

1045 A Appendix

1046 A.1 Benchmark Details

1048 In this section, we include the details on the datasets and the tasks in
 1049 RelBench [42] which we use for our evaluation. RelBench consists
 1050 of 7 datasets from diverse relational database domains, including
 1051 e-commerce, clinical records, social networks, and sports, among
 1052 others. These datasets are curated from their respective source do-
 1053 mains and consist a wide range of sizes, from 1.3K to 5.4M records in
 1054 the training set for the prediction tasks, with a total of 47M training
 1055 records. For each dataset, multiple predictive tasks are defined, such
 1056 as predicting a user’s engagement with an advertisement within the
 1057 next four days or determining whether a clinical trial will achieve
 1058 its primary outcome within the next year. In total, RelBench has
 1059 30 tasks across the 7 datasets, covering entity classification, entity
 1060 regression, and recommendation. For our evaluation, we focus on
 1061 21 tasks on entity classification and regression as RELGT primarily
 1062 serves as a node representation learning model in RDL. We exclude
 1063 recommendation tasks in this work since they involve specific con-
 1064 siderations, such as identifying target nodes [54] or using pair-wise
 1065 learning architectures [55] and using RELGT trivially in RDL is
 1066 sub-optimal. We detail the dataset and task statistics in Table 3.

1067 A.1.1 Datasets.

1069 **rel-amazon.** The Amazon E-commerce dataset consists of prod-
 1070 uct details, user information, and review interactions from Amaz-
 1071 on’s platform, including metadata like pricing and categories,
 1072 along with review ratings and content.

1073 **rel-avito.** Avito’s marketplace dataset contains search queries,
 1074 advertisement characteristics, and contextual information from this
 1075 major online trading platform that facilitates transactions across
 1076 various categories including real estate and vehicles.

1078 **rel-event.** The Event Recommendation dataset from Hangtime
 1079 mobile app tracks users’ social planning, capturing interactions,
 1080 event details, demographic data, and social connections to reveal
 1081 how relationships impact user behavior.

1083 **rel-f1.** The F1 dataset provides comprehensive Formula 1 rac-
 1084 ing information since 1950, documenting drivers, constructors, man-
 1085 ufacturers, and circuits with detailed records of race results, stand-
 1086 ings, and specific data on various racing sessions and pit stops.

1087 **rel-hm.** H&M’s dataset contains customer-product interactions
 1088 from their e-commerce platform, featuring customer demographics,
 1089 product descriptions, and purchase histories.

1091 **rel-stack.** The Stack Exchange dataset documents activity from
 1092 this network of Q&A websites, including user biographies, posts,
 1093 comments, edits, votes, and question relationships where users earn
 1094 reputation through contributions.

1095 **rel-trial.** The clinical trial dataset from the AACT initiative
 1096 has study protocols and outcomes, containing trial designs, partici-
 1097 pient information, intervention details, and results metrics, serving
 1098 as a key resource for medical research.

1100 **A.1.2 Tasks.** The following entity classification and regression
 1101 tasks are defined in RelBench for the above datasets.

1102 (1) **rel-amazon**

- (a) **user-churn:** Predict whether a user will discontinue reviewing products within the next three months.
- (b) **item-churn:** Predict if a product will have no reviews in the next three months.
- (c) **user-ltv:** Estimate the total monetary value of merchandise in dollar that a user will purchase and review within the next three months.
- (d) **item-ltv:** Estimate the total monetary value of purchases and reviews a product will receive during the next three months.

1111 (2) **rel-avito**

- (a) **user-visits:** Predict if a user will engage with several (advertisements) ads within the upcoming four days.
- (b) **user-clicks:** Predict whether a user will interact with multiple ads through clicking within the upcoming four days.
- (c) **ad-ctr:** Estimate the interaction probability for an ad, assuming it receives an interaction within four days.

1119 (3) **rel-event**

- (a) **user-attendance:** Estimate the number of events a user will confirm attendance to (RSVP yes or maybe) within the upcoming seven days.
- (b) **user-repeat:** Predict whether a user will join an event (RSVP yes or maybe) within the upcoming seven days, provided they attended in an event during the previous fourteen days.
- (c) **user-ignore:** Predict whether a user will disregard or ignore more than two events invitations within the upcoming seven days.

1132 (4) **rel-f1**

- (a) **driver-dnf:** Predict if a driver will not finish a race within the upcoming month.
- (b) **driver-top3:** Determine if a driver will achieve a top-three qualifying position in a race within the upcoming month.
- (c) **driver-position:** Estimate a driver’s average finishing placement across all races in the upcoming two months.

1141 (5) **rel-hm**

- (a) **user-churn:** Predict whether a customer will not perform any transactions in the upcoming week.
- (b) **item-sales:** Estimate total revenue generated by a product in the upcoming week.

1147 (6) **rel-stack**

- (a) **user-engagement:** Predict whether a user will contribute through voting, posting, or commenting within the upcoming three months.
- (b) **user-badge:** Predict whether a user will secure a new badge within the upcoming three months.
- (c) **post-votes:** Estimate the number of votes a user’s post will accumulate over the upcoming three months.

1155 (7) **rel-trial**

- (a) **study-outcome:** Predict whether a clinical trial will achieve its principal outcome within the upcoming year.

1156 1157 1158 1159 1160

Table 3: Dataset and task statistics from RelBench used for our evaluation.

Dataset	Task	Task type	#Rows of training table			#Unique Entities	%train/test Entity Overlap
			Train	Validation	Test		
rel-amazon	user-churn	classification	4,732,555	409,792	351,885	1,585,983	88.0
	item-churn	classification	2,559,264	177,689	166,842	416,352	93.1
	user-ltv	regression	4,732,555	409,792	351,885	1,585,983	88.0
	item-ltv	regression	2,707,679	166,978	178,334	427,537	93.5
rel-avito	user-clicks	classification	59,454	21,183	47,996	66,449	45.3
	user-visits	classification	86,619	29,979	36,129	63,405	64.6
	ad-ctr	regression	5,100	1,766	1,816	4,997	59.8
rel-event	user-repeat	classification	3,842	268	246	1,514	11.5
	user-ignore	classification	19,239	4,185	4,010	9,799	21.1
	user-attendance	regression	19,261	2,014	2,006	9,694	14.6
rel-f1	driver-dnf	classification	11,411	566	702	821	50.0
	driver-top3	classification	1,353	588	726	134	50.0
	driver-position	regression	7,453	499	760	826	44.6
rel-hm	user-churn	classification	3,871,410	76,556	74,575	1,002,984	89.7
	item-sales	regression	5,488,184	105,542	105,542	105,542	100.0
rel-stack	user-engagement	classification	1,360,850	85,838	88,137	88,137	97.4
	user-badge	classification	3,386,276	247,398	255,360	255,360	96.9
	post-votes	regression	2,453,921	156,216	160,903	160,903	97.1
rel-trial	study-outcome	classification	11,994	960	825	13,779	0.0
	study-adverse	regression	43,335	3,596	3,098	50,029	0.0
	site-success	regression	151,407	19,740	22,617	129,542	42.0

- (b) **study-adverse**: Estimate the number of patients who will experience significant adverse effects or mortality in a clinical trial over the upcoming year.
(c) **site-success**: Estimate the success rate of a clinical trial site in the upcoming year.

A.2 Node initialization for Subgraph GNN PE in RELGT

As described in Section 3.1, we employ a lightweight GNN PE to capture local graph structures that cannot be represented by other elements of the token, particularly the parent-child relationships among nodes in the local subgraph. The GNN is implemented as:

$$h_{\text{pe}}(v_j) = \text{GNN}(A_{\text{local}}, Z_{\text{random}})_j \in \mathbb{R}^d \quad (9)$$

where $\text{GNN}(\cdot, \cdot)_j$ is a lightweight GNN applied to the local subgraph, yielding the encoding for node v_j . Here, $A_{\text{local}} \in \mathbb{R}^{K \times K}$ represents the adjacency matrix of the sampled subgraph containing K nodes, and $Z_{\text{random}} \in \mathbb{R}^{K \times d_{\text{init}}}$ denotes randomly initialized node features for the GNN (with d_{init} as the initial feature dimension). In RELGT, we set $d_{\text{init}} = 1$.

The randomly initialized node features (Z_{random}) provide enhanced properties as discussed in Section 3.1. We investigate the alternative approach of using Laplacian PE (Z_{LapPE}) computed over the subgraph instead of random initialization and report these results in Table 4. For these results, we utilized a positional encoding dimension size of 4. Our findings indicate that Z_{LapPE} consistently

underperforms compared to Z_{random} , while also introducing additional computational overhead ranging from $1.02\times$ to $3.38\times$ across the 8 selected tasks in our study. This shows the challenges of using existing PEs such as Laplacian PE in relational entity graphs and signify the use of GNN PE as part of RELGT’s tokenization strategy.

A.3 HGT Baseline

In the main experiments (Section 4), we use the Heterogeneous Graph Transformer (HGT) [20] as a graph transformer (GT) baseline, and report results for two variants to demonstrate the advantages of RELGT over existing GT models. Specifically, we consider the standard HGT model and an enhanced version, HGT+PE, which incorporates Laplacian positional encodings (LapPE). These positional encodings are computed on sampled subgraphs rather than the full graph.

For implementation, we use the HGTConv layer from PyTorch Geometric [14] and integrate it into the RDL pipeline [42] by replacing the default GNN module. Both variants use 4 attention heads and 2 layers, similar to the configuration of the GNN module in RDL, with residual connections and layer normalization applied between layers. For the HGT+PE variant, we use LapPE of dimension 4 for all tasks, except for rel-amazon item-ltv and rel-hm item-sales, where we use dimension 2. Notably, because the relational entity graphs are heterogeneous, the Laplacian positional encodings is computed multiple times for each node type, unlike the original homogeneous setting for which LapPE was designed [10].

Table 4: Study of node initialization in Subgraph GNN PE. Relative drop is expressed as percentage drop of using Z_{LapPE} vs. Z_{random} and runtime ratio compares the time for Z_{LapPE} vs. Z_{random} .

Dataset	Task (# train)	Performance			Epoch time (m)			Runtime Ratio
		MAE ↓	Z_{random}	Z_{LapPE}	% Rel Drop	Z_{random}	Z_{LapPE}	
rel-avito	ad-ctr	Test	0.035	0.0369	-5.43	0.76	2.57	3.38
		Val	0.0314	0.0314				
rel-trial	site-success	Test	0.326	0.3452	-5.89	32.88	36.09	1.1
		Val	0.359	0.3683				
rel-hm	item-sales	Test	0.0536	0.0573	-6.9	49.26	53.8	1.09
		Val	0.0627	0.0667				
Dataset	Task (# train)	AUC ↑	Z_{random}	Z_{LapPE}	% Rel Drop	Z_{random}	Z_{LapPE}	Runtime Ratio
		Test	0.607	0.583	-3.95	6.42	7.43	1.16
rel-avito	user-clicks	Val	0.656	0.6564				
		Test	0.664	0.6626	-0.21	9.26	10.50	1.13
rel-event	user-ignore	Val	0.699	0.7002				
		Test	0.8	0.7988	-0.15	1.85	2.77	1.5
rel-trial	study-outcome	Val	0.881	0.8916				
		Test	0.674	0.6532	-3.09	1.41	1.52	1.08
rel-amazon	user-churn	Val	0.689	0.6719				
		Test	0.7039	0.7044	0.07	168.00	170.55	1.02
Val			0.7036	0.7036				

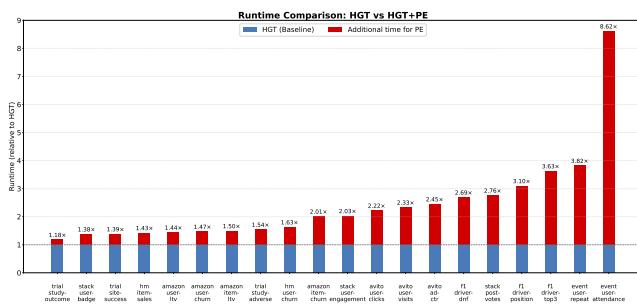


Figure 5: Runtime Comparison of HGT and HGT+PE baseline. Adding the Laplacian Positional Encoding increases computational overhead, with penalties on average training time per epoch. The overhead for PE reaches up to 761% relative to the training time of HGT on the same dataset.

In addition to the main results in Table 1, we report per-epoch runtimes in Figure 5 and Table 5. We observe a significant computational overhead from precomputing Laplacian positional encodings, with slowdowns ranging from $1.8\times$ to $8.62\times$, highlighting the challenge of directly applying existing graph PE techniques *as is* to relational entity graphs, and signifying the contributions of RELGT.

A.4 Detailed Results

In Table 6, we report the full results of different configurations we tuned for RELGT, particularly on the smaller datasets with lesser than a million training nodes. Table 7 provides the full scores for the RELGT component study in Table 2, while Table 8 provides the supporting results for Figure 4. Finally, we provide the elaborated version of the Tables 1a and 1b in Tables 9 and 9, respectively.

Table 5: Relative performance drop (%) when position encoding (PE) is removed from HGT+PE models and average training time per epoch of HGT and HGT+PE. Negative scores suggest the PE is critical, and vice-versa. HGT+PE consistently requires more training time per epoch compared to HGT without PE across all datasets.

Dataset	Task	No PE	HGT(s)	HGT+PE(s)
rel-f1	driver-position	1.79	1.47	4.56
rel-avito	ad-ctr	10.73	1.63	4.00
rel-event	user-attendance	-2.85	4.36	37.57
rel-trial	study-adverse	-2.03	9.72	15.02
rel-trial	site-success	1.29	45.73	63.41
rel-amazon	user-ltv	3.45	73.59	106.21
rel-amazon	item-ltv	-0.93	73.68	110.33
rel-stack	post-votes	0.15	191.23	528.25
rel-hm	item-sales	-2.18	94.66	135.05
rel-f1	driver-dnf	0.46	2.54	6.84
rel-f1	driver-top3	-23.39	0.38	1.38
rel-avito	user-clicks	3.08	11.09	24.66
rel-avito	user-visits	-1.24	17.16	40.07
rel-event	user-repeat	1.93	1.35	5.16
rel-event	user-ignore	2.29	4.49	651.10
rel-trial	study-outcome	-0.21	4.09	4.83
rel-amazon	user-churn	0.29	78.56	115.53
rel-amazon	item-churn	-0.20	75.51	152.06
rel-stack	user-engagement	0.52	175.16	356.07
rel-stack	user-badge	1.57	153.68	212.21
rel-hm	user-churn	4.34	77.73	127.04
Average		-0.05	52.28	128.64

A.5 Resource Information.

We implement RELGT using PyTorch framework [39], PyTorch Geometric framework [14] and adapt the codebase of relational deep learning [42] <https://github.com/snap-stanford/relbench>. All our experiments are conducted on an NVIDIA A100 GPU server with 8 GPU nodes.

Table 8: Ablation of context size K in RELGT.

Dataset	Task (# train)	MAE ↓	RELGT K=100	RELGT K=300	RELGT K=500
rel-avito	ad-ctr		Test	0.0375	0.0374
			Val	0.0329	0.0319
rel-trial	site-success		Test	0.3739	0.3674
			Val	0.3708	0.372
rel-hm	item-sales		Test	0.055	0.0532
			Val	0.0643	0.0619
				0.052	
Dataset	Task (# train)	AUC ↑	RelGT K=100	RelGT K=300	RelGT K=500
rel-avito	user-clicks		Test	0.6628	0.6491
			Val	0.6437	0.6622
rel-avito	user-visits		Test	0.6664	0.6653
			Val	0.7013	0.701
rel-event	user-ignore		Test	0.7674	0.8105
			Val	0.8682	0.8853
rel-trial	study-outcome		Test	0.7078	0.6526
			Val	0.6575	0.663
rel-amazon	user-churn		Test	0.7038	0.7054
			Val	0.7033	0.7044
				0.7043	

Table 9: Results on the entity regression tasks in RelBench. Lower is better. Best values are in bold. Relative gains are expressed as percentage improvement over RDL baseline.

Dataset	Task	MAE ↓	RDL Baseline	HGT	HGT +PE	RelGT (ours)	% Rel. Gain
rel-f1	driver-position		Test	4.022	4.1598	4.2358	3.9170 2.61
			Val	3.193	3.3517	2.9894	3.3257
rel-avito	ad-ctr		Test	0.041	0.0441	0.0494	0.0345 15.85
			Val	0.037	0.0409	0.0456	0.0314
rel-event	user-attendance		Test	0.258	0.2635	0.2562	0.2502 2.79
			Val	0.255	0.2617	0.2574	0.2548
rel-trial	study-adverse		Test	44.473	43.3253	42.4622	43.9923 1.08
			Val	46.290	45.9957	45.7966	46.2148
rel-trial	site-success		Test	0.400	0.4374	0.4431	0.3263 18.43
			Val	0.401	0.4198	0.4245	0.3593
rel-amazon	user-ltv		Test	14.313	15.3804	15.9296	14.2665 0.32
			Val	12.132	13.1017	13.5599	12.1151
rel-amazon	item-ltv		Test	50.053	56.1384	55.6211	48.9222 2.26
			Val	45.1401	51.2139	50.3468	43.8161
rel-stack	post-votes		Test	0.065	0.0679	0.0680	0.0654 -0.62
			Val	0.059	0.0617	0.0618	0.0592
rel-hm	item-sales		Test	0.056	0.0655	0.0641	0.0536 4.29
			Val	0.065	0.0749	0.0735	0.0627

Table 10: Results on the entity classification tasks in RelBench. Higher is better. Best values are in bold. Relative gains are expressed as percentage improvement over RDL baseline.

Dataset	Task	AUC ↑	RDL Baseline	HGT	HGT +PE	RelGT (ours)	% Rel. Gain
rel-f1	driver-dnf	Test	0.7262	0.7142	0.7109	0.7587	4.48
		Val	0.7136	0.7678	0.7318	0.6804	
	driver-top3	Test	0.7554	0.6389	0.8340	0.8352	10.56
		Val	0.7764	0.6659	0.6079	0.7958	
rel-avito	user-clicks	Test	0.6590	0.6584	0.6387	0.6830	3.64
		Val	0.6473	0.5977	0.5656	0.6649	
	user-visits	Test	0.6620	0.6426	0.6507	0.6678	0.88
		Val	0.6965	0.6696	0.6732	0.7024	
rel-event	user-repeat	Test	0.7689	0.6717	0.6590	0.7609	-1.04
		Val	0.7125	0.6247	0.5974	0.7285	
	user-ignore	Test	0.8162	0.8348	0.8161	0.8157	-0.06
		Val	0.9170	0.8896	0.8940	0.8868	
rel-trial	study-outcome	Test	0.6860	0.5679	0.5691	0.6861	0.01
		Val	0.6818	0.5985	0.5925	0.6678	
rel-amazon	user-churn	Test	0.7042	0.6608	0.6589	0.7039	-0.04
		Val	0.7045	0.6639	0.6622	0.7036	
	item-churn	Test	0.8281	0.7824	0.7840	0.8255	-0.31
		Val	0.8239	0.7845	0.7846	0.8220	
rel-stack	user-engagement	Test	0.9021	0.8898	0.8852	0.9053	0.35
		Val	0.9059	0.8914	0.8847	0.9033	
	user-badge	Test	0.8986	0.8652	0.8518	0.8632	-3.94
		Val	0.8886	0.8760	0.8691	0.8741	
rel-hm	user-churn	Test	0.6988	0.6773	0.6491	0.6927	-0.87
		Val	0.7042	0.6814	0.6502	0.6988	