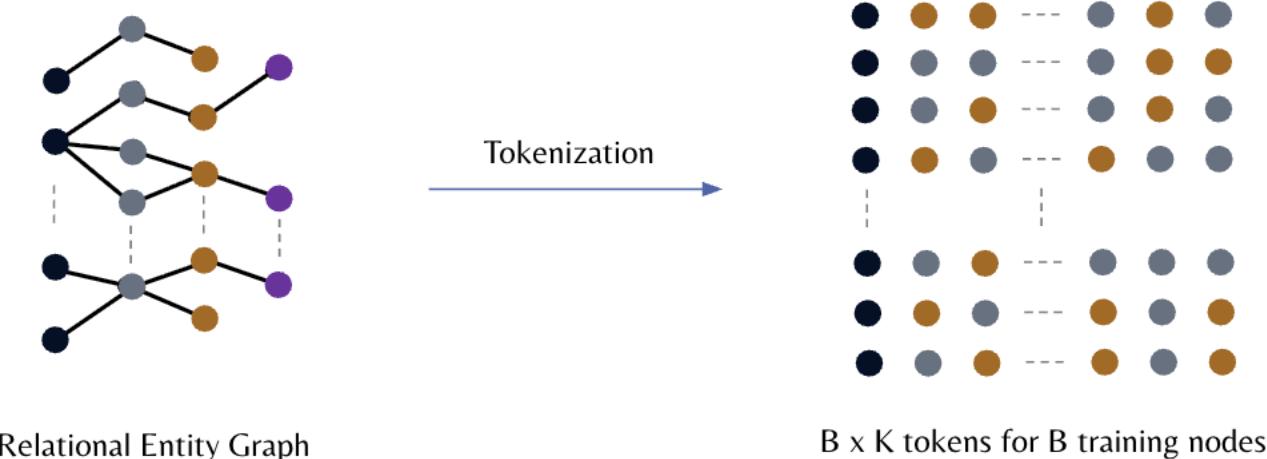
Relational Graph Transformer



Token Preparation (Sampling Stage)

- For each training node (●), a fixed set of K tokens (●●○) from local neighborhood (e.g., up to 2 hops) is selected through a temporal-aware sampling.
- We encode structure with a **GNN positional encoding** for each node in the context + **hop distance from seed node**.

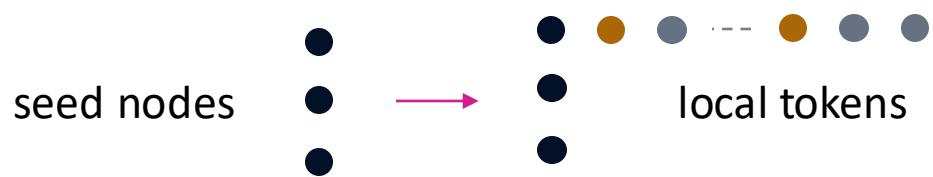
Relational Graph Transformer



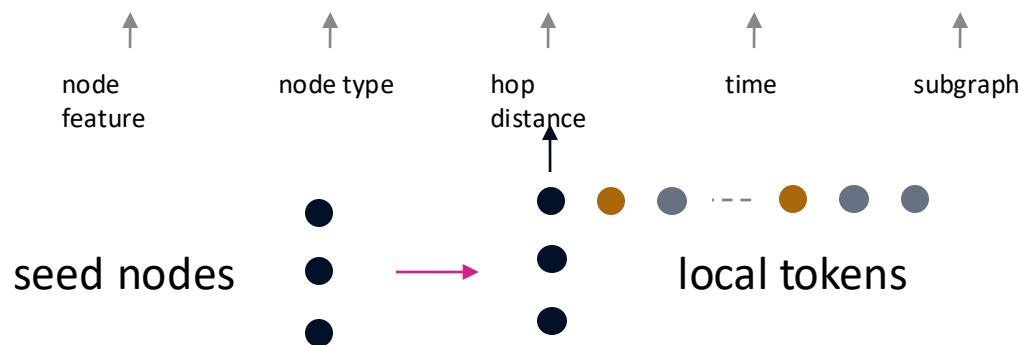
Token Preparation (Sampling Stage)

- For each training node (●), a fixed set of K tokens (●●○) from local neighborhood (e.g., up to 2 hops) is selected through a temporal-aware sampling.
- Each node in K set is represented by a **5-tuple**:
(node feature, node type, hop distance, time, GNN PE)

Relational Graph Transformer



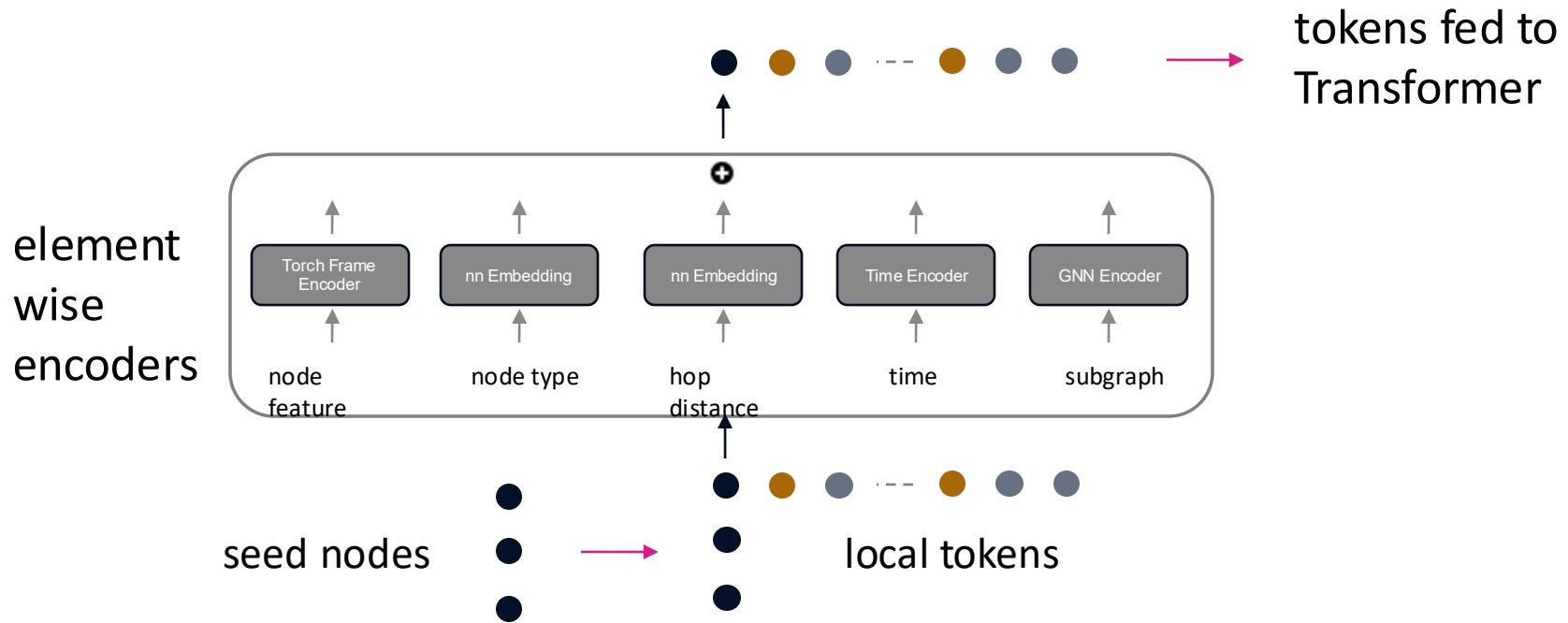
Relational Graph Transformer



Token Feature Encoders

- Each node in K set is represented by a **5-tuple**:
(node feature, node type, hop distance, time, GNN PE)

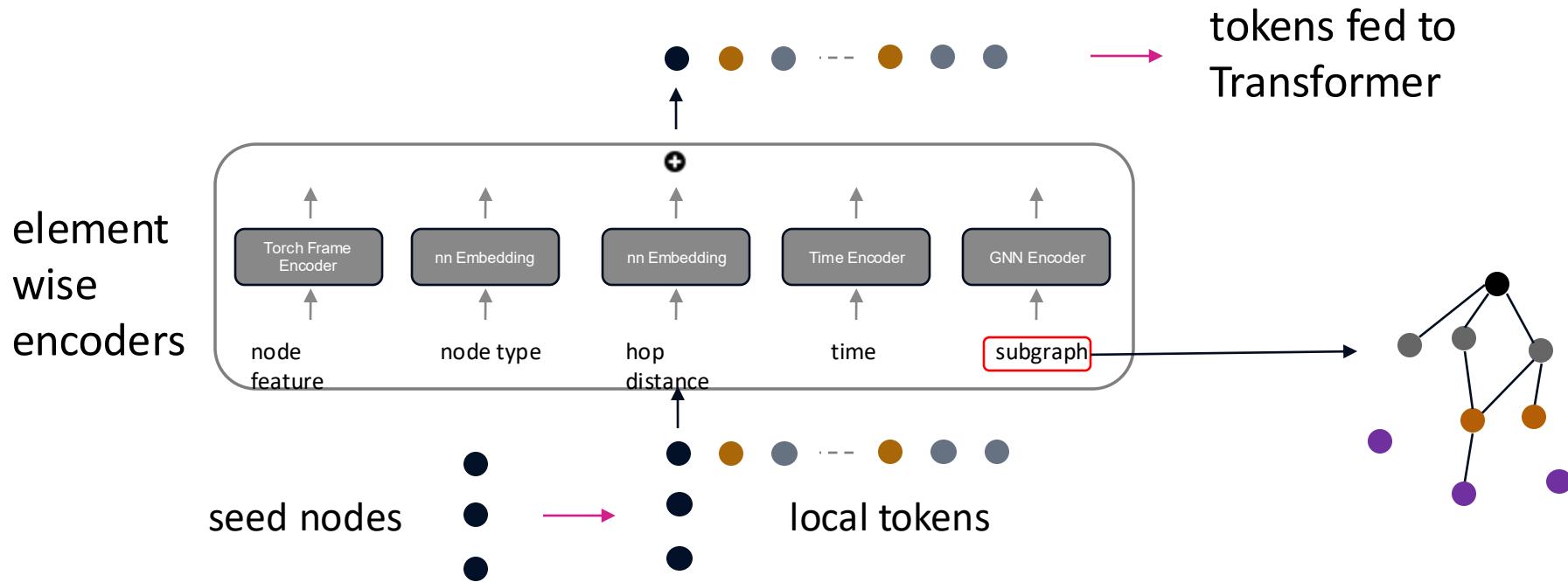
Relational Graph Transformer



Token Feature Encoders

- Each node in K set is represented by a **5-tuple**:
(node feature, node type, hop distance, time, GNN PE)
- Each of these elements is passed to a **feature encoder** and the combination becomes the token features for Transformer Network

Relational Graph Transformer



Token Feature Encoders

- Each node in K set is represented by a **5-tuple**:
(node feature, node type, hop distance, time, GNN PE)
- Each of these elements is passed to a **feature encoder** and the combination becomes the token features for Transformer Network

Ablation studies

Takeaway:

- Contribution of the multi-element tokenization strategy to capture arbitrary structure in relational data

Table 2: Relative drop (%) in performance in RELGT after removing a model component. Negative scores suggest the component is critical in RELGT, and vice-versa. Full results in Table 7.

Dataset	Task	No Global Module	No GNN PE	No Node Type	No Hop Distance	No Relative Time
rel-avito	ad-ctr	-6.00	-1.14	-7.14	-3.43	-9.14
rel-avito	user-clicks	7.79	-15.19	4.96	5.72	8.32
rel-avito	user-visits	-0.32	-2.35	-0.08	0.42	-0.72
rel-event	user-ignore	-1.28	0.15	-0.09	0.69	-0.06
rel-trial	study-outcome	-2.08	-1.66	3.80	-0.37	2.54
rel-trial	site-success	-19.08	-9.23	-2.94	-21.56	-0.77
rel-amazon	user-churn	-0.64	-0.78	0.16	0.06	-2.20
rel-hm	item-sales	-9.33	-17.35	-12.69	0.93	-77.24
Average		-3.87	-5.95	-1.75	-2.19	-9.91

Results

Takeaway:

- Relational Graph Transformers improve over GNNs

Table 1: Test set results on the entity regression and classification tasks in RelBench. Best values are in **bold**. RDL: HeteroGNN baseline [40], HGT: Heterogeneous GT [20], PE: Laplacian Positional Encodings [9]. Relative gains are expressed as percentage improvement over RDL baseline.

(a) MAE for entity regression. Lower is better

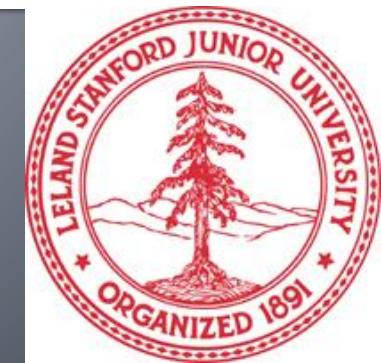
Dataset	Task	RDL	HGT	HGT +PE	RelGT (ours)	% Rel. Gain
rel-f1	driver-position	4.022	4.1598	4.2358	3.9170	2.61
	driver-top3				0.7554	10.56
rel-avito	ad-ctr	0.041	0.0441	0.0494	0.0345	15.85
rel-event	user-attendance	0.258	0.2635	0.2562	0.2502	2.79
rel-trial	study-adverse	44.473	43.3253	42.4622	43.9923	1.08
	site-success	0.400	0.4374	0.4431	0.3263	18.43
rel-amazon	user-ltv	14.313	15.3804	15.9296	14.2665	0.32
	item-ltv	50.053	56.1384	55.6211	48.9222	2.26
rel-stack	post-votes	0.065	0.0679	0.0680	0.0654	-0.62
rel-hm	item-sales	0.056	0.0655	0.0641	0.0536	4.29

(b) AUC for entity classification. Higher is better.

Dataset	Task	RDL	HGT	HGT +PE	RelGT (ours)	% Rel. Gain
rel-f1	driver-dnf	0.7262	0.7142	0.7109	0.7587	4.48
	driver-top3	0.7554	0.6389	0.8340	0.8352	10.56
rel-avito	user-clicks	0.6590	0.6584	0.6387	0.6830	3.64
	user-visits	0.6620	0.6426	0.6507	0.6678	0.88
rel-event	user-repeat	0.7689	0.6717	0.6590	0.7609	-1.04
	user-ignore	0.8162	0.8348	0.8161	0.8157	-0.06
rel-trial	study-outcome	0.6860	0.5679	0.5691	0.6861	0.01
rel-amazon	user-churn	0.7042	0.6608	0.6589	0.7039	-0.04
	item-churn	0.8281	0.7824	0.7840	0.8255	-0.31
rel-stack	user-engagement	0.9021	0.8898	0.8852	0.9053	0.35
	user-badge	0.8986	0.8652	0.8518	0.8632	-3.94
rel-hm	user-churn	0.6988	0.6773	0.6491	0.6927	-0.87

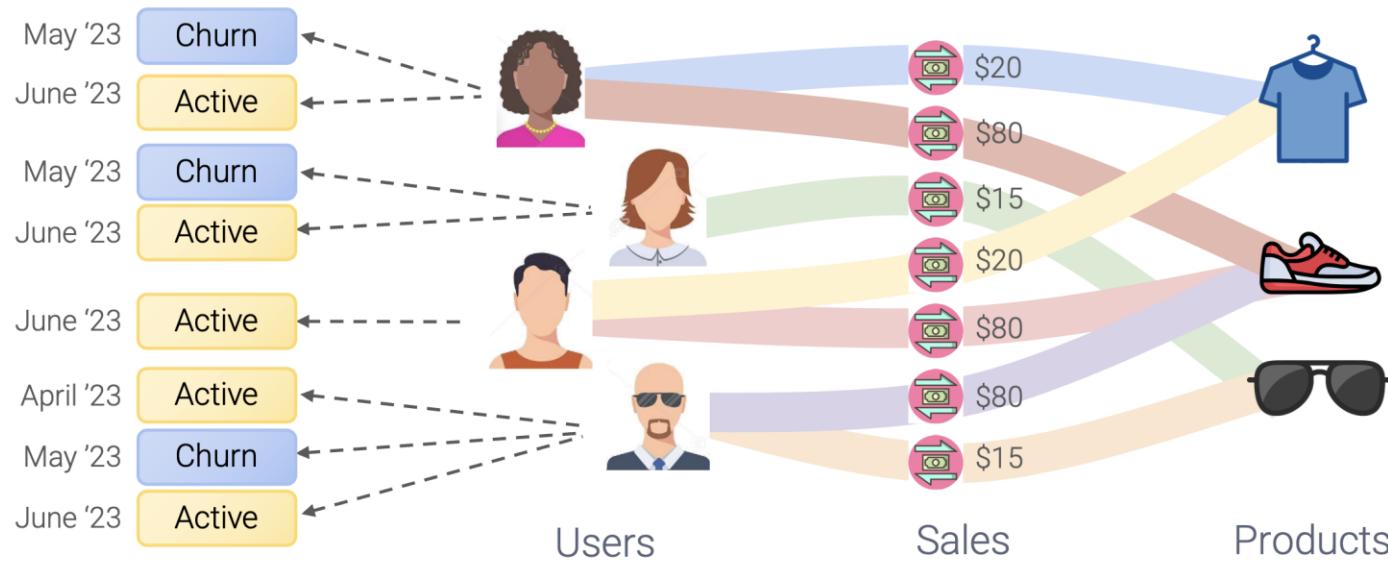
Stanford CS224W: Relational Foundation Model

CS224W: Machine Learning with Graphs
Jure Leskovec, Stanford University
<http://cs224w.stanford.edu>



Limitations of previous models

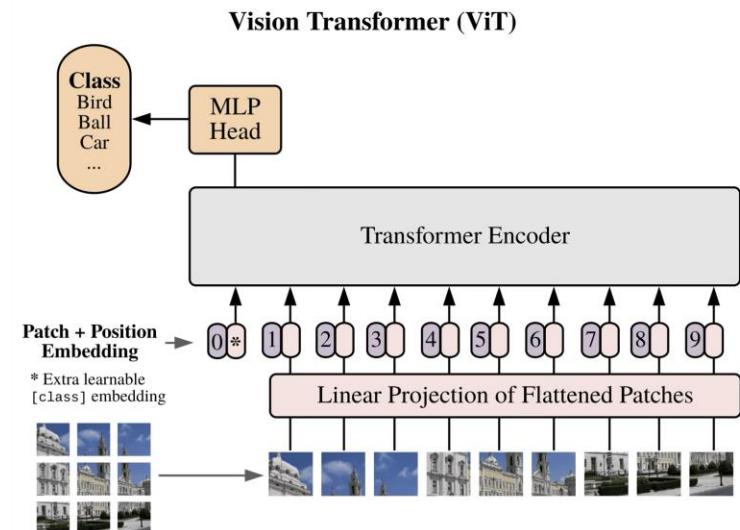
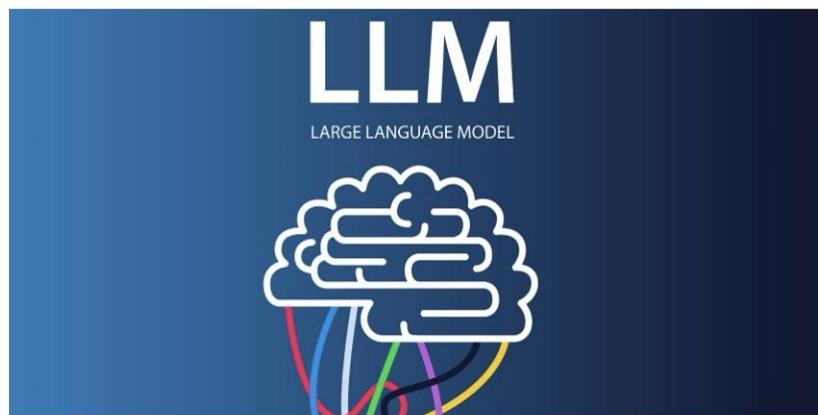
- They are schema- and task- specific



- They do not transfer knowledge from one database to another or one task to another
- Do not generalize to new databases or tasks

Can we do better?

- Goal: Bring the success of foundation models (LLMs, ViT) to relational data.



Relational Foundation Model

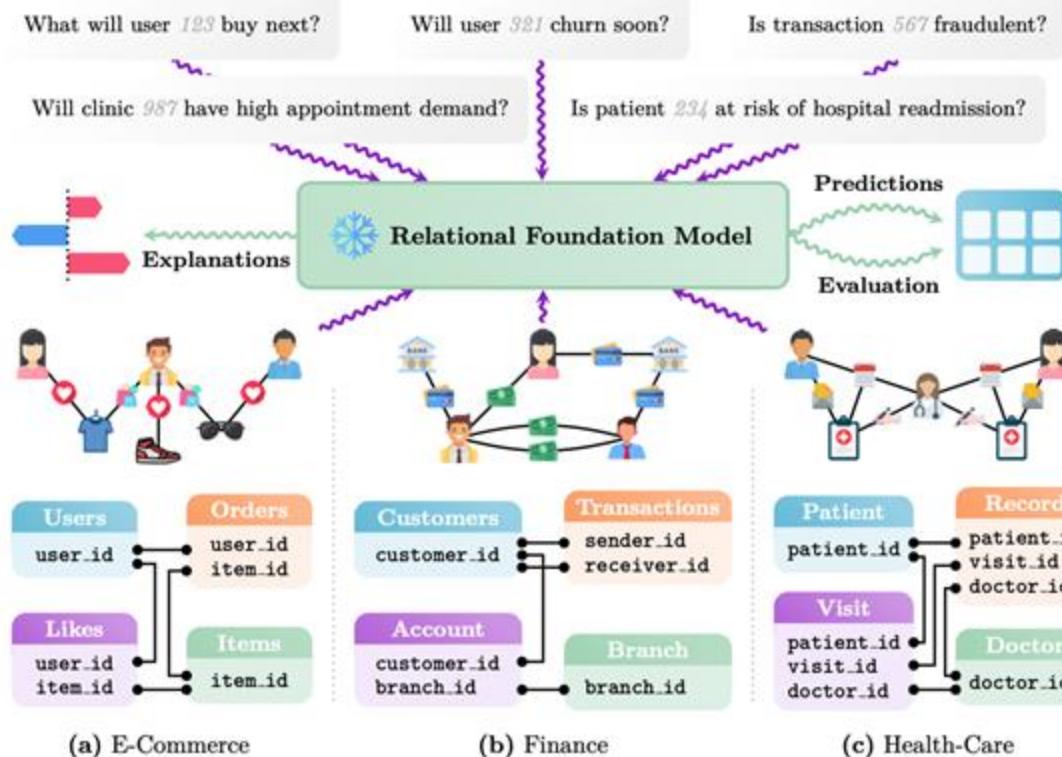


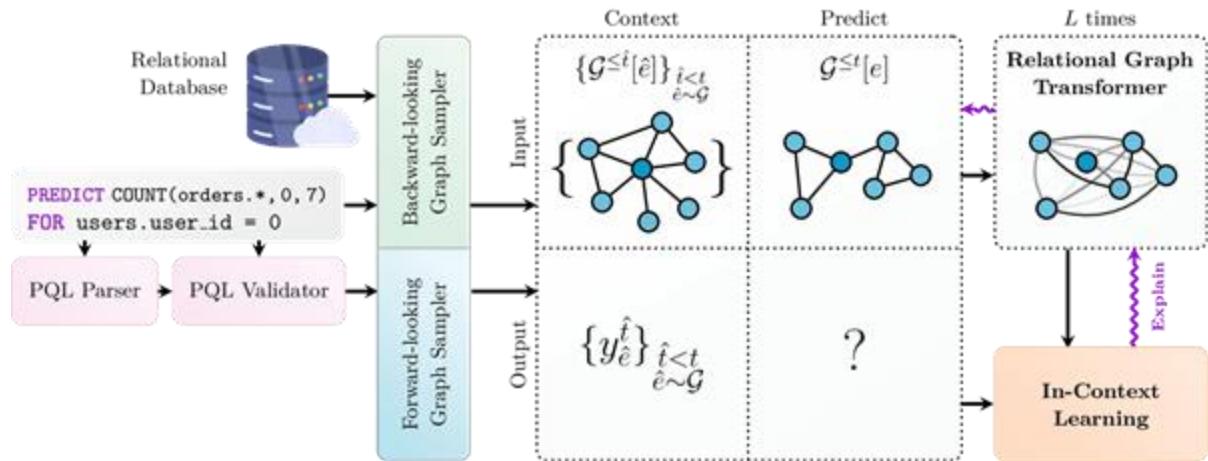
Figure 1: **Key capabilities of Relational Foundation Models.** RFMs can be applied to new/unseen databases and schemas with highly varying structural characteristics, as found in (a) e-commerce, (b) finance, or (c) health care. Secondly, they can be applied to any predictive task type, ranging from one-off assessments (*e.g.* entity-level fraud prediction) to temporal predictive queries (*e.g.* temporal recommendation prediction). Thirdly, they generalize to new predictive tasks and give accurate predictions without any task-specific model tuning. Finally, not only do RFMs support prediction outputs, but they also offer insights into the reasoning processes via explanations, and build trust through extensive quantitative evaluation mechanisms.

RelGT as the backbone for RFM

 News / Product / 

Introducing KumoRFM: A Foundation Model for In-Context Learning on Relational Data

May 20, 2025



Relational Graph Transformer is the backbone for
first Relational Foundation Model
with zero-shot capabilities powered by ICL

Relational Foundation Model

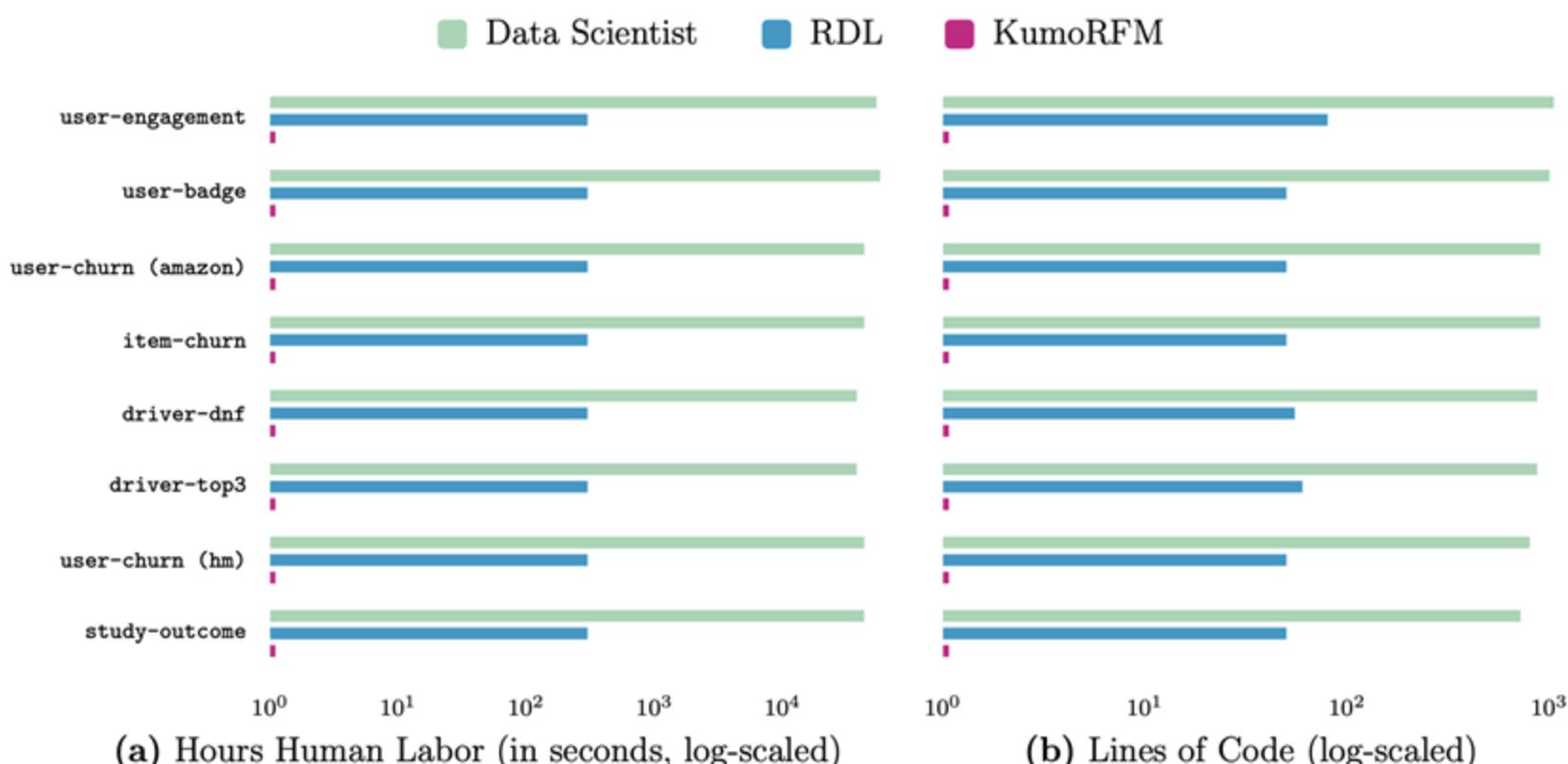


Figure 4: **Time to first prediction in (a) hours of human labor and (b) lines of code (LoC), illustrated on node-level classification tasks.** KumoRFM is orders of magnitude faster (≈ 1 second vs. ≈ 30 minutes vs. ≈ 12.3 hours) and requires zero-code to get to accurate predictions (1 LoC vs. ≈ 56 LoC vs. ≈ 878 LoC) compared to Data Scientist and RDL baselines, respectively.