



Pre-trained Language Model for Temporal Knowledge Graph Completion

Wenying Feng¹, Angxiao Zhao^{1,2}, Jianming Li^{1,2}, Zhiqiang Zhang^{1,2}, Yan Jia³ and Zhaoquan Gu^{*,1,2}

Abstract

Large Language Models (LLMs) have been proven remarkable in natural language processing, which prompts numerous works on knowledge extraction, knowledge fusion, knowledge representation, and knowledge completion. Existing works mainly focus on static multi-relational knowledge graph (KG). Unlike static knowledge graph, temporal knowledge graph (TKG) contains temporal information and evolves over time. Learning and reasoning about the representation of temporal knowledge graph are more difficult. Training or fine-tuning LLMs for temporal graph related tasks incurs significant computational overhead and requires the design of prompts. Conducting tasks of TKG does not necessarily require such complex work. Therefore, we explore temporal knowledge graph completion (TKGC) based on pre-trained “small” language model. We propose **TKG-BERT** by applying BERT for temporal knowledge graph completion and classification. Specifically, We introduce three ways to model temporal knowledge in TKG-BERT: vanilla knowledge embedding (Van.), explicit time modeling (Exp.) and implicit time modeling (Imp.). TKG-BERT(Van.) only adopts static knowledge without embedding time information; TKG-BERT(Exp.) embeds timestamp in quadruple explicitly; TKG-BERT(Imp.) models time implicitly, by dividing the training set and testing set in chronological order. We conduct experiments on ICEWS14 and ICEWS05-15, which are two public temporal knowledge graph datasets. Various experiments of temporal knowledge graph completion and classification tasks show the effectiveness of pre-trained language model for TKG completion. We also compare the performance of TKG-BERT across different time modeling way and proportion of training set.

Keywords: Knowledge graph; knowledge graph representation; temporal knowledge graph; pre-trained language model

Introduction

Knowledge graph is a kind of graph-structured data for representing knowledge, which supports knowledge reasoning. Knowledge graphs are widely used in various applications such as information retrieval, recommendation, and natural language processing. However, most existing knowledge graphs are incomplete and need to be reasoned and completed. Regarding this purpose, Knowledge Graph Embedding (KGE) transforms knowledge into low-dimensional vector space, to obtain the form of quantifiable, computable and reasonable knowledge for knowledge graph completion. Existing KGE methods can be classified into two categories: static KGE and temporal KGE. Static knowledge embedding considers the representation of static triples without considering changes in entities and relations along with time. Temporal knowledge embedding focuses on modeling the dynamic evolutionary properties of knowledge.

Affiliation:

¹Department of New Networks, Pengcheng Laboratory

²School of Computer Science and Technology, Harbin Institute of Technology (Shenzhen), China

³School of Computer Science, National University of Defense Technology

[#]These authors contribute equally

*Corresponding author:

Zhaoquan Gu. School of Computer Science and Technology, Harbin Institute of Technology (Shenzhen), China.

Citation: Wenying Feng, Angxiao Zhao, Jianming Li, Zhiqiang Zhang, Yan Jia and Zhaoquan Gu. Pre-trained Language Model for Temporal Knowledge Graph Completion. Journal of Pharmacy and Pharmacology Research 10 (2026): 25-35.

Received: January 05, 2026

Accepted: January 12, 2026

Published: January 23, 2026

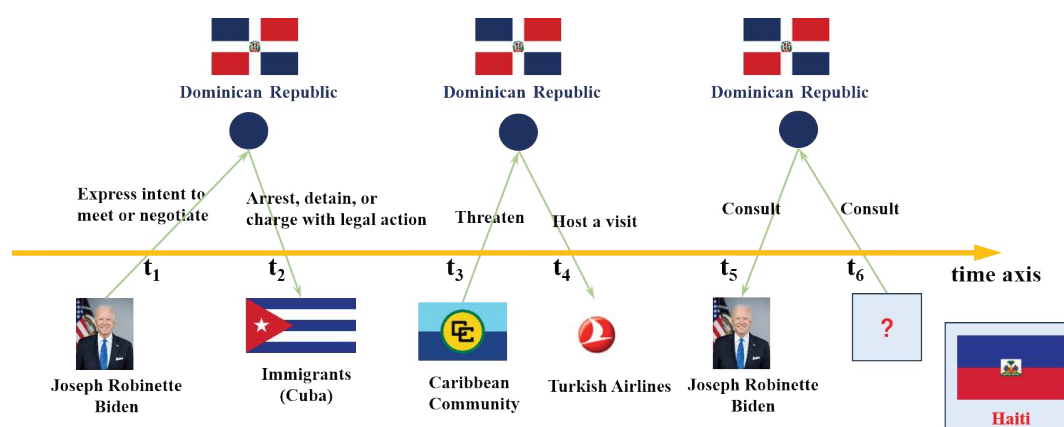


Figure 1: Example of temporal knowledge graph

There are two main categories of approaches to modeling temporal knowledge: timestamp transformation-based models and snapshot-based models. **Timestamp transformation-based** models represent timestamps within text. They incorporate timestamps into embedding vectors of entities or relations, or model time using hyperplane projections. **Snapshot-based models** represent temporal knowledge as a sequence of graphs, and using *time-and-graph* or *time-then-graph* to model graph structures^[1]. Figure 1 is a temporal knowledge graph example, and some of its knowledge is represented in quadruple displayed in Table 1. The relation between entities consistently vary with time. The difficulty in complete temporal knowledge graph is modeling temporal information accurately to help complete historical knowledge or predict future fact.

LLMs have demonstrated remarkable success in natural language processing tasks. This leads to numerous research efforts in knowledge extraction, fusion, representation, and completion using pre-trained language models. However, there are the following challenges in directly applying large

Table 1: Examples of temporal knowledge represented by quadruples

Subject entity	Relation	Object entity	timestamp
Joseph Robinette Biden	Express intent to meet or negotiate	Dominican Republic	t1
Dominican Republic	Arrest, detain, or charge with legal action	Immigrants(Cuba)	t2
Caribbean Community	Threaten	Dominican Republic	t3
Dominican Republic	Host a visit	Turkish Airlines	t4
Joseph Robinette Biden	Consult	Dominican Republic	t5

models to temporal knowledge graphs: (1) Calculation threshold: Training and fine-tuning large models require a significant amount of computational resources; (2) Prompt requirement: The use of large models requires the design of high-quality prompts; (3) Adaptation problem: The graph is generally sparse, and using large models can easily cause overfitting. Recently, researchers have explored the design and application of small pre-trained language models (SLMs), and various small model architectures have emerged^[2]. The performance of small models in certain tasks is not inferior to that of large models, indicating that the task of completing temporal knowledge graphs may be relatively effective through simplified language model architecture.

Driven by this motivation, we explored the application of BERT, the most basic pre-trained language model (PLM) architecture whose parameter number is far less than LLM, in tasks related to temporal knowledge graphs. We propose TKG-BERT (Temporal Knowledge Graph by Bidirectional Encoder Representations from Transformers). Our approach extends the capabilities of pre-trained language models to capture temporal aspects by introducing three methods for modeling temporal knowledge: vanilla knowledge embedding, explicit time modeling and implicit time modeling. These methods leverage the pre-trained representations learned by BERT and aim to enhance the understanding and completion of TKGs. We conduct experiments with various temporal models and explore the capacity of TKG-BERT to model temporal knowledge. The contributions of this paper are summarized as follows:

- Application of pre-trained language models to temporal knowledge graphs. This work explores the utilization of PLMs, specifically BERT, for temporal knowledge graph completion. It extends the application of PLMs beyond static KGs and investigates their effectiveness in capturing and modeling temporal aspects.
- Three ways for temporal knowledge embedding. We introduce three methods: vanilla knowledge embedding,

explicit time modeling and implicit time modeling for embedding temporal knowledge. Main fine-tuning tasks includes masked entity modeling, masked relation modeling, and tuple classification modeling.

- Evaluation of PLM capacity for temporal knowledge modeling. We conducts experiments on two temporal KG datasets to assess the capacity of TKG-BERT in temporal knowledge completion. It evaluates the ability of TKG-BERT to conduct knowledge reasoning under different time modeling ways.

The rest of this paper is organized as follows: Section 2 gives the notations and problem definition. Section 3 presents the model architecture, and three time modeling methods of our proposed approach TKG-BERT. Section 4 describes the experimental setup and presents the results and analysis. Finally, Section 5 gives the discussions and conclusions of this research.

Notations and Problem Definition

We first list the general notations used in model description in Table 2, and define the problem of temporal knowledge graph.

A temporal knowledge graph G can be formalized by quadruple $q = (s, r, o, t)$, which is a triple with timestamp t . The quadruple denote an fact that happen at timestamp t , where $t \in T$. $s \in V$ and $o \in V$ are subject entity and object entity, respectively. $r \in R$ denote a relationship fact between s and o . The notations used in this seciton are listed in Table 2.

Table 2: Notation description

Notation	Description
G	temporal knowledge graph
V	entity set
R	relation set
T	timestamp set
q	quadruple (s, r, o, t)
s	subject entity
r	relation
o	object entity
t	timestamp

Table 3: Abbreviation of models and tasks

Model	Description	Tasks	Description
TKG-BERT (Van.)	The vanilla version of TKG-BERT, without time modeling	<i>sop</i>	subject object prediction, namely entity prediction
TKG-BERT(Exp.)	TKG-BERT with explicit time modeling	<i>rp</i>	relation prediction
TKG-BERT(Imp.)	TKG-BERT with implicit time modeling	<i>tc</i> and <i>qc</i>	tuple classification (<i>tuc</i>), including triplet classification (<i>tc</i>) and quadruple classification(<i>qc</i>)

There are three common tasks of temporal knowledge embedding: entity prediction, relation prediction, and tuple classification. Entity prediction is to predict the missing subject entity s in the incomplete quadruple $(?, r, o, t)$ or the missing object entity o in $(s, r, ?, t)$. Relation prediction is to predict the missing relation r in $(s, ?, o, t)$. Tuple classification is to determine whether the given tuple (s, r, o, t) is correct or incorrect.

For ease of reading, the abbreviations and meanings of the three models proposed in this study are listed in Table 3, as well as the abbreviations and explanations of the knowledge graph completion tasks that were adopted.

Methodology

This section introduce our proposed temporal knowledge graph embedding approach: TKG-BERT. We illustrate model architecture of TKG-BERT in section 3.1. Then, we introduced three temporal knowledge modeling methods based on TKG-BERT in sequence, including vanilla, explicit, and implicit modeling. We provide a detailed introduction on how to design and perform input-output tasks for each of these three different modeling approaches.

Overview of Temporal Knowledge Graph BERT

BERT (Bidirectional Encoder Representations from Transformers) [3] is a pre-trained language model based on multi-layer Transformer encoder [4]. BERT learns deep bidirectional representations from unlabeled text by jointly conditioning on both the left and right context in all layers. The same as other language model, BERT consists of two steps: pre-training and fine-tuning. During pre-training, BERT is trained on large scale unlabeled general domain corpus. In fine-tuning phase, BERT is initialized with pre-trained parameters and then is fine-tuned on specific domain corpus and tasks such as named entity recognition, question answering, and sentence pair classification. To leverage the rich language patterns and contextual representations effectively, we fine-tune the pre-trained BERT model for temporal knowledge completion tasks. We concatenate the entity tokens, relation tokens, and timestamp tokens as word sequences into BERT for fine-tuning. Such architecture is called **TKG-BERT** (Temporal Knowledge Graph based on BERT). TKG-BERT utilize pre-trained BERT (*BERT_base*) and are fine-tuned on sequence classification with temporal knowledge graph corpus.

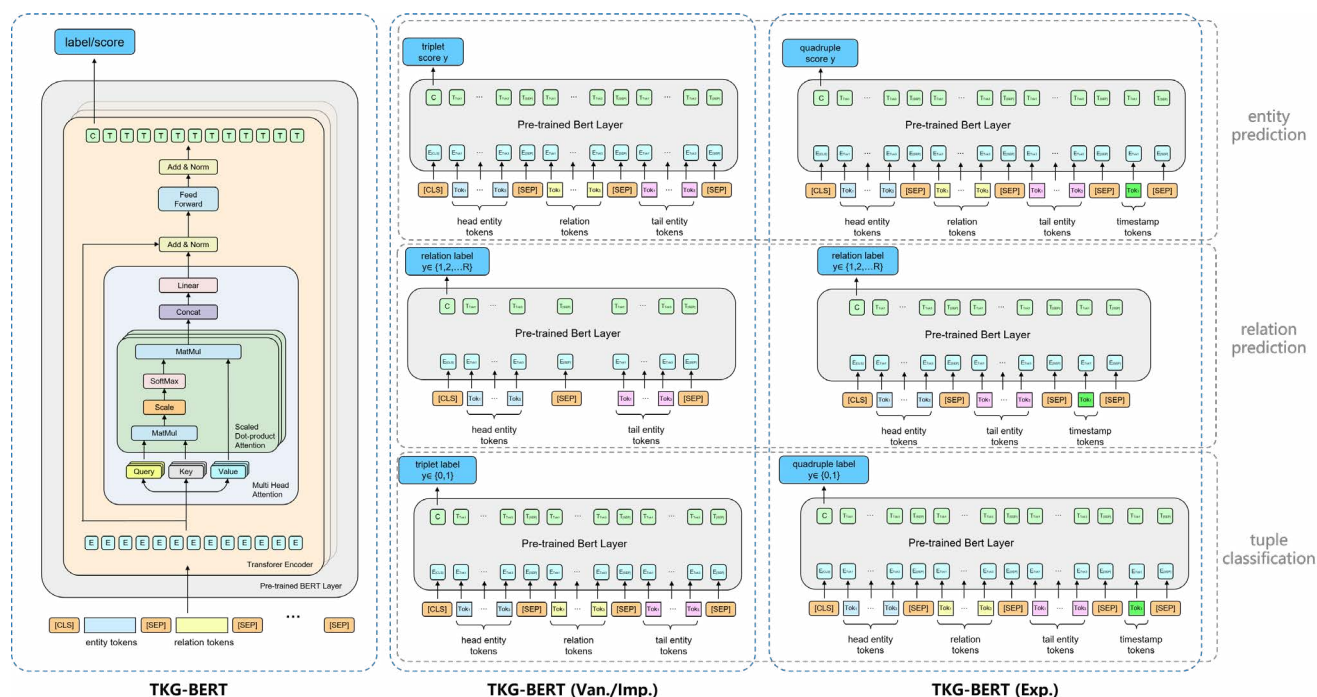


Figure 2: Illustrations of fine-tuning TKG-BERT with different time moedling ways on various tasks.

The left part of Figure 2 shows the architecture of TKG-BERT for modeling knowledge represented by tuple (triple or quadruple). For each tuple, we represent the entities and relation as their text word sequences. TKG-BERT take entity and relation word sequences as the input sentence for fine-tuning. As shown in Figure 2, we concatenate the word sequences of (s, r, o) as a single input sequence, i.e., the input token sequence to BERT. This is the general universal architecture, because the inpput tokens and the output labels maybe different according to different modeling modes. For example, there is an temporal quadruple:

(Islamic Rebirth Party, Make a visit, Tatarstan, 2014-03-21)

KG-BERT takes the following token sequences as an input:

$([CLS], \text{Islamic, Rebirth, Party, [SEP], Make, visit, [SEP], Tatarstan, [SEP], 2014-03-21, [SEP]})$.

In original BERT, “[CLS]” is the special symbol for classification output, and “[SEP]” is the special symbol to separate non-consecutive token sequences. In our TKG-BERT, the first token of each input sequence is always “[CLS]”, denoting the tuple representation is fed into an output layer for classification. The word sequences of entities and relations are separated by “[SEP]”.

Token sequences of knowledge are input into pre-trained

BERT to generate embeddings (blue blocks marked with “E” in Figure 2). The embedding of each given token is generated by summing token embedding, segment embedding, and position embedding. Token embedding is the original word embedding of current token. Segment embedding is the embedding to distinguish the tokens in different segment. The tokens seperated by “[SEP]” have different segment embeddings, whereas tokens within one same entity or relation have the same segment embedding. Position embedding aims to fuse position information, so different tokens at the same position have the same position embedding. The token embeddings are fed into BERT, generate the final hidden vectors (green blocks marked with “C” and “T” in Figure 2) after Transformer encoding. The final hidden vector “C” is used for aggregating sequence representation for computing the final label. Other hidden vectors marked “T” corresponds to entity tokens, relation tokens, and “[SEP]” tokens. Label denotes the final output given input triple, which is different due to different training task and mode.

$$Attention(Q, K, V) = \text{soft max}(\frac{QK^T}{\sqrt{d_k}})V \quad (1)$$

$$MultiHead(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^0 \quad (2)$$

$$T = \text{Norm}(\text{FeedForward}(\text{Mid}) + \text{Mid}) \quad (3)$$

$$\text{Mid} = \text{Norm}(E + \text{MultiHead}(E)) \quad (4)$$

The pre-trained BERT layer consists of 12 bidirectional Transformer encoders. Each bidirectional Transformer encoder implements a multi-head self-attention. The multi-

head attention generate multiple sets of (Q, K, V) according to different weight matrices. Q, K, V refers to query, key, and value in multi-head self-attention. Transformer calculate attention according to equation (1). The output of Multi-head attention is equation (2). The final hidden vector T are calculated by equation (3), wherein Mid is the output of normalized multi-head after residual, calculated by equation (4). Building on the above, we designed three approaches for fine-tuning TKG-BERT on temporal knowledge graph reasoning tasks. This enables us to investigate whether temporal information plays a role when using language models for knowledge graph reasoning, as well as the extent to which different temporal modeling methods affect the reasoning outcomes.

Vanilla knowledge embedding of TKG-BERT

The vanilla knowledge embedding design of TKG-BERT (Abbreviated as TKG-BERT (Van.)) intends to investigate its performance on temporal knowledge graph tasks without incorporating temporal information. The task modeling approach for TKG-BERT (Van.) is illustrated in the middle section of Figure 2. Under this configuration, TKG-BERT does not model the temporal information present in the temporal knowledge graph but instead trains and predicts using only the static triple components of the temporal quadruples: (s, r, o) . The tasks are identical to those performed on static knowledge graphs.

Original BERT randomly masks some tokens of the input sequences and then predict those masked token. Inspired by this masked language modeling, TKG-BERT adopts masking entity or relation in triple to learning their embeddings. As depicted, the three tasks include entity prediction, relation prediction, and triple classification. For the entity prediction task, masked entity modeling is employed, while for the relation prediction task, masked relation modeling is used.

- **Masked entity modeling** is to construct positive and negative tuple samples by randomly corrupt the subject entity s or the object entity o . TKG-BERT will learning the optimal embeddings to make the triple scores of positive and negative samples seperated as far as possible. Then during the test phase, the masked entity would be predicted towards the correct scoring.
- **Masked relation modeling** is to delete the relation in input tuple sequence. Only the subject entity and object entity are input into fine-tuning. The relations are regarded as labels. TKG-BERT learns to embedding the entities towards fitting the relation label representations.

The architecture of TKG-BERT(Van.) for triple classification mode is shown in Figure 2. On this task, TKG-BERT also take the concatenation of word sequences of entities and relation as token sequence input, whereas the output label denotes the quadruple is true or false.

Explicit Time Modeling of TKG-BERT

Explicit Temporal Modeling of TKG-BERT (Abbreviated as TKG-BERT (Exp.)) refers to the process of explicitly incorporating time-related information into models designed for handling data that has a temporal component. This modeling method typically involves the explicit representation and utilization of timestamps or other temporal features in the learning and inference mechanisms of the model. Temporal knowledge graph is usually formally represented as quadruple: (s, r, o, t) , wherein t is the time that the triple fact happens. Explicit time modeling is to treat the timestamp as individual elements as entity and relation, and learn the embedding of the timestamp. Compared to TKG-BERT(Van.) with no temporal modeling, TKG-BERT(Exp.) embeds timestamps alongside entities and relations, appending the timestamp token after the entity-relation triple, thereby inputting the temporal quadruples into the model. Tasks under explicit temporal modeling include entity prediction with timestamps, relation prediction with timestamps, and quadruple classification. The inputs and outputs for these three tasks are illustrated in the right part of Figure 2.

- TKG-BERT(Exp.) for predicting entity takes the concatenation of subject entity, relation, object entity, and timestamp in quadruple as token sequence input. Embedded by the pre-trained BERT layer, the token embeddings are transformed to final hidden vectors. The hidden vector C of the special token “[CLS]” aggregates the sequence representation, then calculate the quadruple score as the model output.
- TKG-BERT(Exp.) for predicting relation only use the tokens of subject entity s , the object entity o , and timestamp t to predict the relation r between them. In preliminary of KG-BERT^[5], predicting with two entities directly is better than using entity prediction mode with relation corruption. So we adopt the same way for predicting relation. After encoding and final hidden vector generating, the model output the relation label $y \in R$ of given entity pair.

TKG-BERT(Exp.) for quadruple classification takes the quadruple as token sequence input. The only difference between quadruple classification and triple classification is the addition of the timestamp. Quadruple classification also adopt binary classification, distinguish positive and negative quadruple samples.

Implicit Time Modeling of TKG-BERT

In previous research on temporal knowledge graphs, the modes of knowledge prediction include interpolation and extrapolation. As shown in the left part of Figure 3, “interpolation” involves randomly selecting a portion of the knowledge for model training and speculating on the missing knowledge. In this mode, the model may infer missing

historical knowledge based on knowledge from future time points. “Interpolation” mode corresponds to vanilla knowledge modeling of TKG-BERT. “Extrapolation”, on the other hand, involves training the model using historical data and then reasoning or predicting knowledge at future time points. In the previous subsection on explicit temporal modeling, we adopted the interpolation setting. However, in practice, the extrapolation mode is more aligned with practical applications. Therefore, we designed a time modeling approach under the extrapolation setting, which is the implicit temporal modeling (Abbreviated as TKG-BERT (Imp.)).

As illustrated in the right part of Figure 3, we restructured the two datasets used in this study according to their temporal order, selecting 80% of the historical data for the training set and the more recent data for the test set. Under this setting, similar to the TKG-BERT(Van.), we do not explicitly embed temporal information such as timestamps. Instead, we implicitly model time through the restructuring of the dataset, learning from history to predict the future.

TKG-BERT(Imp.) captures temporal dynamics within a KG without explicitly encoding or representing time-related information. In this method, the model learns to infer temporal patterns and dependencies from the input data itself, rather than relying on explicit timestamps or time intervals. For the given temporal knowledge graph, we reconstruct the graph, create training set by selecting the fact quadruples that occurred relatively earlier, and conduct entity or relation prediction of the fact quadruples which occur in a relatively future time.

Experiments

We conduct abundant experiments to evaluate the performance of TKG-BERT on public temporal KG datasets. In the following sections, we first introduce experiment

settings, including dataset selection, evaluation tasks and metrics, and hyperparameter settings. Then we present each part of the experiment results, analyze the role of time and temporal information in knowledge reasoning tasks.

Experimental Settings

This section introduces the temporal knowledge graph dataset used in this research, the experimental tasks, and the evaluation metrics, as well as the of hyperparameter settings during model training.

Datasets

We evaluate TKG-BERT on two temporal KG dataset: ICEWS14 and ICEWS05-15. They are constructed from ICEWS (Integrated Crisis Early Warning System)^[6], an periodic updated event graph. ICEWS consists of events that represent interactions between the socio-political actors (for instance, cooperative or hostile actions between individuals, groups, sectors and nation states). ICEWS has been discontinued, changed to POLECAT. POLECAT contains event data from 2018 to the present, updated weekly, and previously events of terminated years stored monthly. Event knowledge are stored in fields according to an event standard. However, as most of the research and experiments on temporal knowledge graphs are based on the ICEWS dataset, we also used this series of datasets for research purposes in order to facilitate performance comparison with baseline models.

ICEWS14 and ICEWS05-15 used in this work are subsets of the data present in the ICEWS repository as created by Garcia-Duran et al.^[7]. The statistics of the 2 datasets are listed in Table 4. ICEWS14 is a short-range dataset consisting of the events that occurred in 2014. ICEWS05-15 is a long-range dataset consisting of the events that occurred between 2005 to 2015. We take the same partition proportion for training/validation/test set as DE^[8], which is a classic temporal KGE model.

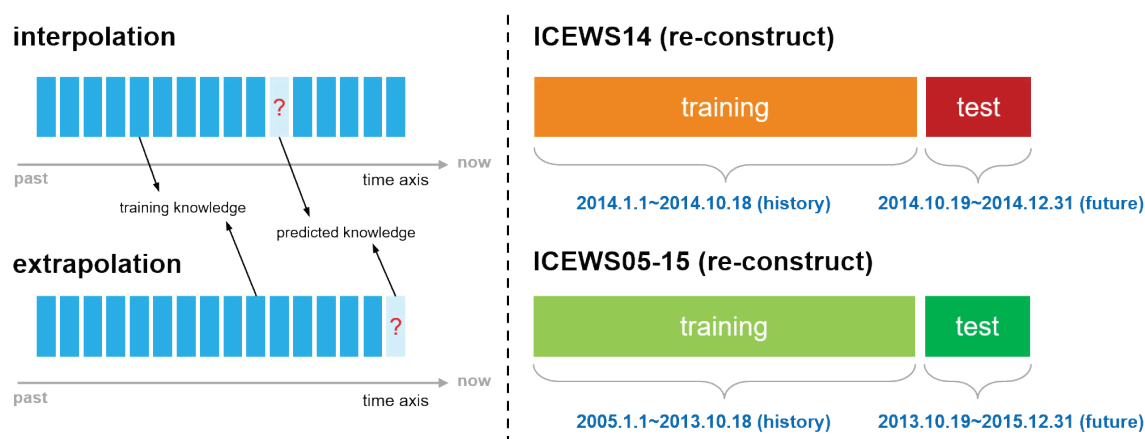


Figure 3: The task setting of interpolation and extrapolation and the reconstruction of datasets.

Tasks and Metrics

Four kinds of tasks are adopted to evaluate different time modeling way of TKG-BERT, including *sop*, *rp*, *tc*, and *qc*. The most principal task for temporal knowledge graph completion (TKGC) is *sop*, namely entity prediction. It aims to forecast the absent entity within a test quadruple (*s*, *r*, *o*, *t*). Specifically, it means predicting the missing *s* when presented with (*?*, *r*, *o*, *t*), or predicting the missing *o* when given (*s*, *r*, *?*, *t*). For each test quadruple, we substitute the unknown entity with every candidate entity from the entity set *V*, thereby creating a collection of corrupted quadruples. We then evaluate these corrupted quadruples using temporal KG embedding models to assign scores. A superior model should position the genuine quadruple ahead of the corrupted ones in its ranking. Consequently, we ascertain the rank of the genuine quadruple among the corrupted set. The mean rank across all test quadruples is denoted as MR, while the mean reciprocal rank is referred to as MRR. Hits@*n* measures the percentage of correctly predicted entities that appear within the top *n* ranks. Lower MR values, higher MRR values, and increased Hits@*n* percentages all indicate enhanced model performance. To conduct a quantitative comparison of different models, we utilize the widely accepted evaluation metrics: Mean Rank (MR), Mean Reciprocal Rank (MRR), and Hits@*n* (where the value of *n* is 1, 3, and 10).

Hyper-parameter Setting

TKG-BERT are implemented with deep learning framework PyTorch and is trained on GPU of NVIDIA GeForce RTX 4090. Unless otherwise specified, we take the default setting for all hyper-parameters, which are listed as follows: learning rate is $5e^{-5}$, training batch size is 32, evaluation batch size is 135, the max token sequence is limited to 15, the training epochs are 5.0, negative ratio is 2, embedding size is 50, output embedding dimension is 100; we use multi-head attention in fine-tuning, and set the number of head to 2, the margin of hinge loss is 5. We use pre-trained BERT word embeddings *bert-base-uncased*.

Temporal Knowledge Graph Completion

This section is the main experimental part of TKGC, which includes three temporal knowledge modeling methods for subject entity and object entity prediction, relation prediction, and tuple classification tasks.

Comparison on Entity Prediction

For entity prediction task, we compare TKG-BERT with

two categories of models: static KGE model and temporal KGE model. We compare TKG-BERT (Van.) with static KGE models to explore their completion capabilities without temporal modeling. Meanwhile, TKG-BERT(Exp.) and TKG-BERT(Imp.) are contrasted with temporal models to investigate the performance gains achieved through time-aware information modeling. Compared static models include DistMult^[9], ComplEx^[10], R-GCN^[11], ConvE^[12], ConvTransE^[13], RotatE^[14]. Compared temporal models under the interpolation setting include HyTE^[15], TTransE^[16], TA-DistMult^[7], and DE^[8]. The interpolation setting is appropriate for the comparison of TKG-BERT(Exp.). Compared temporal models under the extrapolation setting include RG-CRN^[17], CyGNet^[18], RE-NET^[19], RE-GCN^[20]. The extrapolation setting is appropriate for comparison of TKG-BERT(Imp.).

Table 5 present the entity prediction results on the test sets of ICEWS14 and ICEWS05-15. The results of baseline models are from the paper^[21]. As demonstrated in the table, the three time modeling approaches of TKG-BERT consistently outperform the baselines across most metrics, showing significant improvements. Specifically, for static knowledge graph embedding, which there is no time information modeling, TKG-BERT(Van.) far exceed all baselines. On ICEWS14, TKG-BERT(Van.) achieves the improvements of 20.42%(51.92-31.50) on MRR, 14.0% (36.26-22.46) on Hits@1, 27.9% (62.88-34.98) on Hits@3, 28.17% (78.20-50.03) on Hits@10, compared with suboptimal model (ConvTransE). On ICEWS05-15, TKG-BERT(Van.) achieves the improvements of 35.06%(85.34-30.28) on MRR, 37.51% (59.07-21.56) on Hits@1, 32.36% (68.06-35.70) on Hits@3, 21.73% (72.69-50.96) on Hits@10 compared with suboptimal model (ConvTransE). There is a huge improvement in both ICEWS14 and ICEWS05-15, and the improvement in ICEWS05-15 is greater than that in ICEWS14.

For temporal knowledge graph embedding under interpolation setting, corresponding to explicit time modeling, TKG-BERT(Exp.) also show superiority on entity prediction. On ICEWS14, TKG-BERT(Exp.) achieves the improvements of 2.07%(52.17-50.10) on MRR, 0.92% (40.12-39.20) on Hits@1, 1.37% (58.27-56.90) on Hits@3, 0.11% (70.91-70.80) on Hits@10 compared with suboptimal model (DE-DistMult). On ICEWS05-15, TKG-BERT(Exp.)

Table 4: Dataset statistical information.

Dataset	Entity	Relation	Time Span	Time Gap	Training	Validation	Test	Total
ICEWS14	7,128	230	2014	365	72,826	8,941	8,963	90,730
ICEWS05-15	10,488	251	2005-2015	4,017	3,86,962	46,275	46,092	4,79,329

Table 5: Performance for the entity prediction task on ICESW14 and ICEWS05-15 with raw metrics (in percentage)

Method	ICEWS14				ICEWS05-15			
	MRR	Hits@1	Hits@3	Hits@10	MRR	Hits@1	Hits@3	Hits@10
DistMult	20.32	6.13	27.59	46.61	19.91	5.63	27.22	47.33
CompLex	22.61	9.88	28.93	47.57	20.26	6.66	26.43	47.31
R-GCN	28.03	19.42	31.95	44.83	27.13	18.83	30.41	43.16
ConvE	30.3	21.3	34.42	47.89	31.4	21.56	35.7	50.96
ConvTransE	31.5	22.46	34.98	50.03	30.28	20.79	33.8	49.95
RotatE	25.71	16.41	29.01	45.16	19.01	10.42	21.35	36.92
TKG-BERT(Van.)	51.92	36.46	62.88	78.2	65.34	59.07	68.06	72.69
HyTE	16.78	2.13	24.84	43.94	16.05	6.53	20.2	34.72
TTransE	12.86	3.14	15.72	33.65	16.53	5.51	20.77	39.26
TA-DistMult	26.22	16.83	29.72	45.23	27.51	17.57	31.46	47.32
DE-TransE	32.6	12.4	46.7	68.6	31.4	10.8	45.3	68.5
DE-DistMult	50.1	39.2	56.9	70.8	48.4	36.6	54.6	71.8
TKG-BERT(Exp.)	52.17	40.12	58.27	70.91	53.58	38.54	63.37	75.99
RG-CRN	33.31	24.08	36.55	51.54	35.93	26.23	40.02	54.63
CyGNet	34.68	25.35	38.88	53.16	35.46	25.44	40.2	54.47
RE-NET	35.77	25.99	40.1	54.87	36.86	26.24	41.85	57.6
RE-GCN	37.78	27.17	42.5	58.84	38.27	27.43	43.06	59.93
TKG-BERT(Imp.)	49.35	35.64	57.33	75.76	50.57	39.43	54.97	69.25

achieves the improvements of 5.18%(53.58-48.40) on MRR, 1.94% (38.54-36.60) on Hits@1, 8.77% (63.37-54.60) on Hits@3, 4.19% (75.99-71.80) on Hits@10 compared with suboptimal model (DE-DistMult). The effect improvement of TKG-BERT(Exp.) is relatively small compared to TKG-BERT(Van.), and the improvement on ICEWS05-15 is slightly higher than that on ICEWS14.

For temporal knowledge graph embedding under extrapolation setting, corresponding to implicit time modeling. TKG-BERT(Imp.) also show superiority on entity prediction. On ICEWS14, TKG-BERT(Imp.) achieves the improvements of 11.57%(49.35-37.78) on MRR, 8.47% (35.64-27.17) on Hits@1, 14.83% (57.33-42.50) on Hits@3, 16.92% (75.76-58.84) on Hits@10 compared with suboptimal model (RE-GCN). On ICEWS05-15, TKG-BERT(Exp.) achieves the improvements of 12.3%(50.57-38.27) on MRR, 12% (39.43-27.43) on Hits@1, 11.91% (54.97-43.06) on Hits@3, 9.32% (69.25-59.93) on Hits@10 compared with suboptimal model (RE-GCN). The effect improvement of TKG-BERT(Imp.) is relatively significant compared to TKG-BERT(Exp.), whereas the improvement on ICEWS14 and that on ICEWS05-15 is almost equivalent.

The excellent performance of TKG-BERT shows the powerful contextual prediction capability of BERT.

According to the statistics of the datasets in Table 4, ICEWS05-15 has more complicated graph structure than ICEWS14. Correspondingly, the entity prediction results on ICEWS05-15 is better compared with ICEWS14. These results indicates that TKG-BERT take good advantage of BERT because it is expert in dealing with complex graph structure. Furthermore, From the comparative results of the three types of KGE models, it can be seen that without utilizing time information, the entity prediction capability of TKG-BERT far exceeds that of static models. However, under conditions where time modeling is employed, the improvement in TKG-BERT's performance is limited. TKG-BERT with implicit time modeling (TKG-BERT(Imp.)) has a slight advantage over TKG-BERT with explicit time modeling (TKG-BERT(Exp.)) in entity prediction. This is because implicit time modeling and explicit time modeling are equivalent to different tasks, with implicit modeling being more difficult and therefore performing worse.

Mean Rank of Entity and Relation

We have visualized the Mean Rank (MR) values for entities and relations in our entity prediction and relation prediction experiments, and the results are shown in Figure 4. It can be seen from the figure that, TKG-BERT reduces the entity MR of entity prediction task for both datasets to 6 and

7 respectively. This means that the vast majority of entities are trained to have ranked mean values of about 6 and 7 in entity prediction task. Since other research works usually do not focus on this metric, the baseline under this metric is not found for comparison in this study. It is known that the MR of entities in static knowledge graphs is generally reduced to a few hundred to a few thousand after training. Given the entity number of these two temporal datasets in this study (7,128 entities in ICEWS14 and 10,488 entities in ICEWS05-15), reducing the MR of entities to less than 10 proves the strong entity semantic learning and reasoning capability of TKG-BERT.

The lower part of Figure 4 represents the Mean Rank (MR) on ICEWS14 using different proportions of the training set. As the size of the training set increases, TKG-BERT (Van.) consistently reduces the overall MR value, lowering the rank of the correct entities and relations. This demonstrates that, without time modeling, continuously increasing the training set helps improve TKG-BERT's overall predictive performance for both entities and relations.

Explicit Time Modeling v.s. Implicit Time Modeling

We compared the performance of TKG-BERT on different tasks under different time modeling settings, including TKG-BERT(Van.), TKG-BERT(Exp.) and TKG-BERT(Imp.). The results are in Figure 5. Horizontal coordinates represent

the metrics for the tasks, including MR, Hits@3, precision, recall, f1, and accuracy. Vertical coordinates represent the results of TKG-BERT on corresponding metric. The metric values except MR are in percentage. The histograms show following findings:

- On a fixed dataset, whether or not time information is embedded does not significantly affect the performance of reasoning tasks.
- For the *sop* task, the model performance ranking is: TKG-BERT(Van.) > TKG-BERT(Exp.) > TKG-BERT(Imp.). For *rp* and *tc* tasks, there is almost no difference between TKG-BERT(Van.) and TKG-BERT(Exp.), both performing slightly better than TKG-BERT(Imp.).
- The differences in the various time modeling approaches of TKG-BERT are more pronounced on ICEWS05-15.

We re-construct ICEWS14 to model temporal information implicitly. Specifically, the fact quadruples that happens in relatively previous time are used for training, those happens later are used for testing. The proportion of the training, validating, and test set is the same as original dataset. This setting increases the difficulty of the model in performing prediction tasks, so the results show worse. Besides, there are some other reasons for the above phenomena: (1) BERT uses pre-trained word vectors, which do not capture the implicit temporal information in digital text. TKG-

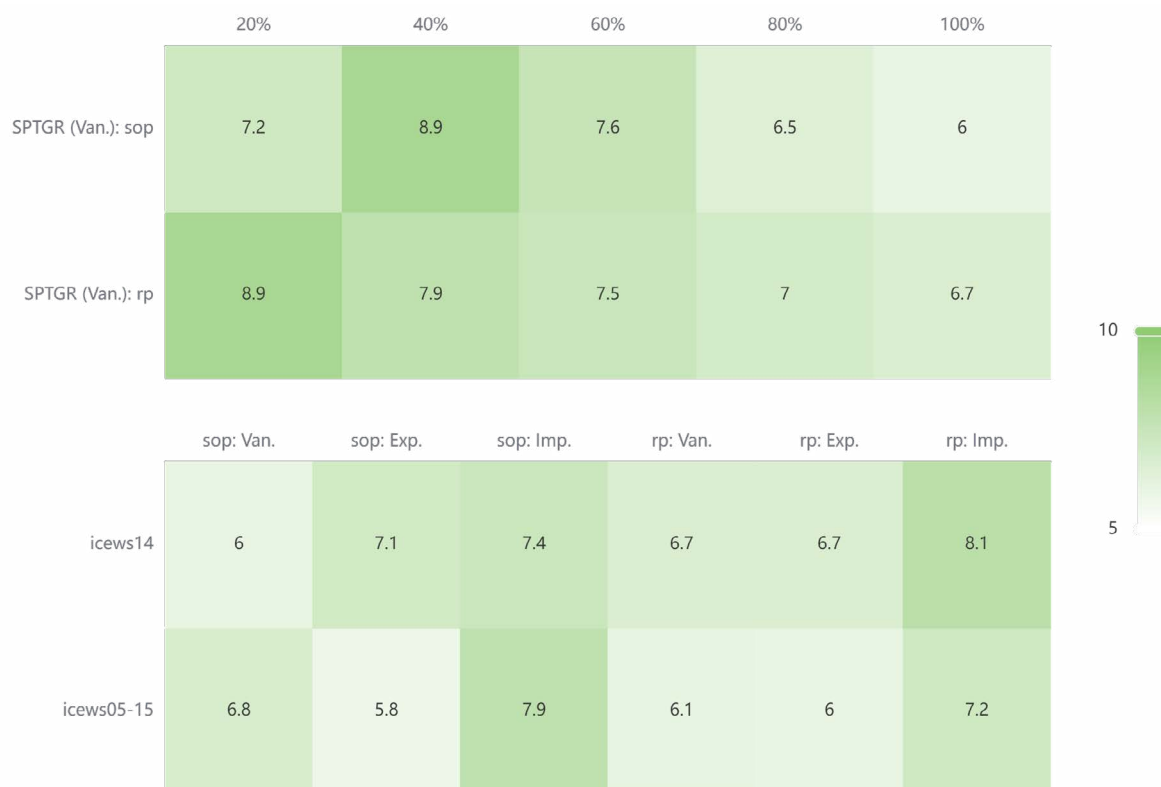


Figure 4: Mean Rank of TKG-BERT on ICEWS datasets.

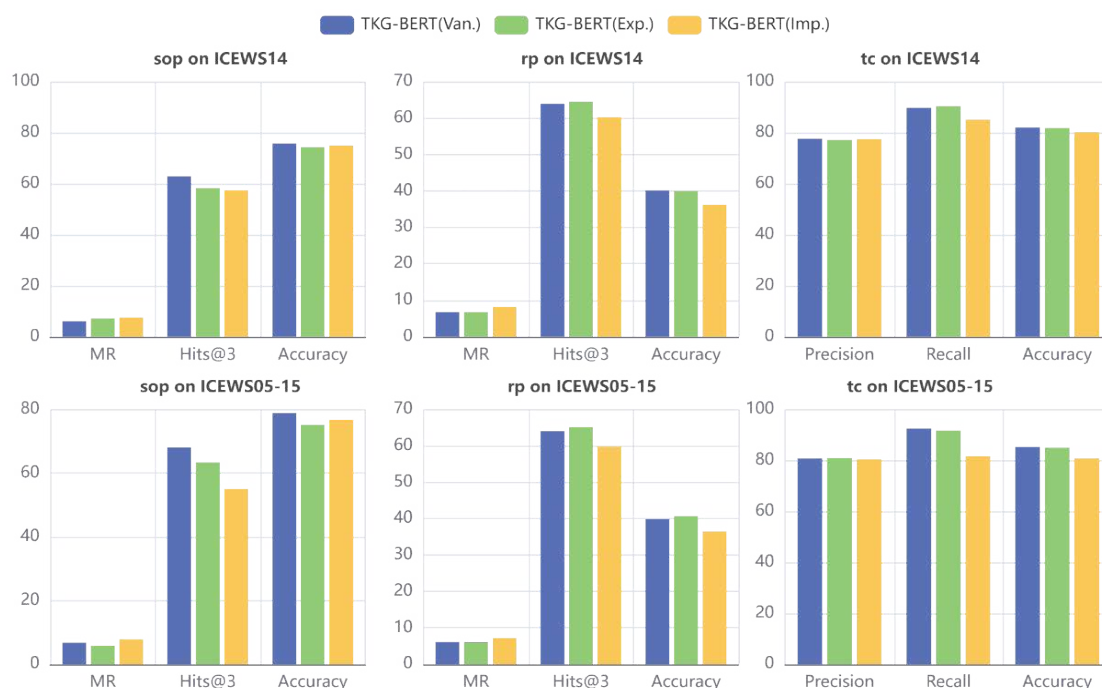


Figure 5: Comparison of TKG-BERT time modeling on different tasks.

BERT(Exp.) embeds timestamp information directly, has little contribution compared with conventional no-temporal modeling approach (2). We use reconstructed dataset for implicit temporal modeling. In original dataset, there is a situation where some event knowledge with later timestamp is introduced into the training set (test leakage), while knowledge with former timestamp is classified into test set, due to data randomization. However, implicit temporal modeling, on the other hand, divides the training set and the test set according to the chronological order, so it does not suffer from such problem. Correspondingly, TKG-BERT (Imp.) performs slightly worse under, which is more realistic. (3)The inclusion of timestamp information had a weakening effect on the model's performance in entity prediction tasks, while it provided a slight improvement in relation prediction tasks. This might be due to the large time span covered by the ICEWS05-15 dataset, making TKG-BERT more sensitive to changes in the textual information of timestamps.

Conclusions

This paper proposes a novel approach called TKG-BERT for temporal knowledge graph embedding. The paper investigates the role of temporal information in knowledge completion tasks. TKG-BERT utilizes pre-trained language models, specifically BERT, to model temporal knowledge in three different ways: vanilla static knowledge modeling, explicit time modeling, and implicit time modeling. The proposed approach is evaluated through various KG completion task to explore the capacity of pre-trained language models to handle temporal knowledge. Experimental results suggest the following conclusions:

- The entity prediction task is the most complex, leading to the least stable model performance. In contrast, the model shows relatively stable performance in the tuple classification task.
- Temporal information is necessary, and explicitly incorporating temporal information can lead to a slight improvement in performance. But there is a need to find more effective ways of modeling timing, in order to effectively utilize temporal information.
- When temporal information is lacking, the model can choose to mine existing knowledge information, thereby improving the inference of unknown knowledge.
- Explicit temporal modeling is suitable for interpolation scenarios where timestamps are available, while implicit temporal modeling is more applicable to stream data and extrapolation scenarios where timestamps are not accessible, making it more aligned with practical needs.

Overall, the TKG-BERT approach presented in this paper fills the gap in the research on temporal knowledge graph completion using pre-trained language models. The experimental results demonstrate the effectiveness and potential of TKG-BERT in temporal knowledge graph representation.

References

1. Gao J, Ribeiro B. On the equivalence between temporal and static equivariant graph representations. In Proceedings of the International Conference on Machine Learning. PMLR (2022): 7052–7076.

2. Nguyen CV, Shen X, Aponte R, et al. A Survey of Small Language Models (2024).
3. Devlin J, Chang MW, Lee K, et al. Pre-training of deep bidirectional transformers for language understanding (2018).
4. Vaswani A, Shazeer N, Parmar N, et al. Attention is all you need. *Advances in neural information processing systems* 30 (2017).
5. Yao L, Mao C, Luo Y. KG-BERT: BERT for knowledge graph completion (2019).
6. Boschee E, Lautenschlager J, O'Brien S, et al. ICEWS coded event data. *Harvard Dataverse* 12 (2015).
7. García-Durán A, Dumančić S, et al. Learning sequence encoders for temporal knowledge graph completion (2018). Goel R, Kazemi SM, Brubaker M, et al. Diachronic embedding for temporal knowledge graph completion. In *674 Proceedings of the Proceedings of the AAAI conference on artificial intelligence* 34 (2020): 3988–3995.
8. Yang B, Yih Wt, He X, et al. Embedding entities and relations for learning and inference in knowledge bases (2014).
9. Trouillon T, Welbl J, Riedel S, et al. Complex embeddings for simple link prediction. In *Proceedings of the International conference on machine learning*. PMLR (2016): 2071–2080.
10. Schlichtkrull M, Kipf TN, Bloem P, et al. Modeling relational data with graph convolutional networks. In *Proceedings of the The Semantic Web: 15th International Conference, ESWC 2018, Heraklion, Crete, Greece, June 3–7, 2018, Proceedings* 15 (2018): 593–607.
11. Dettmers T, Minervini P, Stenetorp P, Riedel S. Convolutional 2d knowledge graph embeddings. In *Proceedings of the Proceedings of the AAAI conference on artificial intelligence* 32 (2018).
12. Shang C, Tang Y, Huang J, et al. End-to-end structure-aware convolutional networks for knowledge base completion. In *Proceedings of the Proceedings of the AAAI conference on artificial intelligence* 33 (2019): 3060–3067.
13. Sun Z, Deng ZH, Nie JY, et al. Knowledge graph embedding by relational rotation in complex space (2019).
14. Dasgupta SS, Ray SN, Talukdar P. Hyte: Hyperplane-based temporally aware knowledge graph embedding. In *Proceedings of the Proceedings of the 2018 conference on empirical methods in natural language processing* (2018): 2001–2011.
15. Jiang T, Liu T, Ge T, et al. Towards time-aware knowledge graph completion. In *Proceedings of the Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers* (2016): 1715–1724.
16. Seo Y, Defferrard M, Vandergheynst P, et al. Structured sequence modeling with graph convolutional recurrent networks. In *Proceedings of the Neural Information Processing: 25th International Conference, ICONIP 2018, Siem Reap, Cambodia, December 13-16, 2018, Proceedings, Part I* 25 (2018): 362–373.
17. Zhu C, Chen M, Fan C, et al. Learning from history: Modeling temporal knowledge graphs with sequential copy-generation networks. In *Proceedings of the Proceedings of the AAAI conference on artificial intelligence* 35 (2021): 4732–4740.
18. Jin W, Qu M, Jin X.; Ren, X. Recurrent Event Network: Autoregressive Structure Inference over Temporal Knowledge Graphs. In *Proceedings of the Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2020, pp. 6669–6683.
19. Jin W, Qu M, Jin X, et al. Recurrent event network: Autoregressive structure inference over temporal knowledge graphs (2019).
20. Li Z, Jin X, Li W, et al. Temporal knowledge graph reasoning based on evolutionary representation learning. In *Proceedings of the Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval* (2021): 408–417.



This article is an open access article distributed under the terms and conditions of the [Creative Commons Attribution \(CC-BY\) license 4.0](https://creativecommons.org/licenses/by/4.0/)