

T-GAP: Learning to Walk across Time for Temporal Knowledge Graph Completion

Jaehun Jung¹, Jinhong Jung², U Kang¹

¹Department of Computer Science, Seoul National University

²Department of Computer Science, Jeonbuk National University

sharkmir1@snu.ac.kr, jinhongjung@jbnu.ac.kr, ukang@snu.ac.kr

Abstract

Temporal knowledge graphs (TKGs) inherently reflect the transient nature of real-world knowledge, as opposed to static knowledge graphs. Naturally, automatic TKG completion has drawn much research interests for a more realistic modeling of relational reasoning. However, most of the existing models for TKG completion extend static KG embeddings that do not fully exploit TKG structure, thus lacking in 1) accounting for temporally relevant events already residing in the local neighborhood of a query, and 2) path-based inference that facilitates multi-hop reasoning and better interpretability. In this paper, we propose T-GAP, a novel model for TKG completion that maximally utilizes both temporal information and graph structure in its encoder and decoder. T-GAP encodes query-specific substructure of TKG by focusing on the temporal displacement between each event and the query timestamp, and performs path-based inference by propagating attention through the graph. Our empirical experiments demonstrate that T-GAP not only achieves superior performance against state-of-the-art baselines, but also competently generalizes to queries with unseen timestamps. Through extensive qualitative analyses, we also show that T-GAP enjoys from transparent interpretability, and follows human intuition in its reasoning process.

Introduction

Knowledge graph (KG), due to its expressiveness over structured knowledge, has been widely used in various applications including recommender system (Wang et al. 2019), information retrieval (Liu et al. 2018), and question answering (Zhang et al. 2018). Moreover, the inherent sparseness of KGs gave rise to research interests over automatic knowledge graph completion, predicting missing entity for incomplete queries in form of (subject, predicate, ?).

Recent advancements in KG completion tasks have extended to a more challenging domain of temporal knowledge graphs (TKGs), as they model realistic events that are only temporarily valid. Triples in temporal graphs are annotated with corresponding time token, taking form of (subject, predicate, object, timestamp). Naturally, TKG completion task can be formulated as predicting missing tail entity for queries in form of (subject, predicate, ?, timestamp).

Majority of existing approaches to TKG completion propose a straightforward extension of conventional KG embeddings onto the temporal graphs (Dasgupta et al. 2018;

Query: (COVID-19, infects, ?, 12/20)

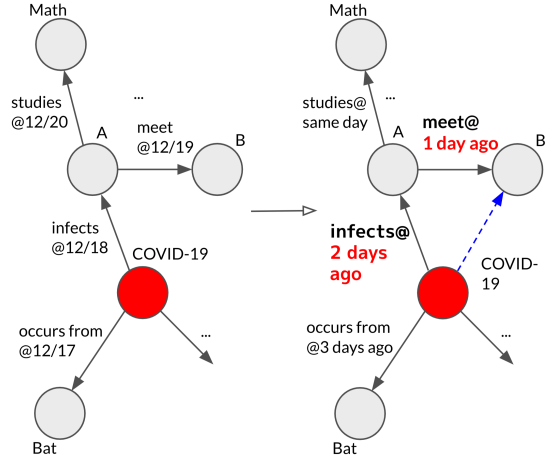


Figure 1: Example of temporal displacement. Edges in bold are important events relevant to the input query.

Lacroix, Obozinski, and Usunier 2019). Although these approaches are generally successful in TKG completion, we observe two major rooms of improvement from the existing models, each from the encoding phase, and the decoding phase.

In the encoding phase, a model could benefit from the rich neighborhood information residing in the structure of TKGs. Extracting, and encoding query-relevant information from the neighborhood nodes and their associated edges would help in fine-grained modeling of entity representation. The importance of neighborhood encoding has already been appreciated in static KGs (Chang et al. 2018; Nathani et al. 2019), but extension of these models to TKG is non-trivial, due to the additional time dimension in each triple.

Next, in the decoding phase, relational reasoning on TKG could leverage path-based inference. Several works adopted path-traversal model in static KGs (Xu et al. 2019; Das et al. 2018), showing preferable performance in relational reasoning compared to embedding-based models. Although path-based inference helps in capturing long-term dependency between nodes and gives better interpretability over model’s reasoning process, these approaches are yet to be examined in TKG completion tasks.

To this end, we propose T-GAP (Temporal GNN with Attention Propagation), a novel model for TKG completion, that tackles both challenges stated above. In the encoder, we introduce a new type of temporal graph neural network (GNN), which attentively aggregates query-relevant information from each entity’s local neighborhood. Specifically, we focus on encoding the temporal displacement between the timestamps of the input query and each edge being encoded. An intuitive example is presented in Figure 1. Evidently, the two most important facts to answer the given query (COVID-19, infects, ?, 12/20), are that A has been infected to COVID-19 at 12/18, and A met B at 12/19. Here, one should note that the valuable information lies in the fact that A got infected *2 days before* the time of interest, not that he was infected at a *specific day* of 12/18. What matters most when accounting for temporal events, is the relative displacement between the event and the time of interest, rather than the absolute time of the event. To effectively capture the temporal displacement, our proposed encoder separately encodes both the sign of the displacement (i.e. whether the time of the event belongs to past, present, or future), and the magnitude of the displacement (i.e. how far is the event from the time of our interest).

Also, T-GAP performs a generalized path-based inference over TKG, based on the notion of Attention Flow (Xu et al. 2018). In each decoding step, our model explores KG by propagating attention value at each node to its reachable neighbor nodes, rather than sampling one node to walk from the neighborhood. The soft approximation of path traversal with attention propagation not only allows our model to be easily trained with end-to-end supervised learning, but also provides better interpretation over its reasoning process, compared to embedding-based models.

In summary, our contributions are as follows:

- We introduce a new GNN encoder that effectively captures query-relevant information from temporal KGs.
- Based on the encoder, we present T-GAP, a novel path-based TKG reasoning model. We examine T-GAP in 3 benchmark datasets in TKG completion task, and the quantitative metrics show clear improvement in all benchmarks compared to the state-of-the-art baselines.
- By analyzing the inferred attention distribution, we show that T-GAP possesses clear interpretability over its reasoning process, which has not yet been extensively discussed in TKG domain.

Related Work

Various approaches have been made toward automatic completion of static KGs. Majority of conventional approaches propose embedding-based models, including translative (Bordes et al. 2013; Wang et al. 2014), and factorization-based models (Yang et al. 2014; Trouillon et al. 2016). To complement for the weak representation power of KG embeddings, several recent works incorporate neural network either to the scoring function, or as an additional encoding layer. ConvE (Dettmers et al. 2018) adopts convolutional layer to model sophisticated interaction between entities in the scoring function. KBGAT (Nathani et al. 2019),

and RGhat (Zhang et al. 2020) adopt a variant of graph attention network to contextualize entity embedding with the corresponding neighborhood structure. DPMPN (Xu et al. 2019) employs two GNNs to encode both the original graph and an induced subgraph, for a scalable learning of KG structure. We extend the neighborhood encoding scheme of these prior works to temporal graphs, especially focusing on query dependent encoding of KG structure.

Existing works on TKG completion primarily focus on extending static KG embedding to dynamic graphs. Different models mainly differ in how to represent independent timestamps, and incorporate it to their scoring functions. HyTE (Dasgupta et al. 2018) extends TransH (Wang et al. 2014), projecting entity and relation embedding to time-specific hyperplane. García-Durán et al. (2018) propose to represent temporal relation as a sequence of relation type and characters in the timestamp, and encode the sequence using RNN. TComplEx (Lacroix, Obozinski, and Usunier 2019) considers the score of each triple as canonical decomposition of order 4 tensors in complex domain, adding time embedding to the order 3 decomposition of ComplEx. Goel et al. (2019) suggest to learn entity representation that changes over time, transforming part of the embedding with sinusoidal activation of learned frequencies.

Meanwhile, path-based reasoning has been actively employed for node prediction on knowledge graphs. Lin et al. (2015) infuse multi-hop path information into entity representation, using additive composition of relation embeddings. Das et al. (2018) and Lin, Socher, and Xiong (2018) consider KG reasoning problem as partially observed Markov Decision Process, training a policy network that starts traversing from the query head and reaches at the predicted tail entity. Xu et al. (2018) propose a soft approximation of path traversal with attention distribution. T-GAP aligns with these line of works by exploring informative paths relevant to the input query with attention propagation, to maximally utilize the KG structure during the decoding process.

Lastly, we find connections to our work from recently suggested GNNs for dynamic graphs. Pareja et al. (2020) propose a variant of graph convolutional network (GCN) suited for TKG, employing RNN to evolve GCN parameters across time. Han et al. (2020) present Graph Hawkes Neural Network, modeling temporal dependency between events with Hawkes process. While these works focus on modeling the evolution of the graph as whole, we discuss a new methodology of encoding temporal displacement between each event and the input query, which better suits our goal to explore query-relevant paths and reach the answer node.

Proposed Method

Overview

First, we denote TKG as $G_{KG} = \{(v, r, u, t)\} \subseteq V_{KG} \times R_{KG} \times V_{KG} \times T_{KG}$, where V_{KG} is a set of entities, R_{KG} is a set of relations, and T_{KG} is a set of timestamps associated with the relations. Given the graph G_{KG} and query $\mathbf{q} = (v_{query}, r_{query}, ?, t_{query})$, TKG completion is formulated as predicting $u \in V_{KG}$ most probable to fill in the query. We

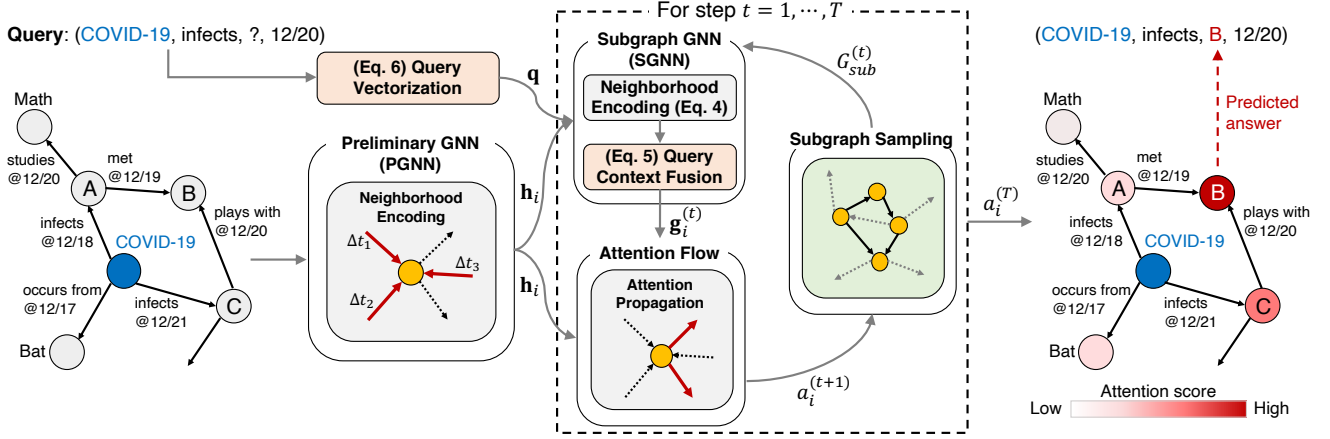


Figure 2: Overview of T-GAP. Starting from the query head, T-GAP explores relevant nodes and edges by iteratively propagating attention, and reaches at the target entity after the final propagation step.

also denote \overleftarrow{N}_i as a set of incoming neighbor nodes of v_i , i.e. nodes possessing edges toward v_i , and \overrightarrow{N}_i as a set of outgoing neighbor nodes of v_i .

Figure 2 illustrates an overview of T-GAP’s reasoning process. T-GAP consists of 4 sub-modules: Preliminary GNN (PGNN), Subgraph GNN (SGNN), Subgraph Sampling, and Attention Flow. Given G_{KG} and the query, in the encoding phase, T-GAP first uses PGNN to create preliminary node feature \mathbf{h}_i for all entities in the G_{KG} .

Next, at each decoding step $t = 1, \dots, T$, T-GAP iteratively samples a subgraph $G_{sub}^{(t)}$ from G_{KG} , that consists only of query-relevant nodes and edges. For each entity i included in $G_{sub}^{(t)}$, SGNN creates query-dependent node feature $\mathbf{g}_i^{(t)}$, incorporating the query vector \mathbf{q} and the preliminary feature \mathbf{h}_i . As noted in Xu et al. (2019), this additional encoding of induced subgraph not only helps in extracting only the query-relevant information from the original graph, but also scales up GNN by input-dependent pruning of irrelevant edges. Using both \mathbf{h}_i and $\mathbf{g}_i^{(t)}$, Attention Flow computes transition probability to propagate the attention value of each node to its reachable neighbor nodes, creating the next step’s node attention distribution $a_i^{(t+1)}$. After the final propagation step T , the answer to the input query is inferred as the node with the highest attention value $a_i^{(T)}$.

Preliminary GNN

Given G_{KG} , T-GAP first randomly initializes node feature \mathbf{h}_i for all $v_i \in V_{KG}$. Then, to contextualize the representation of entities in G_{KG} with the graph structure, each layer in PGNN updates node feature \mathbf{h}_i of entity v_i by attentively aggregating v_i ’s neighborhood information. The important intuition underlying PGNN is that the temporal displacement between timestamps of the query and each event is integral to capture the time-related dynamics of each entity. Therefore, for each timestamp t_e of edge e in G_{KG} , we resolve to separately encode the sign and magnitude of the temporal displacement $\Delta t_e = t_e - t_{query}$. Concretely, PGNN computes message \mathbf{m}_{ij} from entity v_i to v_j as fol-

lows:

$$\mathbf{m}_{ij} = \mathbf{W}_{\lambda(\Delta t_{ij})}(\mathbf{h}_i + \rho_{ij} + \tau_{|\Delta t_{ij}|})$$

$$\text{where } \mathbf{W}_{\lambda(\Delta t_{ij})} = \begin{cases} \mathbf{W}_{past} & \text{if } \Delta t_{ij} < 0 \\ \mathbf{W}_{present} & \text{if } \Delta t_{ij} = 0 \\ \mathbf{W}_{future} & \text{if } \Delta t_{ij} > 0 \end{cases} \quad (1)$$

ρ_{ij} is a relation-specific parameter associated with r_{ij} which denotes the relation that connects v_i to v_j . In addition to the entity and relation, we learn the discretized embedding of size of temporal displacement, i.e. $\tau_{|\Delta t_{ij}|}$. We take account of the sign of displacement by applying sign-specific weight for each event.

Next, the new node feature \mathbf{h}'_j is computed by attention weighted sum of all incoming messages to v_j :

$$\mathbf{h}'_j = \sum_{i \in \overleftarrow{N}_j} a_{ij} \mathbf{m}_{ij},$$

$$a_{ij} = \text{softmax}_i(\alpha_{ij}),$$

$$\alpha_{ij} = \text{LeakyReLU}((\mathbf{W}_Q \mathbf{h}_j)^\top (\mathbf{W}_K \mathbf{m}_{ij}))$$

$$(2)$$

The attention values are computed by applying softmax over all incoming edges of v_j , with \mathbf{h}_j as query and \mathbf{m}_{ij} as key.

In addition, we extend this attentive aggregation scheme to multi-headed attention, which helps to stabilize the learning process and jointly attend to different representation subspaces (Veličković et al. 2017). Hence our message aggregation scheme is modified to:

$$\mathbf{h}'_j = \parallel \sum_{k=1}^K a_{ij}^k \mathbf{m}_{ij}^k \quad (3)$$

concatenating independently aggregated neighborhood features from each attention heads, where K is a hyperparameter indicating the number of attention heads.

Subgraph GNN

At each decoding step t , SGNN updates node feature \mathbf{g}_i for all entities that are included in the induced subgraph of current step, $G_{sub}^{(t)}$. We present the detailed procedure of subgraph sampling in upcoming section. Essentially, SGNN not

only contextualizes \mathbf{g}_i with respective incoming edges, but also infuses the query context vector with the entity representation. First, the subgraph features for entities newly added to the subgraph, are initialized to their corresponding preliminary features \mathbf{h}_j . Next, SGNN performs message propagation, using the same message computation and aggregation scheme as PGNN (Eq. 1-3), but with separate parameters:

$$\tilde{\mathbf{g}}'_j = \big\| \sum_{k=1}^K \sum_{i \in \tilde{N}_j} a_{ij}^k \mathbf{m}_{ij}^k \quad (4)$$

This creates an intermediate node feature $\tilde{\mathbf{g}}'_j$. The intermediate features are then concatenated with query context vector \mathbf{q} , and linear-transformed back to the node embedding dimension, creating new feature \mathbf{g}'_j :

$$\mathbf{g}'_j = \mathbf{W}_g [\tilde{\mathbf{g}}'_j \parallel \mathbf{q}] \quad (5)$$

$$\mathbf{q} = \mathbf{W}_c \times \text{LeakyReLU}(\mathbf{W}_{\text{present}}(\mathbf{h}_{\text{query}} + \rho_{\text{query}})) \quad (6)$$

where $\mathbf{h}_{\text{query}}$ is the preliminary feature of v_{query} , and ρ_{query} is the relation parameter for r_{query} .

Attention Flow

T-GAP models path traversal with the soft approximation of attention flow, iteratively propagating the attention value of each node to its outgoing neighbor nodes. Initially, the node attention is initialized to 1 for v_{query} , and 0 for all other entities. Hereafter, at each step t , Attention Flow propagates edge attention $\tilde{a}_{ij}^{(t)}$ and aggregates them to node attention $a_j^{(t)}$:

$$\begin{aligned} \tilde{a}_{ij}^{(t+1)} &= \mathcal{T}_{ij}^{(t+1)} a_i^{(t)}, \quad a_j^{(t+1)} = \sum_{i \in \tilde{N}_j} \tilde{a}_{ij}^{(t+1)} \\ \text{s.t. } \sum_i a_i^{(t+1)} &= 1, \quad \sum_j \tilde{a}_{ij}^{(t+1)} = 1 \end{aligned}$$

The key here is the transition probability \mathcal{T}_{ij} . In this work, we define \mathcal{T}_{ij} as applying softmax over the sum of two scoring terms, regarding both the preliminary feature \mathbf{h} , and the subgraph feature \mathbf{g} :

$$\begin{aligned} \mathcal{T}_{ij}^{(t+1)} &= \text{softmax}_j(\text{score}(\mathbf{g}_i^{(t)}, \mathbf{g}_j^{(t)}, \rho_{ij}, \tau_{|\Delta t_{ij}|}) + \\ &\quad \text{score}(\mathbf{g}_i^{(t)}, \mathbf{h}_j, \rho_{ij}, \tau_{|\Delta t_{ij}|})), \\ \text{score}(\mathbf{i}, \mathbf{j}, \mathbf{r}, \tau) &= \sigma((\mathbf{W}_Q \mathbf{i})^\top (\mathbf{W}_K (\mathbf{j} + \mathbf{r} + \tau))) \end{aligned}$$

The first scoring term accounts only for subgraph feature \mathbf{g}_i and \mathbf{g}_j , giving additional point to entities that are already included in the subgraph (note that \mathbf{g}_i is initialized to zero for entities not yet included in the subgraph). Meanwhile, the second scoring term could be regarded as *exploring term*, as it relatively prefers entities not included in the subgraph, by modeling the interaction between \mathbf{g}_i and \mathbf{h}_j .

As T-GAP consists only of differentiable operations, the path traversal of T-GAP can be trained end-to-end by directly supervising on the node attention distribution after T propagation steps. We train T-GAP to maximize the log probability of the answer entity u_{label} at step T .

$$\mathcal{L} = -\log a_{u_{\text{label}}}^{(T)}$$

Subgraph Sampling

The decoding process of T-GAP depends on the iterative sampling of query-relevant subgraph $G_{\text{sub}}^{(t)}$. The initial subgraph $G_{\text{sub}}^{(0)}$ before the first propagation step contains only one node, v_{query} . As the propagation step proceeds, edges with high relevance to the input query, measured by the size of attention value assigned to the edges, are added to the previous step's subgraph. Specifically, the subgraph sampling at step t proceeds as follows:

- Find x number of *core nodes* with highest (nonzero) node attention value $a_i^{(t-1)}$ at the previous step.
- For each of the core node, sample y number of edges that originate from the node.
- Among $x \cdot y$ sampled edges, find z number of edges with highest edge attention value $\tilde{a}_{ij}^{(t)}$ at the current step.
- Add the z edges to $G_{\text{sub}}^{(t-1)}$.

In this module, x, y, z are hyperparameters. Intuitively, we only collect 'important' events that originate from 'important' entities (core nodes) with respect to the query, while keeping the subgraph size under control (edge sampling). We provide an illustrative example on subgraph sampling in Appendix A.

Note that although edge sampling brings in stochasticity to T-GAP's inference, this does not hinder the end-to-end training of the model. Since the sampling is not parameterized and we only use node feature \mathbf{g} from the sampled subgraph, gradients back-propagate through \mathbf{g} , not through the sampling operation.

Experiment

Datasets

We evaluate our proposed method on three benchmark datasets for TKG completion: ICEWS14, ICEWS05-15, and Wikidata11k, all suggested by García-Durán, Dumančić, and Niepert (2018). ICEWS14 and ICEWS05-15 are subsets of ICEWS¹, each containing socio-political events in 2014, and from 2005 to 2015 respectively. Wikidata11k is a subset of Wikidata², possessing facts of various timestamps that span from A.D. 20 to 2020. All facts in Wikidata11k are annotated with additional temporal modifier, *occurSince* or *occurUntil*. For the sake of consistency and simplicity, we follow García-Durán, Dumančić, and Niepert (2018) to merge the modifiers into predicates rather than modeling them in separate dimension (e.g. (A, loves, B, since, 2020) transforms to (A, loves-since, B, 2020)). Detailed statistics of the three datasets are provided in Appendix B.

Baselines

We compare T-GAP with representative static KG embeddings - TransE and DistMult, and state-of-the-art embedding-based baselines on temporal KGs, including ConT (Ma et al. 2019), TTransE (Jiang et al. 2016), HyTE (Dasgupta

¹<https://dataverse.harvard.edu/dataverse/icews>

²https://www.wikidata.org/wiki/Wikidata:Main_Page

Model	ICEWS14				ICEWS05-15				Wikidata11k			
	MRR	Hits@1	Hits@3	Hits@10	MRR	Hits@1	Hits@3	Hits@10	MRR	Hits@1	Hits@3	Hits@10
TransE [▼]	0.280	9.4	-	63.7	0.294	9.0	-	66.3	0.316	18.1	-	65.9
DistMult [▼]	0.439	32.3	-	67.2	0.456	33.7	-	69.1	0.316	18.1	-	66.1
ConT [◊]	0.185	11.7	20.5	31.5	0.163	10.5	18.9	27.2	-	-	-	-
TTransE [▼]	0.255	7.4	-	60.1	0.271	8.4	-	61.6	0.488	33.9	-	80.6
HyTE [◊]	0.297	10.8	41.6	65.5	0.316	11.6	44.5	68.1	-	-	-	-
TA-TransE [▼]	0.275	9.5	-	62.5	0.299	9.6	-	66.8	0.484	32.9	-	80.7
TA-DistMult [▼]	0.477	36.3	-	68.6	0.474	34.6	-	72.8	0.700	65.2	-	78.5
DE-Simple [◊]	0.526	41.8	59.2	72.5	0.513	39.2	57.8	74.8	0.310	18.4	31.8	62.5
TComplEx	0.560	47.0	61.0	73.0	0.580	49.0	64.0	76.0	0.731	67.3	76.2	84.5
TNTComplEx	0.560	46.0	61.0	74.0	0.600	50.0	65.0	78.0	0.718	65.4	74.9	85.6
T-GAP	0.610	50.9	67.7	79.0	0.670	56.8	74.3	84.5	0.778	69.7	84.4	90.3

Table 1: T-GAP outperforms baselines in all three benchmarks, over all metrics. We used official implementation of DE-Simple, and T(NT)ComplEx for Wikidata11k. Other results with [▼] are from García-Durán, Dumančić, and Niepert (2018), results with [◊] are from Goel et al. (2019), and results on T(NT)ComplEx are from Lacroix, Obozinski, and Usunier (2019).

et al. 2018), TA (García-Durán, Dumančić, and Niepert 2018), DE-Simple (Goel et al. 2019), and T(NT)ComplEx (Lacroix, Obozinski, and Usunier 2019).

Experimental Setting

For each dataset, we create G_{KG} with only the triples in the train set. We add inverse edges to G_{KG} for proper path-based inference on reciprocal relations. Also, we follow Xu et al. (2019) by adding self-loops to all entities in the graph, allowing the model to stay at the ‘answer node’ if it reaches an optimal entity in $t < T$ steps. To measure T-GAP’s performance in head entity prediction, we add reciprocal triples to valid and test sets too. For all datasets, we find through empirical evaluation that setting the maximal path length $T = 3$ results in the best performance. Following previous works, we fix the dimension of entity / relation / displacement embedding to 100. Except for the embedding size, we search for the best set of hyperparameters using grid-based search, choosing the value with the best *Hits@1* while all other hyperparameters are fixed. We implement T-GAP with PyTorch and DGL, and plan to make the code publicly available. We provide further implementation details including hyperparameter search bounds and the best configuration in Appendix C.

Results

Table 1 shows the overall evaluation result of T-GAP against baseline methods. Along with Hits@1, 3, 10, we report MRR of the ground-truth entity, compared to the baseline methods. As seen in the table, T-GAP outperforms baseline models in all the benchmarks, improving up to 10% relative performance consistently over all metrics. We find through ablation study that the improvements are mainly contributed from resolving the two shortcomings of pre-existing TKG embeddings, which we indicate in earlier sections - absence of 1) neighborhood encoding scheme, 2) path-based inference with scalable subgraph sampling.

Model Variants and Ablation To further examine the effect of our proposed method in solving the aforementioned two problems, we conduct an ablation study as shown in Table 2. First, we consider T-GAP without temporal displacement encoding. In both PGNN and SGNN, we do not con-

Model	MRR	Hits@1	Hits@3	Hits@10
No Displacement	0.477	35.7	53.9	71.5
No Subgraph	0.598	49.3	66.8	78.5
No PGNN	0.59	49.9	66.5	78.4
2-layer PGNN	0.605	50.2	67.3	78.9
T-GAP	0.607	50.4	67.5	78.7

Table 2: Ablation study results on ICEWS14, compared to the best configuration. Results annotated ‘T-GAP’ are the best performance of T-GAP with temporal displacement encoding, 1-layer PGNN, and subgraph sampling.

Model	MRR	Hits@1	Hits@3	Hits@10
DE-Simple	0.434	33.3	49.2	62.4
TComplEx	0.443	34.8	49.2	62.5
TNTComplEx	0.444	34.6	49.4	63.5
T-GAP	0.483	37.2	54.6	69.0

Table 3: Generalization performance over unseen timestamps in ICEWS14. Accounting for relative displacement rather than independent timestamps, our model is the most robust to queries with unseen timestamps.

sider the sign and magnitude of temporal displacement, and simply learn the embedding of each timestamp as is. While computing the message \mathbf{m}_{ij} , the two GNNs simply add the time embedding $\tau_{t_{ij}}$ to $\mathbf{h}_i + \rho_{ij}$. No sign-specific weight is multiplied, and all edges are linear-transformed with a same weight matrix. In this setting, T-GAP’s performance on ICEWS14 degrades about 30% in Hits@1, and performs similar to TA-DistMult in Table 1. The result attests to the importance of temporal displacement for the neighborhood encoding in temporal KGs.

Next, to analyze the effect of subgraph sampling on overall performance, we resort to a new configuration of T-GAP where no subgraph sampling is applied, and SGNN creates node feature \mathbf{g}_i for all entities in G_{KG} . Here, T-GAP’s performance slightly degrades about 1 percent in all metrics. This implies the importance of subgraph sampling to prune query-irrelevant edges, helping T-GAP to concentrate on the plausible substructure of the input graph.

Finally, we analyze the effect of PGNN by training T-

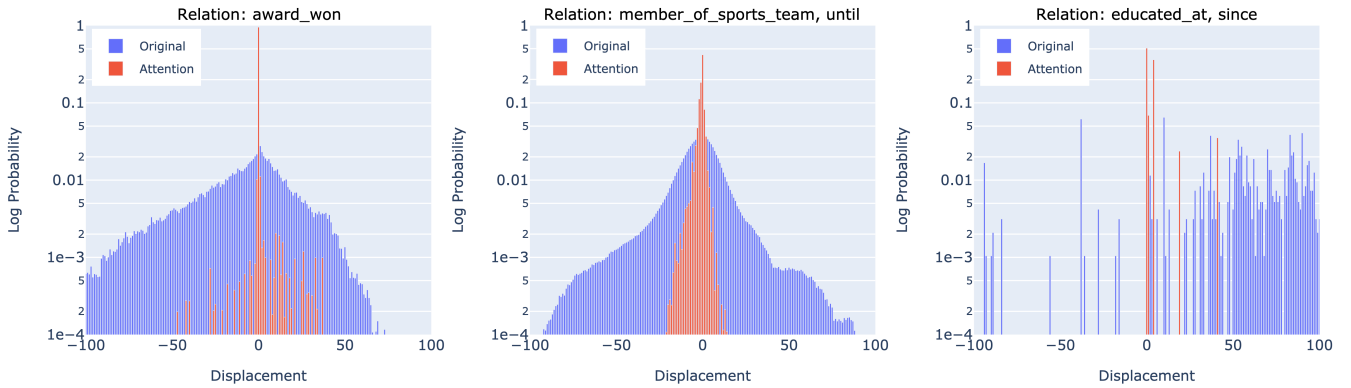


Figure 3: Edge attention distribution over different temporal displacements. Note that y-axis is in log scale. For better visualization, we cut the range of temporal displacement to $[-100, 100]$ in all charts. T-GAP learns to effectively attend on events with specific temporal distance from the query time, depending on the relation type of the query.

GAP with different numbers of PGNN layers. We find that T-GAP, trained with 1-layer PGNN, performs superior to the model without PGNN by absolute gain of 1% in MRR. However, adding up more layers in PGNN gives only a minor gain, or even aggravates the test set accuracy, mainly owing to early overfitting on the triples in train set.

Generalizing to Unseen Timestamps We conduct an additional study that measures the performance of T-GAP in generalizing to queries with unseen timestamps. Following Goel et al. (2019), we modify ICEWS14 by including all triples except those on 5th, 15th, 25th day of each month in the train set, and creating valid and test sets using only the excluded triples. The performance of T-GAP against the strongest baselines in this dataset are presented in Table 3. In this setting, DE-Simple and T(NT)Complex perform much more similar to each other than in Table 1, while T-GAP performs superior to all baselines. DE-Simple shows strength in generalizing over time, as it represents each entity as a continuous function over temporal dimension. However, the model is weak when the range of timestamps is large and sparse, as shown for Wikidata in Table 1. Meanwhile, TComplex and TNTComplex show fair performance in Wikidata, but poorly infer for unseen timestamps, as they only learn independent embeddings of discrete timestamps. On the contrary to these models, T-GAP not only shows superior performance in all benchmarks but also is robust to unseen timestamps, by accounting for the temporal displacement, not the independent time tokens.

Interpretation

We provide a detailed analysis on the interpretability of T-GAP in its relational reasoning process.

Relation Type and Temporal Displacement Intuitively, the query relation type, and the temporal displacement between relevant events and the query are closely correlated. For a query such as $(PersonX, member_of_sports_team, ?, t_1)$, events that happened 100 years before t_1 or 100 years after t_1 will highly likely be irrelevant. On the contrary, for a query given as $(NationX, wage_war_against, ?, t_2)$,

one might have to consider those events far-off the time of interest. To verify whether T-GAP understands this implicit correlation, we analyze the attention distribution over edges with different temporal displacements, when T-GAP is given input queries with a specific relation type.

The visualization of the distributions for three relation types are presented in Figure 3. For all queries in the test set of WikiData11k with a specific relation type, we visualize the average attention value assigned to edges with each temporal displacement (red bars). We compare this with the original distribution of temporal displacement, counted for all edges reachable in T steps from the head entity v_{query} (blue bars). Remarkably, on the contrary with the original distribution of high variance over wide range of displacement, T-GAP tends to focus most of the attention to edges with specific temporal displacement, depending on the relation type.

For queries with relation *award_won*, the attention distribution is extremely skewed, focusing over 90% of the attention to events with displacement = 0 (i.e. events in the same year with the query). Note that we have averaged the distribution for all queries with *award_won*, including both temporal modifiers *occurSince* and *occurUntil*. The skewed distribution mainly results from the fact that the majority of the ‘award’ entities in Wikidata11k are annual awards, such as *Latin Grammy Award*, or *Emmy Award*. The annual property of the candidate entities naturally makes T-GAP to focus on clues such as the nomination of awardees, or significant achievement of awardees in the year of interest.

Next, we test T-GAP for queries with relation *member_of_sports_team-occurUntil*. In this case, the attention is more evenly distributed than the former case, but slightly biased toward past events. We find that this phenomenon is mainly due to the existence of temporally reciprocal edge in G_{KG} , which is a crucial key in solving the given query. Here, T-GAP sends more than half of the attention value (on average) to an event with relation *member_of_sports_team-occurSince*, that happened few years before the time of interest. The inference follows our intuition to look for the last sports club where the player became member of, be-

Query	(North_Korea, threaten, ?, 2014/04/29)	
1st Step	(North_Korea, <i>threaten</i> , Japan, 2014/05/12)	0.057
	(North_Korea, <i>threaten</i> , Japan, 2014/12/18)	0.054
	(North_Korea, make_statement, South_Korea, 2014/04/29)	0.044
	(North_Korea, <i>make_an_appeal</i> , Japan, 2014/10/01)	0.041
	(North_Korea, <i>threaten</i> , South_Korea, 2014/08/01)	0.040
2nd Step	(South_Korea, <i>threaten</i> , North_Korea, 2014/04/22)	0.15
	(Japan, <i>release_person</i> , North_Korea, 2014/07/28)	0.078
	(South_Korea, <i>criticize_or_denounce</i> , North_Korea, 2014/04/28)	0.051
	(Japan, <i>express_intent_to_cooperate</i> , North_Korea, 2014/02/28)	0.039
3rd Step	(South_Korea, make_statement, North_Korea, 2014/04/28)	0.037
	(North_Korea, make_statement, South_Korea, 2014/04/29)	0.189
	(North_Korea, <i>accuse</i> , South_Korea, 2014/04/15)	0.121
	(North_Korea, <i>criticized_or_denounced_by</i> , South_Korea, 2014/04/28)	0.052
	(North_Korea, <i>deny_responsibility</i> , South_Korea, 2014/04/25)	0.032
Answer	(South_Korea, <i>release_person</i> , North_Korea, 2014/04/14)	0.021
	South_Korea	

Table 4: List of predominant edges for a case study. Numbers in the right are the corresponding edge attention assigned to each edge. Predicates in red color carry negative meaning, while predicates in blue color hold positive meaning. We find through the case study that the attention propagation allows T-GAP to fix its misleading focus on sub-optimal entities.

fore the timestamp of the query. The third case with relation *educated_at-occurSince* is the opposite of the second case. Majority of the attention have been concentrated on events in 1-5 years after the query time, searching for the first event with relation *educated_at-occurUntil*, that happened after the time of interest.

As the analysis suggests, T-GAP discovers important clues for each relation type, adequately accounting for the temporal displacement between the query and related events, while aligning with human intuition.

Case Study We resort to a case study, to provide a detailed view on T-GAP’s attention-based traversal. In this study, our model is given an input query (*North_Korea, threaten, ?, 2014/04/29*) in ICEWS14 where the correct answer is *South_Korea*. For each propagation step, we list top-5 edges that received the highest attention value in the step.

The predominant edges and their associated attention values are shown in Table 4. In the first step, T-GAP attends to various events pertinent to *North_Korea*, that mostly include negative predicates against other nations. As seen in the table, the two plausible query-filling candidates are *Japan*, and *South_Korea*. *Japan* receives slightly more attention than *South_Korea*, as it is associated with more relevant facts such as “*North_Korea* threatened *Japan* at May 12th”.

In the second step, however, T-GAP discovers additional relevant facts, that could be crucial in answering the given query. As these facts have either *Japan* or *South_Korea* as head entity, they could not be discovered in the first propagation step, which only propagates the attention from the query head *North_Korea*. T-GAP attends to the events (*South_Korea, threaten / criticize_or_denounce, North_Korea*) that happened only a few days before our time of interest. These facts imply the strained relationship between the two nations around the query time. Also, T-GAP finds that most of the edges that span from *Japan* to *North_Korea* before/after few months the time of interest, tend to be positive events. As a result, in the last step, T-GAP propagates most of the node attention in *North_Korea*

to the events associated with *South_Korea*. Highest attention have been assigned to the relation *make_statement*. Although the relation itself does not hold a negative meaning, in ICEWS14, *make_statement* is typically accompanied by *threaten*, as entities do formally threaten other entities by making statements.

Through the case study, we find that T-GAP leverages the propagation-based decoding as a tool to fix its traversal over wrongly-selected entities. Although *Japan* seemed like an optimal answer in the first step, it understands through the second step that the candidate was sub-optimal with respect to the query, propagating the attention assigned to *Japan* back to *North_Korea*. T-GAP fixes its attention propagation in the last step, resulting in a completely different set of attended events compared to the first step. Such an amendment would not have been possible with conventional approaches to path-based inference, which greedily select an optimal entity to traverse at each decoding step.

Conclusion

In this paper, we propose a novel approach to TKG completion named T-GAP, which explores a query-relevant sub-structure of TKG with attention propagation. Unlike other embedding-based models, the proposed method effectively gathers useful information from the existing KG, by accounting for the temporal displacement between the query and respective edges. Quantitative results show that T-GAP not only achieves state-of-the-art performance consistently over all three benchmarks, but also competently generalizes to queries with unseen timestamps. Through extensive analysis, we also show that the propagated attention distribution well serves as an interpretable proxy of T-GAP’s reasoning process that aligns with human intuition.

References

Bordes, A.; Usunier, N.; Garcia-Duran, A.; Weston, J.; and Yakhnenko, O. 2013. Translating embeddings for model-

ing multi-relational data. In *Advances in neural information processing systems*, 2787–2795.

Chang, J.; Gu, J.; Wang, L.; Meng, G.; Xiang, S.; and Pan, C. 2018. Structure-aware convolutional neural networks. In *Advances in neural information processing systems*, 11–20.

Das, R.; Dhuliawala, S.; Zaheer, M.; Vilnis, L.; Durugkar, I.; Krishnamurthy, A.; Smola, A.; and McCallum, A. 2018. Go for a Walk and Arrive at the Answer: Reasoning Over Paths in Knowledge Bases using Reinforcement Learning. In *International Conference on Learning Representations*.

Dasgupta, S. S.; Ray, S.; Nath; and Talukdar, P. 2018. Hyte: Hyperplane-based temporally aware knowledge graph embedding. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, 2001–2011.

Dettmers, T.; Minervini, P.; Stenetorp, P.; and Riedel, S. 2018. Convolutional 2d knowledge graph embeddings. In *Thirty-Second AAAI Conference on Artificial Intelligence*.

García-Durán, A.; Dumančić, S.; and Niepert, M. 2018. Learning Sequence Encoders for Temporal Knowledge Graph Completion. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, 4816–4821. Brussels, Belgium: Association for Computational Linguistics. URL <https://www.aclweb.org/anthology/D18-1516>.

Goel, R.; Kazemi, S. M.; Brubaker, M.; and Poupart, P. 2019. Diachronic embedding for temporal knowledge graph completion. *arXiv preprint arXiv:1907.03143*.

Han, Z.; Wang, Y.; Ma, Y.; Günnemann, S.; and Tresp, V. 2020. Graph Hawkes Neural Network for Future Prediction on Temporal Knowledge Graphs. In *Automated Knowledge Base Construction*.

Jiang, T.; Liu, T.; Ge, T.; Sha, L.; Chang, B.; Li, S.; and Sui, Z. 2016. Towards time-aware knowledge graph completion. In *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*, 1715–1724.

Lacroix, T.; Obozinski, G.; and Usunier, N. 2019. Tensor Decompositions for Temporal Knowledge Base Completion. In *International Conference on Learning Representations*.

Lin, X. V.; Socher, R.; and Xiong, C. 2018. Multi-Hop Knowledge Graph Reasoning with Reward Shaping. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, 3243–3253. Brussels, Belgium: Association for Computational Linguistics.

Lin, Y.; Liu, Z.; Luan, H.; Sun, M.; Rao, S.; and Liu, S. 2015. Modeling Relation Paths for Representation Learning of Knowledge Bases. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, 705–714. Lisbon, Portugal: Association for Computational Linguistics.

Liu, Z.; Xiong, C.; Sun, M.; and Liu, Z. 2018. Entity-Duet Neural Ranking: Understanding the Role of Knowledge Graph Semantics in Neural Information Retrieval. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*,

2395–2405. Melbourne, Australia: Association for Computational Linguistics.

Ma, Y.; Tresp; Volker; and Daxberger, E. A. 2019. Embedding models for episodic knowledge graphs. *Journal of Web Semantics* 59: 100490.

Nathani, D.; Chauhan, J.; Sharma, C.; and Kaul, M. 2019. Learning Attention-based Embeddings for Relation Prediction in Knowledge Graphs. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Florence, Italy: Association for Computational Linguistics. URL <https://www.aclweb.org/anthology/P19-1466>.

Pareja, A.; Domeniconi, G.; Chen, J.; Ma, T.; Suzumura, T.; Kanezashi, H.; Kaler, T.; Schardl, T. B.; and Leiserson, C. E. 2020. EvolveGCN: Evolving Graph Convolutional Networks for Dynamic Graphs. In *Proceedings of the Thirty-Fourth AAAI Conference on Artificial Intelligence*.

Trouillon, T.; Welbl, J.; Riedel, S.; Gaussier, É.; and Bouchard, G. 2016. Complex embeddings for simple link prediction. In *International Conference on Machine Learning*. International Conference on Machine Learning (ICML).

Veličković, P.; Cucurull, G.; Casanova, A.; Romero, A.; Lio, P.; and Bengio, Y. 2017. Graph attention networks. *arXiv preprint arXiv:1710.10903*.

Wang, X.; Wang, D.; Xu, C.; He, X.; Cao, Y.; and Chua, T.-S. 2019. Explainable reasoning over knowledge graphs for recommendation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, 5329–5336.

Wang, Z.; Zhang, J.; Feng, J.; and Chen, Z. 2014. Knowledge Graph Embedding by Translating on Hyperplanes. In *AAAI*.

Xu, X.; Feng, W.; Jiang, Y.; Xie, X.; Sun, Z.; and Deng, Z.-H. 2019. Dynamically Pruned Message Passing Networks for Large-scale Knowledge Graph Reasoning. In *International Conference on Learning Representations*.

Xu, X.; Zu, S.; Gao, C.; Zhang, Y.; and Feng, W. 2018. Modeling Attention Flow on Graphs. *arXiv preprint arXiv:1811.00497*.

Yang, B.; Yih, W.-t.; He, X.; Gao, J.; and Deng, L. 2014. Embedding entities and relations for learning and inference in knowledge bases. *arXiv preprint arXiv:1412.6575*.

Zhang, Y.; Dai, H.; Kozareva, Z.; Smola, A. J.; and Song, L. 2018. Variational reasoning for question answering with knowledge graph. In *Thirty-Second AAAI Conference on Artificial Intelligence*.

Zhang, Z.; Zhuang, F.; Zhu, H.; Shi, Z.-P.; Xiong, H.; and He, Q. 2020. Relational Graph Neural Network with Hierarchical Attention for Knowledge Graph Completion. In *AAAI*, 9612–9619.

A Subgraph Sampling Example

Figure 4 is an illustrative example for the subgraph sampling procedure in T-GAP. The hyperparameters for the example are as follows: $x = 2$ (the maximum number of core nodes), $y = 3$ (the maximum number of candidate edges, considered

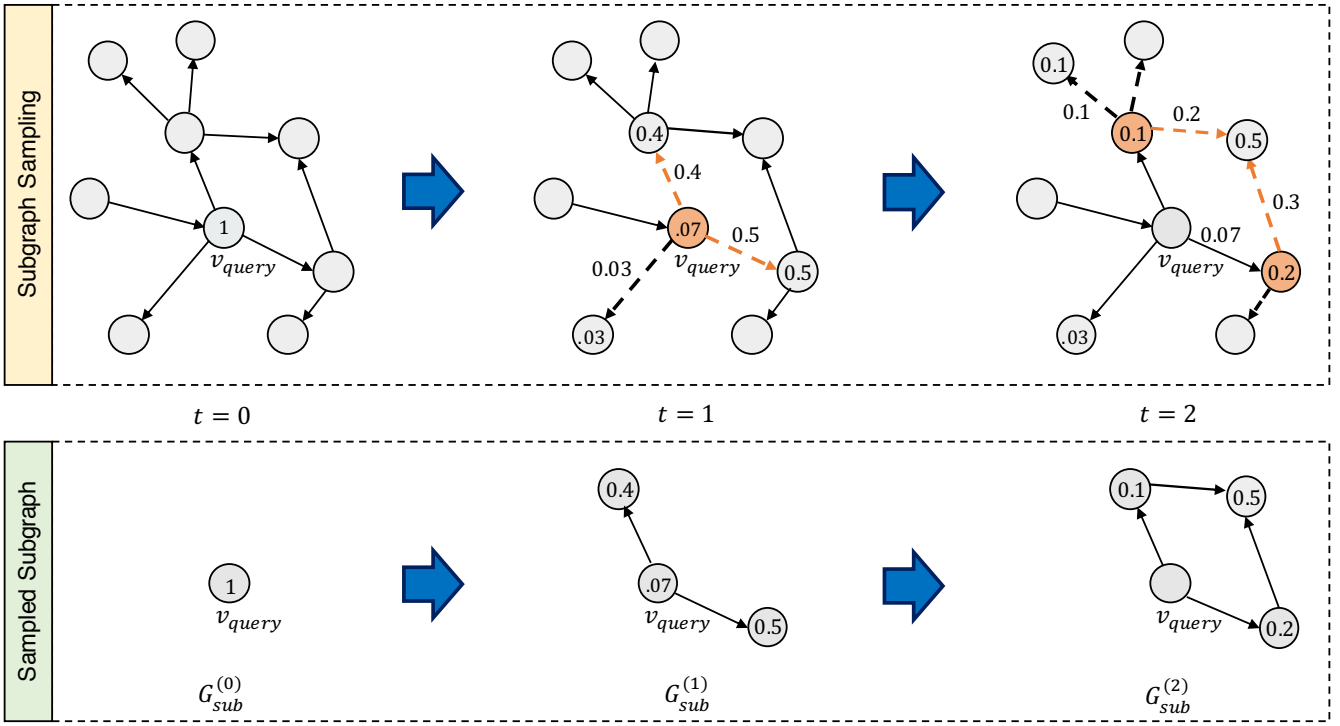


Figure 4: Example of subgraph sampling in T-GAP. The graphs above represent the respective node / edge attention distribution at the initial state ($t = 0$), after the first propagation step ($t = 1$), and after the second propagation step ($t = 2$). The graphs below show the sampled subgraph at each step t . x ($= 2$) orange nodes are the core nodes retrieved at a step, and y ($= 3$) dashed edges from each core node are candidate edges for the sampled subgraph. Among the candidate edges, z ($= 2$) orange edges are newly added to the previous subgraph.

by each core node), and $z = 2$ (the number of sampled edges added to the subgraph).

In the initial state, the only node associated with nonzero attention $a_i^{(0)}$ is the query head v_{query} . Also, the initial subgraph $G_{sub}^{(0)}$ consists only of the node v_{query} . After the first propagation step $t = 1$, T-GAP first finds top- x core nodes (where $x = 2$) w.r.t. nonzero node attention scores $a_i^{(0)}$ of the previous step $t = 0$. Since the only node with nonzero attention value is v_{query} , it is retrieved as the core node. Next, T-GAP randomly samples at most $y = 3$ edges that originate from the core node (e.g. dashed edges with weights). Among the sampled edges, it selects top- z (where $z = 2$) edges in the order of edge attention values $\tilde{a}_{ij}^{(1)}$ at the current step; then, they are added to $G_{sub}^{(0)}$, resulting in the new subgraph $G_{sub}^{(1)}$.

After the second propagation step $t = 2$, T-GAP again finds x core nodes that correspond to highest attention values $a_i^{(1)}$ (e.g. nodes annotated with 0.1 and 0.2, respectively). Then, y outgoing edges for each core node are sampled; among $x \cdot y$ sampled edges, z edges with highest edge attention values $\tilde{a}_{ij}^{(2)}$ are added to $G_{sub}^{(1)}$, creating the new subgraph $G_{sub}^{(2)}$.

As seen in the figure, the incremental subgraph sampling scheme allows our model to iteratively expand the range of nodes and edges to attend, while guaranteeing that the criti-

cal nodes and edges in the previous steps are kept included in the latter subgraphs.

By flexibly adjusting the subgraph related hyperparameters x , y , and z , T-GAP is readily calibrated between reducing computational complexity and optimizing the predictive performance. Intuitively, with more core nodes, more sampled edges, and more edges added to the subgraph, T-GAP can better attend to the substructure of TKG that otherwise might have been discarded. Meanwhile, with small x , y , and z , T-GAP can easily scale-up to large graphs by reducing the number of message-passing operations in SGNN.

B Dataset Statistics

Dataset	ICEWS14	ICEWS05-15	Wikidata11k
$ V_{KG} $	7,128	10,488	11,134
$ R_{KG} $	230	251	95
$ T_{KG} $	365	4017	328
$ G_{KG} $	90,730	479,329	150,079
Time Span	2014	2005-2015	25-2020

Table 5: Dataset Statistics for ICEWS14, ICEWS05-15, and Wikidata11k. The unit of the time span is year.

C Additional Implementation Detail

Computing Infrastructure	Tesla V100 GPU	
Search Strategy	Manual Tuning	
Best Validation <i>Hits@1</i>	52.1 (ICEWS14), 56.8 (ICEWS05-15), 70.7 (Wikidata11k)	

Hyperparameter	Search Bound	Experimental Setting
<i>max path length T</i>	<i>choice</i> [2, 3, 4]	3
<i>number x of core nodes</i>	<i>choice</i> [5, 10, 50, 100]	100 (ICEWS14), 10 (ICEWS05-15), 50 (Wikidata11k)
<i>number y of edges to sample</i>	<i>choice</i> [10, 50, 100, 200, 500]	500 (ICEWS14), 100 (ICEWS05-15), 200 (Wikidata11k)
<i>number z of edges to add</i>	<i>choice</i> [10, 50, 100, 200, 500]	500 (ICEWS14), 100 (ICEWS05-15), 200 (Wikidata11k)
<i>number K of attention heads</i>	<i>choice</i> [3, 4, 5, 6]	5
<i>embedding dimension</i>	100	100
<i>number of epochs</i>	10	10
<i>batch size</i>	<i>choice</i> [8, 16, 32]	16 (ICEWS14), 8 (ICEWS05-15, Wikidata11k)
<i>optimizer</i>	<i>Adam</i>	<i>Adam</i>
<i>learning rate</i>	<i>loguniform-float</i> [5e-2, 5e-5]	5e-4
<i>lr scheduler</i>	<i>reduce_on_plateau</i>	<i>reduce_on_plateau</i>
<i>lr reduction factor</i>	0.1	0.1
<i>gradient clip norm</i>	<i>uniform-integer</i> [1, 5]	3 (ICEWS14, ICEWS05-15), 5 (Wikidata11k)

Table 6: Additional implementation detail of T-GAP. We report the hyperparameter search bounds and the used configurations, along with the experimental environments.