

Heterogeneous Graph Transformer

Ziniu Hu, Yuxiao Dong, Kuansan Wang, Yizhou Sun

Jun 03, 2022

Outline

① Introduction

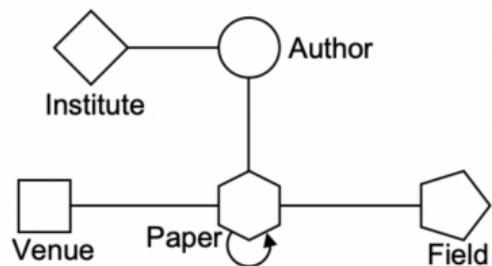
② Method

③ Experiments

④ Futures

Introduction

Heterogeneous Graph (HG) also known as heterogeneous information networks (HIN). A heterogeneous graph can represent as $\mathcal{G} = (\mathcal{V}, \xi)$, where each node \mathcal{V} , and each edge ξ has its own type Γ_v and Γ_e . A heterogeneous graph have two mapping function: $\phi_v : V \rightarrow \Gamma_v$ for node to node types, and $\phi_e : \xi \rightarrow \Gamma_e$ for edge types.



(a) The schema of
heterogeneous academic networks

Author	<i>is_(first/last/other)_author_of</i>	Paper
Author	<i>is_affiliated_with</i>	Institute
Paper	<i>is_published_(conf/journal)_at</i>	Venue
Paper	<i>has_(L1-L5)_field_of</i>	Field
Paper	<i>has_citation_to</i>	Paper

(b) The meta relations of heterogeneous
academic networks

Problem

- Meta-path need domain knowledge.
- Different types of nodes/edges share features.
- Different types of nodes/edges keep different non-shared weights
- Ignore the dynamic of heterogeneous graph
- Incapable of modeling Web-scale (large) heterogeneous graph

Heterogeneous graph transformer (HGT)

- Node and edge type dependent attention mechanism.
 - Not parameterizing each type of edges
 - use meta relation triplet $e = (s, t)$, where s is source node, t is target node
- Relative temporal encoding (RTE) strategy for dynamic graph
- HGSampling for Web-scale graph data.

Symbols

Graph $G = (\mathcal{V}, \mathcal{E}, \mathcal{A}, \mathcal{R})$

Node $v \in \mathcal{V}$, also s, t

Edge $e \in \mathcal{E}$

Node Type $\tau(v) : \mathcal{V} \rightarrow \mathcal{A}$

Edge Type $\phi(e) : \mathcal{E} \rightarrow \mathcal{R}$

edge, source node, target node $e = (s, t)$

meta relation triplet $< \tau(s), \phi(e), \tau(t) >$

Method

Use the **meta-relations** fo heterogeneous graph to parameterize weight matrices for heterogeneous mutual attention, message passing, and propagation steps.

Three steps:

- Heterogeneous Mutual Attention
 - input embedding of s_1, s_2, t
 - output attention matrix of $\phi(e)$.
- Heterogeneous Message Passing
 - output message of $\phi(e)$
- Target-Specific Aggregation

Method

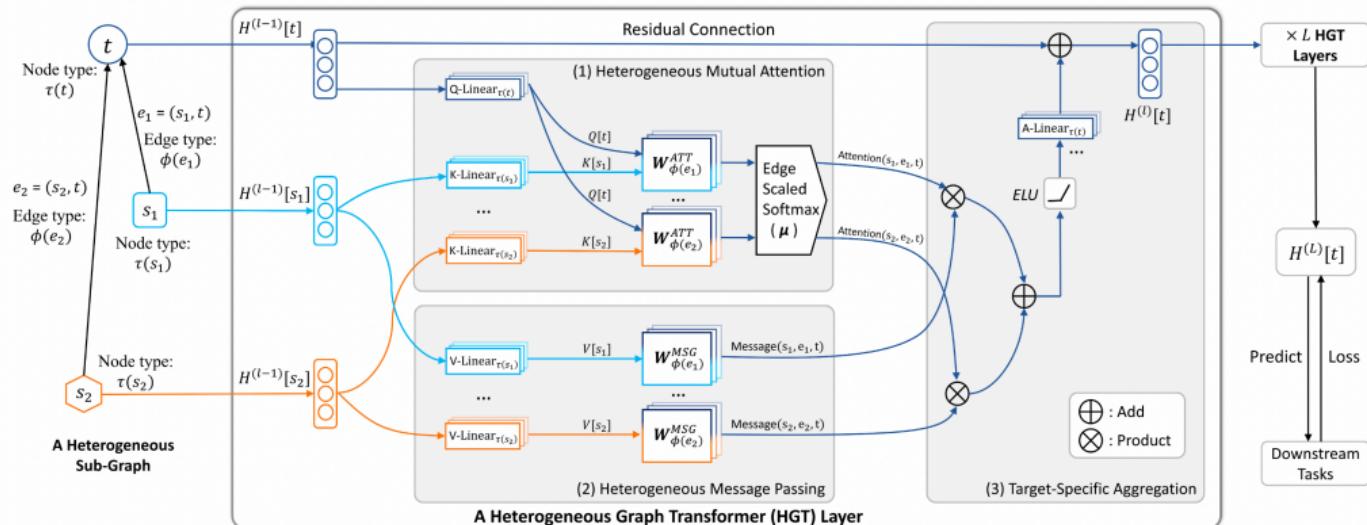


Figure 2: The Overall Architecture of Heterogeneous Graph Transformer. Given a sampled heterogeneous sub-graph with t as the target node, s_1 & s_2 as source nodes, the HGT model takes its edges $e_1 = (s_1, t)$ & $e_2 = (s_2, t)$ and their corresponding meta relations $<\tau(s_1), \phi(e_1), \tau(t)>$ & $<\tau(s_2), \phi(e_2), \tau(t)>$ as input to learn a contextualized representation $H^{(L)}$ for each node, which can be used for downstream tasks. Color decodes the node type. HGT includes three components: (1) meta relation-aware heterogeneous mutual attention, (2) heterogeneous message passing from source nodes, and (3) target-specific heterogeneous message aggregation.

Heterogeneous Mutual Attention

GAT:

$$H^l[t] \leftarrow \underset{\forall s \in N(t), \forall e \in E(s,t)}{\text{Aggregate}} \left(\text{Attention}(s, t) \cdot \text{Message}(s) \right) \quad (2)$$

$$\text{Attention}_{GAT}(s, t) = \text{Softmax} \left(\vec{a} \left(WH^{l-1}[t] \parallel WH^{l-1}[s] \right) \right)$$

$$\text{Message}_{GAT}(s) = WH^{l-1}[s]$$

$$\text{Aggregate}_{GAT}(\cdot) = \sigma \left(\text{Mean}(\cdot) \right)$$

Attention Importance of each source node.

Message Extracts the message by using only the source node.

Aggregate Aggregate the neighborhood message by the attention weight.

Heterogeneous Mutual Attention

Transformer: W_q, W_k, W_v

HGT:

$$\text{Attention}_{HGT}(s, e, t) = \text{Softmax} \left(\parallel_{\substack{\forall s \in N(t) \\ i \in [1, h]}} \text{ATT-head}^i(s, e, t) \right) \quad (3)$$

$$\text{ATT-head}^i(s, e, t) = \left(K^i(s) \ W_{\phi(e)}^{ATT} \ Q^i(t)^T \right) \cdot \frac{\mu_{\langle \tau(s), \phi(e), \tau(t) \rangle}}{\sqrt{d}}$$

$$K^i(s) = \text{K-Linear}_{\tau(s)}^i \left(H^{(l-1)}[s] \right)$$

$$Q^i(t) = \text{Q-Linear}_{\tau(t)}^i \left(H^{(l-1)}[t] \right)$$

- $W_{\phi(e)}^{ATT}$
- $\mu_{\langle \tau(s), \phi(e), \tau(t) \rangle}$

Message passing

$$\textbf{Message}_{HGT}(s, e, t) = \parallel_{i \in [1, h]} MSG\text{-}head^i(s, e, t)$$

$$MSG\text{-}head^i(s, e, t) = \text{M-Linear}_{\tau(s)}^i \left(H^{(l-1)}[s] \right) W_{\phi(e)}^{MSG}$$

- Edge dependent: $W_{\tau(e)}^{MSG}$
- Incorporate the meta relations of edges into the message passing process to alleviate the distribution differences of nodes and edges of different types.

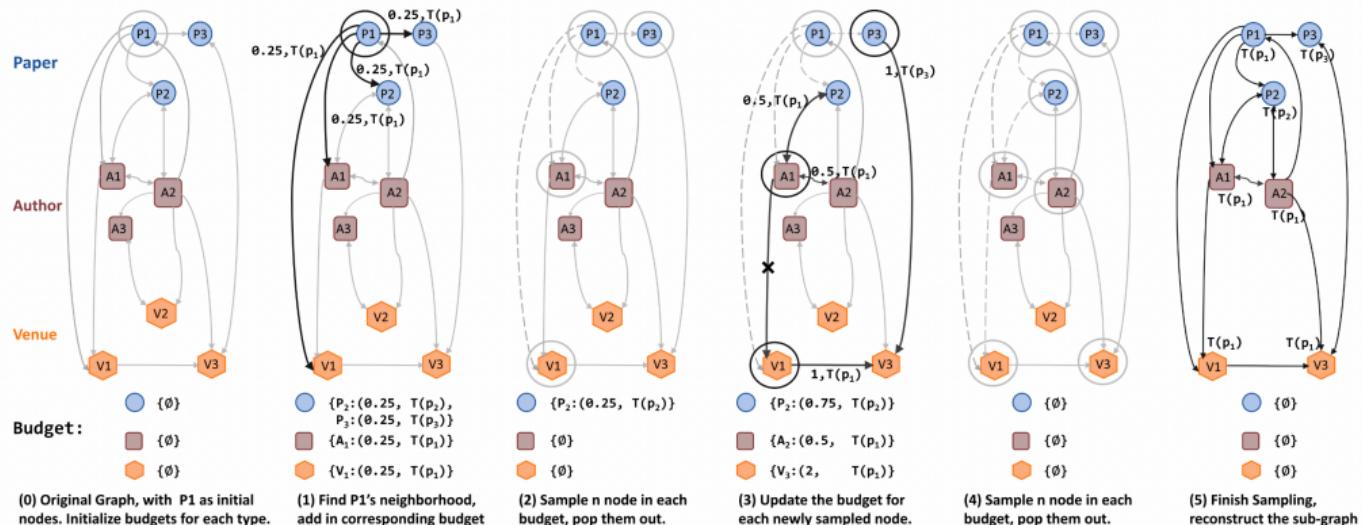
Target-Specific Aggregation

$$\tilde{H}^{(l)}[t] = \bigoplus_{\forall s \in N(t)} \left(\mathbf{Attention}_{HGT}(s, e, t) \cdot \mathbf{Message}_{HGT}(s, e, t) \right).$$

$$H^{(l)}[t] = \text{A-Linear}_{\tau(t)} \left(\sigma(\tilde{H}^{(l)}[t]) \right) + H^{(l-1)}[t]. \quad (5)$$

- A-Linear $_{\tau(t)}$ to map target node t to type specific distribution and update the l -th HGT layers embedding.

HGSampling



- keep a similar number of nodes and edges for each type, and keep the sampled sub-graph dense to minimize the information loss and reduce the sample variance.

Relative Temporal Encoding (RTE)

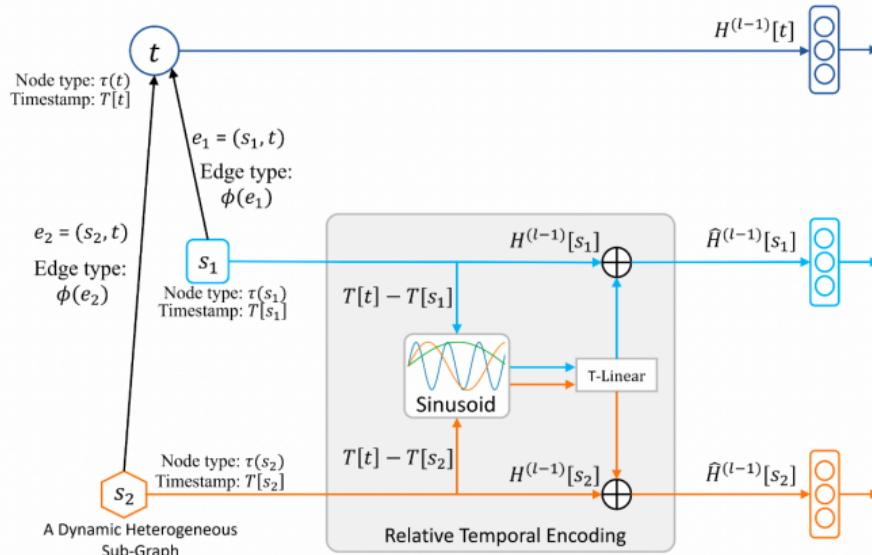


Figure 3: Relative Temporal Encoding (RTE) to model graph dynamic. Nodes are associated with timestamps $T(\cdot)$. After the RTE process, the temporal augmented representations are fed to the HGT model.

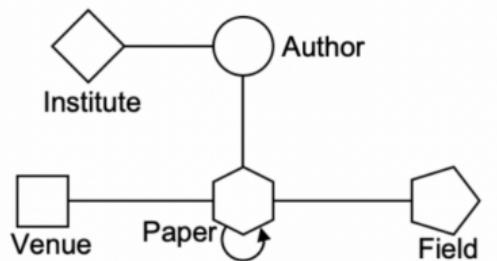
Relative Temporal Encoding (RTE)

$$\widehat{H}^{(l-1)}[s] = H^{(l-1)}[s] + RTE(\Delta T(t, s)) \quad (9)$$

- $\Delta T(s, t) = T(s) - T(t)$



OAG Data



(a) The schema of
heterogeneous academic networks

Author	<i>is_(first/last/other)_author_of</i>	Paper
Author	<i>is_affiliated_with</i>	Institute
Paper	<i>is_published_(conf/journal)_at</i>	Venue
Paper	<i>has_(L_1-L_5)_field_of</i>	Field
Paper	<i>has_citation_to</i>	Paper

(b) The meta relations of heterogeneous
academic networks

Figure 1: The schema and meta relations of Open Academic Graph (OAG). Given a Web-scale heterogeneous graph, e.g., an academic network, HGT takes only its one-hop edges as input without manually designing meta paths.

OAG Data

- OAG
 - All
 - Computer Science (CS)
 - Medicine (Med)

Baseline models

- Graph Convolutional Networks (GCN)
- Graph Attention Networks (GAT)
- Relational Graph Convolutional Networks
 - Keep a different weight for each relationship (edge).
 - $$h_i^{(l+1)} = \sigma \left(\sum_{r \in \mathcal{R}} \sum_{j \in \mathcal{N}_i^r} \frac{1}{c_{i,r}} W_r^{(l)} h_j^{(l)} + W_0^{(l)} h_i^{(l)} \right)$$
- Heterogeneous Graph Neural Networks
 - Adopt different BiLSTM for node type and neighbor information
- Heterogeneous Graph Attention Networks (HAN)
 - Hierarchical attentions to aggregate neighbor via meta-paths

Results

GNN Models		GCN [9]	RGCN [14]	GAT [22]	HetGNN [27]	HAN [23]	HGT ^{-RTE} -Heter	HGT ^{+RTE} -Heter	HGT ^{-RTE} +Heter	HGT ^{+RTE} +Heter	
# of Parameters		1.69M	8.80M	1.69M	8.41M	9.45M	3.12M	3.88M	7.44M	8.20M	
Batch Time		0.46s	1.24s	0.97s	1.35s	2.27s	1.11s	1.14s	1.48s	1.50s	
CS	Paper-Field (L_1)	NDCG MRR	.608±.062 .679±.069	.603±.065 .683±.056	.622±.071 .694±.065	.612±.063 .689±.060	.618±.058 .691±.051	.662±.051 .751±.036	.689±.042 .779±.027	.705±.036 .799±.023	.718±.014 .823±.019
	Paper-Field (L_2)	NDCG MRR	.344±.021 .353±.053	.322±.053 .340±.061	.357±.058 .382±.057	.346±.071 .373±.051	.352±.051 .388±.065	.362±.048 .394±.072	.371±.043 .397±.064	.379±.047 .414±.076	.403±.041 .439±.078
	Paper-Venue	NDCG MRR	.406±.081 .215±.066	.412±.076 .216±.105	.437±.082 .239±.089	.431±.074 .245±.069	.449±.072 .254±.074	.456±.069 .258±.085	.461±.066 .265±.090	.468±.074 .275±.089	.473±.054 .288±.088
	Author	NDCG Disambiguation	.826±.039 .661±.045	.835±.042 .665±.054	.864±.051 .694±.052	.850±.056 .668±.061	.859±.053 .688±.049	.867±.048 .703±.036	.875±.046 .712±.032	.886±.048 .727±.038	.894±.034 .732±.038
	Paper-Field (L_1)	NDCG MRR	.560±.056 .465±.055	.571±.061 .470±.082	.584±.076 .493±.069	.598±.068 .509±.054	.607±.054 .575±.057	.654±.048 .620±.066	.667±.045 .642±.062	.683±.037 .659±.055	.709±.029 .688±.048
	Paper-Field (L_2)	NDCG MRR	.334±.035 .337±.061	.337±.051 .343±.063	.344±.063 .370±.058	.342±.048 .373±.061	.350±.059 .379±.052	.359±.053 .385±.071	.365±.047 .397±.069	.374±.050 .408±.071	.384±.046 .417±.074
	Paper-Venue	NDCG MRR	.377±.059 .211±.045	.383±.062 .217±.058	.388±.065 .244±.091	.412±.057 .259±.072	.416±.068 .271±.056	.421±.083 .277±.081	.432±.078 .282±.085	.446±.083 .288±.074	.445±.085 .291±.062
	Author	MRR NDCG	.776±.042 .614±.051	.779±.048 .625±.049	.828±.044 .663±.046	.824±.058 .659±.061	.834±.056 .667±.053	.838±.047 .683±.055	.844±.041 .691±.046	.864±.043 .708±.041	.871±.040 .718±.043
OAG	Paper-Field (L_1)	NDCG MRR	.508±.141 .556±.136	.511±.128 .565±.105	.534±.103 .610±.096	.543±.084 .616±.076	.544±.096 .622±.092	.571±.089 .649±.081	.578±.086 .657±.078	.595±.089 .675±.082	.615±.084 .702±.081
	Paper-Field (L_2)	NDCG MRR	.318±.074 .322±.067	.328±.046 .332±.052	.339±.049 .348±.045	.336±.062 .350±.053	.342±.051 .358±.049	.350±.045 .362±.057	.354±.046 .369±.058	.358±.052 .371±.064	.367±.048 .378±.071
	Paper-Venue	NDCG MRR	.302±.066 .194±.070	.313±.051 .193±.047	.317±.057 .196±.052	.309±.071 .192±.059	.327±.062 .214±.067	.334±.058 .229±.061	.341±.059 .233±.060	.353±.064 .243±.048	.355±.062 .247±.061
	Author	NDCG MRR	.738±.042 .612±.064	.755±.048 .619±.057	.797±.044 .645±.063	.803±.058 .649±.052	.821±.056 .660±.049	.835±.043 .668±.059	.841±.041 .674±.058	.847±.043 .683±.066	.852±.048 .688±.054
	Disambiguation										

Table 2: Experimental results of different methods over the three datasets.

Futures

- Generate heterogeneous graphs
 - predict new papers and title
- Pre-train HGT to benefit tasks with scarce labels
-
- Downstream Tasks