

Incorporating Structured Sentences with Time-enhanced BERT for Fully-inductive Temporal Relation Prediction

Zhongwu Chen
National University of Defense
Technology
Changsha, China
chenzhongwu20@nudt.edu.cn

Chengjin Xu*
International Digital Economy
Academy
Shenzhen, China
xuchengjin@idea.edu.cn

Fenglong Su*
National University of Defense
Technology
Changsha, China
sufenglong18@nudt.edu.cn

Zhen Huang*
National University of Defense
Technology
Changsha, China
huangzhen@nudt.edu.cn

Yong Dou
National University of Defense
Technology
Changsha, China
douyong@nudt.edu.cn

ABSTRACT

Temporal relation prediction in incomplete temporal knowledge graphs (TKGs) is a popular temporal knowledge graph completion (TKGC) problem in both transductive and inductive settings. Traditional embedding-based TKGC models (TKGE) rely on structured connections and can only handle a fixed set of entities, i.e., the transductive setting. In the inductive setting where test TKGs contain emerging entities, the latest methods are based on symbolic rules or pre-trained language models (PLMs). However, they suffer from being inflexible and not time-specific, respectively. In this work, we extend the fully-inductive setting, where entities in the training and test sets are totally disjoint, into TKGs and take a further step towards a more flexible and time-sensitive temporal relation prediction approach SST-BERT – incorporating Structured Sentences with Time-enhanced BERT. Our model can obtain the entity history and implicitly learn rules in the semantic space by encoding structured sentences, solving the problem of inflexibility. We propose to use a *time masking* MLM task to pre-train BERT in a corpus rich in temporal tokens specially generated for TKGs, enhancing the time sensitivity of SST-BERT. To compute the probability of occurrence of a target quadruple, we aggregate all its structured sentences from both temporal and semantic perspectives into a score. Experiments on the transductive datasets and newly generated fully-inductive benchmarks show that SST-BERT successfully improves over state-of-the-art baselines.

CCS CONCEPTS

• **Computing methodologies** → **Temporal reasoning**.

*Corresponding authors

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

SIGIR '23, July 23–27, 2023, Taipei, Taiwan

© 2023 Copyright held by the owner/author(s). Publication rights licensed to ACM.
ACM ISBN 978-1-4503-9408-6/23/07...\$15.00
<https://doi.org/10.1145/3539618.3591700>

ACM Reference Format:

Zhongwu Chen, Chengjin Xu, Fenglong Su, Zhen Huang, and Yong Dou. 2023. Incorporating Structured Sentences with Time-enhanced BERT for Fully-inductive Temporal Relation Prediction. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '23)*, July 23–27, 2023, Taipei, Taiwan. ACM, New York, NY, USA, 11 pages. <https://doi.org/10.1145/3539618.3591700>

1 INTRODUCTION

Temporal knowledge graphs (TKGs), as one of the most popular ways to store knowledge, are structural fact databases in the form of entities and relations between them over time [7]. Many temporal knowledge-intensive tasks, such as information retrieval [17], question answering [19] and entity alignment [44, 45], can benefit from TKGs. Recently, many temporal knowledge graph completion (TKGC) methods [3, 10, 32, 39, 42, 43] focus on temporal relation prediction, a task of predicting missing links through reasoning over the observed facts in TKGs. In the real world, each fact has temporal constraints indicating that the fact holds true within a specific time period or at a certain time point. The typical formats of a fact are $(s, r, o, t, begin, t_end)$ and (s, r, o, t) in different datasets. Traditional temporal knowledge graph embedding (TKGE) representation learning approaches in TKGs embed entities, relations and time into low-dimensional vector spaces and measure the plausibility of quadruples via inputting their embeddings into a score function [11].

However, knowledge of TKGs is ever-changing and the temporal information makes TKGs highly dynamic. In the real world scenario, the emergence of new entities in the development process over time creates the need for temporal relation prediction in the fully-inductive setting, where entities in training TKGs and test TKGs are totally disjoint. Those traditional TKGE methods optimise the representation for a fixed predefined set of entities, i.e., the transductive setting, so they fail in the more realistic fully-inductive setting. Some inductive relation prediction methods over static KGs [33] have been proposed, but they can not learn time information in TKGs. The rule-based method TLogic [16] applies temporal logical rules and thus obtains the fully-inductive ability, but TLogic suffers from the limitation of the inflexibility of symbolic logic. Due to the prior knowledge in PLMs, some PLM-based models

achieve better results. However, PLMs are usually pre-trained in large-scale corpora, so they are not adapted to particular domains such as temporal signals in TKGs [25]. The texts in their corpora also have no explicit temporal indications, which makes it difficult for PLM-based models to handle temporally-scoped facts in TKGs.

To tackle these problems, we propose SST-BERT, a model that incorporates **Structured Sentences** (constructed from relation paths and historical descriptions) with **Time-enhanced BERT** by our designed *time masking* strategy to solve fully-inductive temporal relation prediction task. Table 1 illustrates the detailed differences. Relation paths between target entities and historical descriptions of target entities make SST-BERT consider rich structural information and temporal knowledge while reasoning. The structured sentences enable SST-BERT to learn in semantic space, overcoming the symbolic restrictions in TLogic. Through our proposed *time masking* pre-training task in a corpus consisting of sentences rich in time tokens, SST-BERT is more time-sensitive to facts in TKGs.

First, in order to consider both structural and linguistic knowledge, we convert relation paths connecting the subject and object into natural language form as a part of the structured sentences and then utilize the language comprehension skill of PLMs to encode the relationships among facts, since LAMA [24] has shown that factual knowledge can be recovered well from PLMs, requiring no schema engineering. Relation paths offer structural information and induce the implicit semantics hidden in the structured connections. Therefore, our model SST-BERT can capture temporal development patterns like rules and reason over TKGs in semantic space.

Secondly, to further enrich background knowledge for the subject and object in a target quadruple, we treat one single edge around the subject or object which is not included in the generated paths as background description texts. Previous PLM-based knowledge graph completion methods, KG-BERT [47] and StAR [35], use the definitions or attributes of entities in external knowledge bases as supporting knowledge. But there exists plenty of information redundant and temporally irrelevant to the target quadruple. To overcome this limitation, our basic idea is to make full use of the more relevant and easily accessible history of entities inside TKGs, instead of relying on noisy external resources outside TKGs. Thus, the edges that happened before the target quadruple can serve as directly relevant information and can be converted into ideal historical description texts of entities. These easily accessible descriptions inside TKGs make up for the inadequate relation paths and are more targeted to the target quadruple than external knowledge bases. Moreover, historical descriptions provide emerging entities with fully-inductive descriptions and integrate them into TKGs via PLMs. As shown in Figure 1, we convert all the target quadruples in the training graph with their relation paths connecting the entities and historical descriptions of the entities into structured sentences.

Thirdly, we propose a *time masking* MLM pre-training task to enhance the time sensitivity of PLMs to changes in facts over time. Different from traditional open-domain corpora, we generate one domain-specific corpus rich in time tokens for each training TKG. In our generated corpus, each sentence is associated with $2n$ or n special time tokens ($n \geq 2$), which depends on the selected dataset. Our proposed pre-training module not only forces PLMs to focus on

the representations of special time tokens but also injects domain-specific knowledge into the parameters of PLMs. In experiments, we use the popular pre-trained language model, BERT [5], and refer to the pre-trained one as *TempBERT*. Compared with BERTRL [48], our pre-trained *TempBERT* is more time-sensitive to the input structured sentences.

Finally, we leverage *TempBERT* to encode different parts of the structured sentences to represent entities and relations in the semantic space and then aggregate all the sentences to compute the probability of occurrence of target quadruples. Traditional TKGC datasets are set up for the transductive setting. To evaluate the ability of SST-BERT to deal with emerging entities, in Section 4.1, we introduce a series of newly generated benchmarks for the temporal relation prediction task in the fully-inductive setting for the first time. The main contributions of this paper are summarized as follows:

- New entities continually emerge in temporal knowledge graphs (TKGs) because of the highly dynamic of TKGs. To the best of our knowledge, this is the first attempt to explore the fully-inductive setting for the temporal relation prediction task. We reconstruct four new fully-inductive benchmarks for each selected dataset.
- We identify the shortcuts of current TKGE, rule-based and PLM-based baselines. Our corpora rich in time tokens are specially generated for TKGs. Therefore, the time-oriented MLM pre-training task, *time masking*, in these corpora and the structured sentences served as inputs of SST-BERT solve the problems of baselines.
- We leverage relation paths and historical descriptions inside TKGs to recover prior knowledge stored in BERT and capture rule-like relation patterns in the semantic space, forming the capability of temporal relation prediction in both transductive and fully-inductive settings. The experiments verify the high performance and robustness of our model SST-BERT.

2 RELATED WORK

2.1 Transductive TKGC Models

Most existing TKGC methods are embedding-based and are extended from distance-based KGC models or semantic matching KGC models to a certain extent. Distance-based KGC models intensively use translation-based scoring functions and measure the distance between two entities. A typical example, TransE [1], defines each relation as a translation from the subject to the object. Semantic matching KGC models, such as ComplEx [34] and DistMult [46], calculate the semantic similarity of representations.

Following TransE, TTransE [14] and HyTE [4] encode time in the entity-relation dimensional spaces with time embeddings and temporal hyperplanes. TA-DistMult [8], a temporal extension of DistMult, learns one core tensor for each timestamp based on Tucker decomposition. Specifically, DE-Simple [9] incorporates time into diachronic entity embeddings and can deal with event-based TKG datasets with timestamps like [2014/12/15], such as ICEWS [2]. But it has issues when facing datasets Wikidata [6] and YAGO [28], since the facts in them have a start time and an end time. In contrast,

Method	Transductive Setting	Fully-inductive Setting	Reasoning Evidence			Time Sensitivity	Explainable
			Relation Paths	Historical Descriptions	Prior Knowledge		
TCompLEx(TKGE)	✓	✗	✗	✗	✗	✓	✗
KG-BERT(PLM-based)	✓	✗	✗	✗	✓	✗	✗
BERTRL(PLM-based)	✓	✓	✓	✗	✓	✗	✓
TLogic(rule-based)	✓	✓	✓	✗	✗	✓	✓
SST-BERT	✓	✓	✓	✓	✓	✓	✓

Table 1: Comparison of our model SST-BERT with other algorithms in terms of their capability of temporal relation prediction in the transductive and fully-inductive settings; considering relation paths, historical descriptions and prior knowledge while reasoning (Reasoning Evidence); whether they can capture temporal changes (Time Sensitivity); and their explainability.

TeRo [41] and TeLM [40] model facts involving time intervals like [2003/##/##, 2005/##/##] and can be generalized to ICEWS, becoming state-of-the-art TKGE models in the embedding-based paradigm. In our model, we can address both two situations. Following ComplEx, TCompLEx [13] extends the regularised CP decomposition to order 4 tensor. However, these TKGE models are naturally transductive, since they embed a fixed set of components while training and cannot be generalized to unseen entities after training.

2.2 Fully-inductive TKGC Models

In contrast to the transductive setting, the inductive setting focuses on continually emerging new entities in TKGs. Graph Neural Network (GNN) is proven to be powerful in capturing non-linear architecture in KGs. However, they cannot handle entirely new graphs, i.e., the fully-inductive setting, since there are no overlapping entities. For the fully-inductive setting, GraIL [33] extracts static KG subgraphs independent of any specific entities and trains a GNN as a score function. Unlike GNN-based methods, the current rule-based TKGC approach TLogic [16] induces probabilistic logic rules and applies learned node-independent rules to completely unseen entities.

Another line of fully-inductive relation prediction is to introduce additional descriptions to embed unseen entities [18, 23, 38]. These methods use definitions or attributes of entities, which are imprecise for the target quadruples and highly costly to obtain. In this paper, we utilize relation paths and historical descriptions in TKGs as texts and construct them into structured sentences for target entities and relations. In this way, we can easily recover their precise knowledge inside TKGs in the semantic space for the fully-inductive setting.

2.3 KG-enhanced Pre-trained Language Models

Pre-trained language models (PLMs) on behalf of BERT [5] are trained in open-domain corpora and show effectiveness in capturing general language representations. But they lack domain-specific knowledge [47]. Recently, many works have investigated how to combine knowledge in KGs with PLMs. A popular approach to make PLMs more entity-aware is to introduce entity-aware objectives while pre-training. ERNIE [30], CokeBERT [27] and CoLAKE [29] predict the entity mentions to entities in texts with a cross-entropy loss or max-margin loss. KG-BERT [47] and StAR [35] fine-tune BERT to incorporate information from the factual triples in KGs. SimKGC [36] and C-LMKE [37] leverage contrastive learning in a batch to model the probability that answers match questions. However, none of them addressed the issue of lacking temporal knowledge in PLMs.

3 METHOD

3.1 Framework

Temporal knowledge graphs (TKGs) consist of a set of edges $(s, r, o, t_{begin}, t_{end})$ or (s, r, o, t) with head, tail entities $s, o \in \mathcal{E}$ (the set of entities) and relation $r \in \mathcal{R}$ (the set of relations). Time period $[t_{begin}, t_{end}]$ and timestamp t indicate when (s, r, o) occurs because the fact may develop from time to time. The temporal relation prediction task in an incomplete temporal knowledge graph \mathcal{G} is to score the probability that a missing edge called target quadruple $(t_{begin}$ and t_{end} are treated as a whole) is true.

Our model scores a target quadruple in four steps as shown in Figure 1. First, we extract the relation paths between target entities s and o and historical descriptions of s and o , denoted as $G(s, o)$ and then we convert them into structured sentences (Section 3.2). Secondly, these structured sentences are used as a pre-training corpus to pre-train BERT by the proposed time-oriented MLM task, *time masking* (Section 3.3). The pre-trained BERT is called *TempBERT*. Thirdly, we encode $G(s, o)$ with *TempBERT* into the representations of entities and relations to compute the plausibility of each individual sentence. The explicit times of the target quadruples are encoded by t2v (Section 3.4). Finally, we design an aggregation function to combine all structured sentences related to a target quadruple for both training (loss) and testing (score). The score is the occurrence probability of the target quadruple (Section 3.5, 3.6).

3.2 Structured Sentences

The structured sentences $G(s, o)$ surrounding entities s and o in a temporal knowledge graph \mathcal{G} provide important insights into predicting missing links between s and o . Traditional methods exploited a small scale of subgraphs to summarize relation patterns. For instance, GraIL [33], CoMPLE [20] and INDIGO [15] discussed subgraph-based relation reasoning and offered first-order logic. However, these subgraph-based approaches involve a paradox of ensuring the necessary patterns of relations are included while producing a sufficiently small graph. In practice, the presence of hub nodes often generates an accumulation of thousands of edges, which does not fit into the GPU memory for PLMs. Furthermore, we find that not all entity pairs have a subgraph, especially in the fully-inductive setting, leading to the failure of the subgraph-based methods.

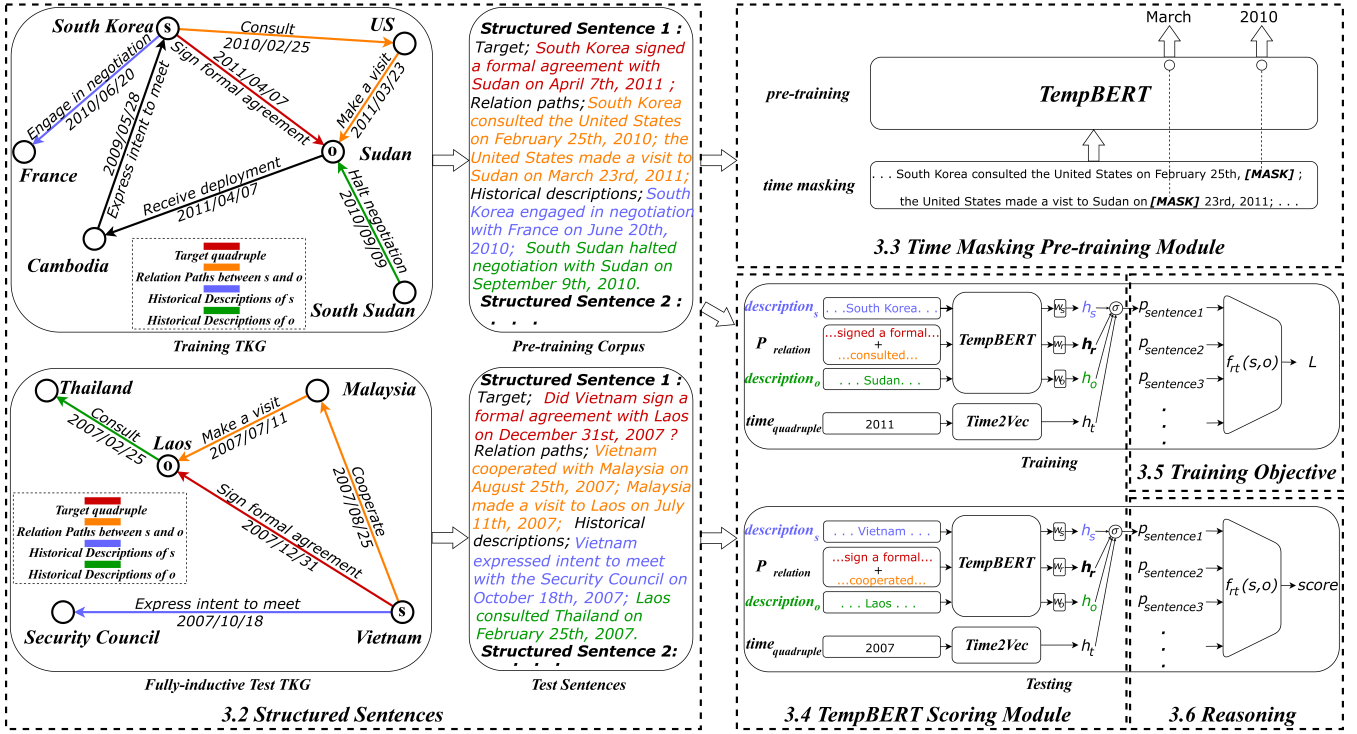


Figure 1: Framework of our model SST-BERT. Each structured sentence consists of four parts: **red** stands for the target quadruple (s, r, o, t) ; **orange** stands for the relation paths between head entity s and tail entity o ; **blue** stands for the historical descriptions of head entity s , $description_s$; **green** stands for the historical descriptions of tail entity o , $description_o$.

Therefore, we adopt a more scalable and universal relation path-based reasoning approach to model the relationship between s and o . Specifically, in Figure 1 we convert symbolic edges into natural language sentences and capture the semantic logic hidden behind the relation paths. In addition, motivated by the prompt-based models [26], we manually define a prompt template for each relation type to express the semantics of quadruples more grammatically, rather than simply splicing their labels in previous works. Compared with symbolic rules explicitly mined by TLogic [16], our constructed structured sentences are more flexible to handle complex associations between facts in semantic space and are not formally restricted.

In addition, temporal reasoning not only needs to consider what happened between target entities (relation paths), but should also pay attention to the history of target entities (historical descriptions). In our PLM-based method, we argue that the old facts in the perspective of time offer vital background information for the target quadruples. Therefore, in Figure 1 we randomly select one edge which is not included in the relation paths as an old fact and regard it as a historical description to reflect the evolution of facts. Usually, these historical descriptions are semantically relevant to target quadruples and act as rich historical knowledge related to target entities. Furthermore, in the fully-inductive setting, the degrees of most entities are small, which means the ideal paths connecting target entities are limited and result in little information. Our proposed incorporation of the history of entities enriches the range

of information available to SST-BERT and makes interpretable inferences. We restrict their timestamps in the structured sentences to precede the timestamp(s) of the target quadruple to avoid information leakage. As Figure 1 shows, the structured sentences in training TKGs serve as both the pre-training corpus for *time masking* and the training inputs in the following sections.

3.3 Time Masking Pre-training Module

There is plenty of temporal information expressed in numerical form (day, year) or terminologies (month) in TKGs, which indicate the happening periods or timestamps of facts. They are represented in the temporal dimension in the embedding-based methods. However, PLMs are pre-trained in noisy and generally purposed open-domain corpora, so PLMs are not effective for applications on time-specific downstream tasks, such as temporal relation prediction. In order to enhance the capability of PLMs on event temporal reasoning over TKGs, we creatively generate a small-scale pre-training corpus based on the training TKG and then pre-train PLMs in the corpus for the specific temporal domain to improve time sensitivity.

Different from the pre-training corpus of the traditional pre-training task, our new corpus is constructed for the temporal domain in TKGs. We regard the whole training quadruple set as a big temporal graph and the relation paths and historical descriptions extracted for each edge correspond to a sentence in the new corpus. Section 3.2 provides the way to produce: we use bidirectional

random walks to find paths and randomly select edges that are not included in the paths as historical descriptions of two target entities. Then we convert them into a natural language form and record the positions of timestamps. We limit the relation paths and historical descriptions of two entities to three hops and one edge, respectively, to restrict the length of a sentence. Because each edge in the training TKGs is considered positive, all the generated sentences for it can be regarded as logically supporting the establishment of the edge to some extent. The timestamp information which is of key importance in our corpus can be easily obtained from all edges, so it becomes a frequent type of element and benefits our pre-training objective.

Different from the masking strategy of the traditional pre-training task, we present a new MLM pre-training task, *time masking*, which replaces the random masking in the original MLM task with a time tokens targeted masking strategy. First, all positions of time tokens in each sentence $\rho = (x_1, x_2, \dots, x_q)$ are recorded in the corpus, so we can easily obtain the positions of temporal expressions denoted by $T = (t_1, t_2, \dots, t_m)$. Secondly, unlike in the case of BERT where 15% of all tokens are randomly sampled, we focus on the time tokens, and 25% of the temporal expressions in T are randomly sampled. We choose 25% as we usually have more than four timestamps in the structured sentences, so we in most cases can mask at least one time token. Thirdly, we continuously randomly sample other tokens in $\rho - T$, until 15% of all the tokens are sampled. Finally, we replace the sampled time tokens in T with [MASK]; 80% of the sampled tokens in $\rho - T$ are replaced with [MASK], 10% with random tokens, and 10% with the original tokens. Note that the sampled time tokens must be masked because recovering them helps PLMs to explore the time order of timestamps in the semantic space and identify the temporal relationships between events.

In practice, we use the popular PLM, BERT [5], which we initialize from $BERT_{BASE}$ (cased) in the publicly available website¹ for simplicity and efficiency. Our pre-trained BERT is called *TempBERT*. Through our proposed *time masking* strategy, *TempBERT* is encouraged to pay attention to the explicit timestamps as well as the temporal changes of facts. The newly generated pre-training corpus rich in time information enables *TempBERT* to comprehend the evolution of events in TKGs from the perspective of time, thus improving its time sensitivity. As a result, our time-enhanced BERT, *TempBERT*, can be well adapted to the temporal relation prediction task.

3.4 TempBERT Scoring Module

In our model, we use *TempBERT* as an encoder to represent entities and relations. As Figure 1 shows, a structured sentence ρ is divided into four parts. The texts converted from target quadruple and relation paths form the rule-like evolution of relation patterns, so they enhance the *TempBERT*'s understanding of the target relation r . We combine token sequences of target quadruple and relation paths as $P_{relation} = ([CLS], p_1^r, p_2^r, \dots, p_\ell^r, [SEP])$ and regard $P_{relation}$ as the enhancement of r . Description texts for head entity s , $description_s = ([CLS], d_1^s, d_2^s, \dots, d_k^s, [SEP])$ and description texts for tail entity o , $description_o = ([CLS], d_1^o, d_2^o, \dots, d_m^o, [SEP])$ are historical descriptions. They can also be considered as the enhancement of entities. In BERT, the output embedding of token '[CLS]'

¹<https://huggingface.co/>

aggregates information of the whole sequence, so we regard the three output embeddings of token '[CLS]' in $description_s$, $P_{relation}$ and $description_o$ as embeddings u_s , u_r and u_o , respectively, to stand for the enhanced encoding by *TempBERT*. We use three independent linear layers to get head entity representation $h_s \in \mathbb{R}^d$, relation representation $h_r \in \mathbb{R}^d$ and tail entity representation $h_o \in \mathbb{R}^d$:

$$\begin{aligned} h_s &= w_s u_s + b_s; \\ h_r &= w_r u_r + b_r; \\ h_o &= w_o u_o + b_o, \end{aligned}$$

where $w_s, w_r, w_o \in \mathbb{R}^{d \times d}$ and $b_s, b_r, b_o \in \mathbb{R}^d$ are learnable weights and biases, d is the dimension of the output of token '[CLS]' in *TempBERT*. Because *TempBERT* has been enhanced by temporal information and each sentence includes rich timestamps, our obtained representations of entities and relations are time-sensitive. In order to learn the time of the target quadruples *time quadruple* t explicitly, we use $t2v$ in Time2Vec [21] to encode t , denoted as $h_t \in \mathbb{R}^d$:

$$h_t[i] = t2v(t)[i] = \begin{cases} \omega_i t + \varphi_i, & \text{if } i = 0 \\ \sin(\omega_i t + \varphi_i), & \text{if } 1 \leq i \leq d - 1 \end{cases}$$

where $h_t[i]$ is the i^{th} element of h_t ; ω_i 's and φ_i 's are learnable parameters. Now, we compute the probability $p_{sentence\rho}$ of the structured sentence ρ leading to the establishment of the target quadruple:

$$p_{sentence\rho} = \sigma \left(w_\theta^T (h_s * h_t + h_r - h_o * h_t) + b_\theta \right),$$

where $w_\theta \in \mathbb{R}^d$, $b_\theta \in \mathbb{R}$ are learnable parameters, σ is the sigmoid function. A target quadruple usually generates a set of sentences, and all of these sentences should be considered together to indicate the truth of the target quadruple while reasoning. Since we take an individual sentence as a knowledge extraction approach, $p_{sentence\rho}$ scores for the individual sentence ρ . In the next section, we aggregate all the probabilities of structured sentences for the target quadruple, making the prediction comprehensively considered.

3.5 Training Objective

While training, the training graph provides positive quadruples. Self-adversarial negative sampling has been proven to be quite effective for TKGC [31], which adopts negative sampling for efficient optimization. We randomly sample n negative quadruples for each positive quadruple by corrupting its head or tail entity and ensuring they do not exist in the training TKGs. These sampled negative entities are restricted within common 3-hop neighbours of the entities in target quadruples. Our self-adversarial loss is as follows:

$$L = -\log \sigma(\gamma - f_{rt}(s, o)) - \sum_{i=1}^n \frac{1}{n} \log \sigma(f_{rt}(s'_i, o'_i) - \gamma),$$

where γ is the margin, σ is the sigmoid function, n is the number of negative samples, $(s'_i, r, o'_i, t_{begin}, t_{end})$ or (s'_i, r, o'_i, t) is the i^{th} negative sample. Note that we generate sentences for both positive and negative quadruples, and $f_{rt}(s, o)$ is the aggregation function for generated structured sentences $\mathcal{D}(s, o)$ to score each quadruple. We further constrain the number of sentences produced by each quadruple, which is a hyperparameter N , to avoid redundant information,

		ICEWS14					ICEWS05-15					YAGO11k					Wikidata12k					
		Entities	Relations	Time	Tokens	Links	Entities	Relations	Time	Tokens	Links	Entities	Relations	Time	Tokens	Links	Entities	Relations	Time	Tokens	Links	
transductive	trans-train	6,869	230		72,826		10,094	251		460,876		10,623	10		16,406		12,554	24		32,497		
	trans-test	2,599	161	365	8,963		4,877	207	4,017	45,858		2,720	10	189	2,051	4,297	20	232	4,062			
fully inductive	v1	ind-train	1,199	133		6,170		1,926	173		23,830		6,772	10		10,067		6,467	24		17,772	
		ind-test	163	35	365	916		814	65	4,017	9,668		712	10	189	1,298	1,093	19	232	2,657		
	v2	ind-train	1,980	166		14,694		3,064	231		57,415		3,964	10		4,587		3,399	23		8,369	
		ind-test	353	57	365	2,657		1,858	108	4,017	20,432		1,504	10	189	1,860	1,509	22	232	3,812		
	v3	ind-train	3,891	213		23,836		6,266	251		191,189		2,455	10		3,890		2,345	23		5,258	
		ind-test	917	102	365	4,884		2,549	196	4,017	38,477		2,312	10	189	2,225	2,400	23	232	4,240		
	v4	ind-train	1,262	147		2,350		2,822	229		46,558		1,738	10		1,914		1,562	24		3,542	
		ind-test	2,492	135	365	7,260		4,107	207	4,017	84,855		3,280	10	189	3,429	3,735	24	232	6,155		

Table 2: Statistics of four transductive datasets, ICEWS14, ICEWS05-15, YAGO11k, Wikidata12k, and four created fully-inductive benchmarks for each of them. ‘trans-’ represents the transductive setting and ‘ind-’ represents the fully-inductive setting.

i.e., $|\mathcal{D}(s, o)| \leq N$. For datasets in the form of $(s, r, o, t_{begin}, t_{end})$:

$$f_{rt}(s, o) = \sum_{\rho \in \mathcal{D}(s, o)} \left(\frac{\exp(t_{\rho} - t_{begin})}{2 \sum_{\hat{\rho} \in \mathcal{D}(s, o)} \exp(t_{\hat{\rho}} - t_{begin})} + \frac{\exp(t_{\rho} - t_{end})}{2 \sum_{\hat{\rho} \in \mathcal{D}(s, o)} \exp(t_{\hat{\rho}} - t_{end})} \right)^{p_{sentence\rho}}$$

where t_{ρ} denotes the earliest timestamp in the sentence ρ and $t_{\rho}(t_{\hat{\rho}}) \leq t_{begin} \leq t_{end}$. For datasets in the form of (s, r, o, t) :

$$f_{rt}(s, o) = \sum_{\rho \in \mathcal{D}(s, o)} \left(\frac{\exp(t_{\rho} - t)}{\sum_{\hat{\rho} \in \mathcal{D}(s, o)} \exp(t_{\hat{\rho}} - t)} \right)^{p_{sentence\rho}}$$

where t_{ρ} denotes the earliest timestamp in the sentence ρ and $t_{\rho}(t_{\hat{\rho}}) \leq t$. The exponential weighting favours sentences with timestamps that are closer to the timestamp(s) of the target quadruple, since they are more likely to be directly relevant to the prediction.

3.6 Reasoning

For a prediction question $(s, r, ?, t_{begin}, t_{end})$ or $(s, r, ?, t)$, we replace the asked tail entity with a set of candidate entities sampled from the whole entity set. Generally, the number of candidate entities can be varied from only a few to all, and it depends on the practical needs. $f_{rt}(s, o)$ is leveraged to score each candidate quadruple through its structured sentences. Then, the entity ranking first is chosen as the answer, and the sentence corresponding to the maximal summation factor in $f_{rt}(s, o)$ provides the most faithful explanation.

4 EXPERIMENT

To demonstrate the temporal relation prediction capability of our model SST-BERT, we compare SST-BERT with some state-of-the-art baselines in the transductive and fully-inductive settings, i.e., Transductive Temporal Relation Prediction and Fully-inductive Temporal Relation Prediction. Moreover, traditional TKG datasets are constructed for the transductive setting, so we are supposed to offer variants derived from them for the fully-inductive setting. By structured sentences, SST-BERT can capture entities’ full-inductive relation paths and historical descriptions inside TKGs. We want to explore the performance of SST-BERT compared with existing PLM-based models which rely on low-quality entity definitions or attributes linking to knowledge bases outside TKGs to offer fully-inductive ability. Another noteworthy point of SST-BERT is that

we specially pre-train BERT on the generated TKG-aimed corpus and use *time masking* strategy to enhance the time sensitivity of the pre-trained *TempBERT*. In contrast, existing PLM-based models are not sensitive to the essential temporal information in TKGs.

4.1 Transductive Datasets and New Fully-inductive Benchmarks

Popular TKG datasets include ICEWS14, ICEWS05-15, YAGO11k, Wikidata12k [4], YAGO15k and Wikidata11k [8]. Time intervals in YAGO15k and Wikidata11k only contain either start dates or end dates, shaped like *occurSince 2000* or *occurUntill 2002*. However, since our model is good at capturing rich and complex temporal information, we prefer to use more challenging YAGO11k and Wikidata12k, where most of the time intervals contain both start dates and end dates, shaped like $[2000/##/##, 2002/##/##]$ or $[2000/01/01, 2002/01/01]$. Similar to the setting in TeRo [41], we only deal with year-level information in YAGO11k and Wikidata12k and treat year timestamps as 189 and 232 special time tokens in BERT. This setting not only balances the numbers of quadruples in different time tokens but also makes the pre-training more targeted. ICEWS14 and ICEWS05-15 are subsets of the event-based database, ICEWS [2], with political events in 2014 and 2005 ~ 2015. These two datasets are filtered by selecting the most frequently occurring entities in graphs. Their time annotations are all day-level, shaped like $[2004/12/24]$. The numbers of special time tokens in them are 365 and 4017.

Our selected original datasets ICEWS14, ICEWS05-15, YAGO11k and Wikidata12k are only suitable for the transductive setting since the entities in the standard test splits are a subset of the entities in the training splits. We create new fully-inductive benchmarks for each dataset by sampling two disjoint subgraphs from its TKG as *ind-train* and *ind-test*. *ind-train* and *ind-test* have a disjoint set of entities, and relations of *ind-test* are entirely included in those of *ind-train*. Specifically, we uniformly sampled several entities to serve as roots and take several random walks to expand them as seen entities in *ind-train*. Then all the quadruples including the selected seen entities create *ind-train*. Finally, we remove these quadruples from the whole TKG and sample *ind-test* using the same way as *ind-train*. Similar to the evaluation setting of GraLL [33], four fully-inductive benchmarks of each dataset are generated with increasing test sizes and different ratios of seen entities in *ind-train* and unseen entities in *ind-test* by adjusting the length of random walks for robust evaluation. In the fully-inductive setting, a model is trained on *ind-train* and tested on *ind-test*. The statistics of four transductive datasets and their sixteen fully-inductive benchmarks are listed in Table 2.

	ICEWS14			ICEWS05-15			YAGO11k			Wikidata12k		
	MRR	Hits@1	Hits@3	MRR	Hits@1	Hits@3	MRR	Hits@1	Hits@3	MRR	Hits@1	Hits@3
HyTE [4]	0.297	10.8	41.6	0.316	11.6	44.5	0.136	3.3	–	0.253	14.7	–
TA-DistMult [8]	0.477	36.3	–	0.474	34.6	–	0.155	9.8	–	0.230	13.0	–
TCompLex [13]	0.610	53.0	66.0	0.660	59.0	71.0	0.185	12.7	18.3	0.331	23.3	35.7
TeRo [41]	0.562	46.8	62.1	0.586	46.9	66.8	0.187	12.1	19.7	0.299	19.8	32.9
TeLM [40]	0.625	54.5	67.3	0.678	59.9	72.8	0.191	12.9	19.4	0.332	23.1	36.0
TLogic [16]	0.430	33.6	48.3	0.470	36.2	53.1	–	–	–	–	–	–
KG-BERT [47]	0.523	50.9	65.2	0.539	52.0	63.2	0.462	43.5	49.6	0.552	52.9	59.6
StAR(Self-Adp) [35]	0.565	53.2	66.8	0.564	53.1	65.6	0.496	48.9	52.4	0.593	54.6	53.8
C-LMKE [37]	0.576	47.0	65.8	0.659	63.2	68.8	0.497	46.3	50.1	0.648	60.0	66.8
BERTRL [48]	0.635	60.9	68.9	0.648	62.4	65.3	0.513	47.1	57.9	0.573	55.8	58.6
SimKGC [36]	0.396	31.4	42.9	0.605	50.7	66.3	0.196	10.7	20.5	0.389	29.3	42.4
SST-BERT	0.688	62.4	72.1	0.693	66.3	76.9	0.558	50.5	58.2	0.684	65.9	69.5

Table 3: Temporal relation prediction results (% for Hits@ k) on four datasets in the transductive setting. Dashes: results are not reported in the corresponding works. Other results are obtained from our experiments. Bold numbers denote the best results.

4.2 Baselines

We compare SST-BERT with typical TKGE models, HyTE [4], TA-DistMult [8], TCompLex [13], TeRo [41] and TeLM [40], for these TKGE models are the state-of-the-art methods in the transductive setting on our four selected datasets. In addition, these TKGE models embed time information into the temporal dimension from different views, becoming a traditional and popular paradigm. However, TKGE models are trained and tested within a fixed TKG with no new entities, so they can not be used directly in the fully-inductive setting. SST-BERT is also compared with a state-of-the-art symbolic framework for TKGC, TLogic [16], which is based on temporal logical rules extracted via temporal random walks. Different from TKGE models, the learned rules in TLogic are entity-independent, making TLogic a fully-inductive method. The results of TLogic on YAGO11k and Wikidata12k can not be obtained, since TLogic focuses on facts with one timestamp rather than facts in YAGO11k and Wikidata12k with start dates and end dates. In SST-BERT, the pre-trained *TempBERT* encodes the structured sentences consisting of relation paths and historical descriptions in semantic space to implicitly form rules, instead of in restricted symbolic space.

Since our model SST-BERT brings BERT and TKG together, we compare it with some current PLM-based models. These approaches leverage BERT to encode entities using description texts, so they are naturally inductive. Some PLM-based TKGC models, such as KG-BRET [47], StAR [35] and BERTRL [48], input external supplementary descriptions or path information into PLMs. Others, such as SimKGC [36] and C-LMKE [37], combine PLMs with contrastive learning. We thus choose them as baselines. Since PLMs in these baselines are designed to encode all kinds of information, different from their original static datasets, we adapt them to our temporal datasets. The external texts of entities required by KG-BRET, StAR and SimKGC are retrieved from Wikipedia², because entity definitions or attributes outside TKGs are a core part of the three models.

²<https://en.wikipedia.org/>

4.3 Implementation Details

For each dataset, we generate a small-scale but domain-specific corpus following Section 3.2 and pre-train *BERT_{BASE}* (cased) with 12 layers and 110M parameters on it for 10 epochs using our proposed *time masking* strategy to enhance the time sensitivity of BERT. The maximum sequence length of the input tokens is 512. The batch size is 8. We use AdamW [12] as the optimizer and set the learning rate to be $5e^{-4}$, with gradient accumulation equal to 32. The above pre-training phase was carried out on 32G Tesla V100 GPUs.

After pre-training, the obtained *TempBERT* is utilized as an encoder to embed entities and relations for scoring. In the training phase, the batch size is chosen from {16, 32, 64} and the learning rate is selected from $\{2e^{-4}, 1e^{-4}, 5e^{-5}, 1e^{-5}\}$. The number of negative quadruples for training n is tuned in {5, 10, 20, 50}, and the maximum number of sentences produced by each quadruple N is in a range of {10, 20, 50, 100}. The output dimension d of *TempBERT* and $t2v$ is 768.

4.4 Experimental Setup

PLM-based methods use BERT as an encoder throughout the process of running. In order to speed up the evaluation process, we generate 50 negative samples for each quadruple in the validation set in advance. These negative samples are generated by corrupting head or tail entities in each quadruple and then selecting other entities in 3-hop neighbourhoods. While testing, we select all entities as candidates. The final results are obtained by ranking the ground truth among the negative quadruples by decreasing scores. We compute the mean reciprocal rank (MRR) and Hits@ k for $k \in \{1, 3\}$.

4.5 Transductive Temporal Relation Prediction

In this section, we explore the performance of our model SST-BERT in the transductive setting in Table 3. TKGE models embed time into vector space and are originally designed for the transductive setting. They form powerful baselines in the temporal relation prediction task. Compared with TKGE models, SST-BERT considers time and easily accessible knowledge inside TKGs in the semantic space and

	ICEWS14				ICEWS05-15				YAGO11k				Wikidata12k			
	v1	v2	v3	v4	v1	v2	v3	v4	v1	v2	v3	v4	v1	v2	v3	v4
TLogic [16]	0.153	0.173	0.228	0.284	0.268	0.283	0.346	0.392	–	–	–	–	–	–	–	–
KG-BERT [47]	0.461	0.437	0.454	0.422	0.565	0.536	0.594	0.538	0.465	0.495	0.435	0.423	0.399	0.485	0.468	0.342
StAR(Self-Adp) [35]	0.498	0.453	0.481	0.464	0.592	0.562	0.584	0.524	0.482	0.462	0.521	0.487	0.389	0.431	0.494	0.405
C-LMKE [37]	0.523	0.542	0.582	0.535	0.724	0.777	0.751	0.706	0.535	0.552	0.527	0.518	0.517	0.596	0.624	0.546
BERTRL [48]	0.342	0.292	0.256	0.237	0.685	0.640	0.623	0.652	0.530	0.556	0.549	0.512	0.532	0.545	0.569	0.467
SimKGC [36]	0.064	0.028	0.030	0.039	0.073	0.033	0.046	0.021	0.070	0.068	0.064	0.059	0.052	0.081	0.074	0.083
SST – BERT	0.569	0.588	0.622	0.613	0.732	0.798	0.765	0.739	0.547	0.599	0.562	0.584	0.551	0.619	0.644	0.571

Table 4: Temporal relation prediction MRR results on the fully-inductive benchmarks. Bold numbers denote the best results.

	ICEWS14				ICEWS05-15				YAGO11k				Wikidata12k			
	v1	v2	v3	v4	v1	v2	v3	v4	v1	v2	v3	v4	v1	v2	v3	v4
TLogic [16]	13.2	15.3	19.5	26.6	23.5	27.5	30.9	35.5	–	–	–	–	–	–	–	–
KG-BERT [47]	45.9	42.5	43.7	40.1	52.9	52.6	55.7	50.4	42.8	46.2	40.1	38.9	32.8	41.2	42.5	30.1
StAR(Self-Adp) [35]	45.2	41.6	45.3	42.9	52.4	51.6	54.9	50.9	45.2	49.3	50.4	43.7	34.6	35.4	35.6	38.6
C-LMKE [37]	46.9	53.8	56.7	46.2	62.3	63.7	62.2	58.5	49.3	51.1	47.5	44.9	47.0	53.9	51.7	54.2
BERTRL [48]	22.9	20.2	24.5	15.4	61.2	53.4	60.3	62.9	50.2	57.9	52.0	49.6	47.5	51.1	43.2	42.6
SimKGC [36]	2.2	0.6	0.4	0.5	4.1	1.4	2.7	0.4	2.2	2.4	1.8	1.9	2.0	3.1	2.8	3.2
SST-BERT	54.4	56.9	60.8	59.6	67.1	66.9	65.9	63.2	53.2	59.1	54.9	52.6	51.1	57.2	53.4	56.9

Table 5: Temporal relation prediction Hits@1 results on the fully-inductive benchmarks. Bold numbers denote the best results.

outperforms them by a stable margin. TLogic [16] is a current rule-based model which utilizes entity-independent rules. Our generated relation paths and historical descriptions deal with rule-like reasoning more flexibly and the semantic understanding capability makes the results of SST-BERT higher than TLogic. Compared with PLM-based baselines which directly use BERT as an encoder without adaptation, SST-BERT enhances itself with time sensitivity by pre-training on the corpus generated from the target dataset TKGs and achieves considerable performance gains over PLM-based baselines.

4.6 Fully-inductive Temporal Relation Prediction

The fully-inductive setting requires models not only to understand the logical patterns of relations but also to be capable of generalizing to unseen entities. Our proposed framework captures the intrinsic evolutionary patterns of relations in the semantic space via the pre-trained *TempBERT* and encodes the emerging entities over time through generated structured sentences. Therefore, our model can be naturally applied to the fully-inductive setting. In reality, since we can not control the number or the ratio of newly emerging entities, we conduct experiments on the four generated benchmarks of each dataset with increasing test sizes and different ratios of seen entities and unseen entities to verify the robustness of our model.

From Table 4 and Table 5, we can find our model SST-BERT outperforms all baselines in the fully-inductive setting. In TKGE models, new entities can not be sufficiently trained, so we exclude them. TLogic [16] relies on matching the data in the test set with the found rules, so the larger the amount of data in the test set, the more likely TLogic is to give more correct answers. However, the emerging entities are usually much fewer than the total entities in the transductive setting and the edges linking to the emerging entities are also limited. Therefore, compared with the transductive setting, the performance of TLogic decreases in the fully-inductive setting. Previous PLM-based models, KG-BERT [47], StAR [35] and C-LMKE [37],

focus on triplet-level information and heavily rely on co-reference resolution and entity linking. They require external data (texts from Wikipedia), causing the accumulation of errors [22]. Although BERTRL [48] considers paths, the important historical descriptions for temporal relation prediction are still ignored. In-batch Negatives (IB) module is a key part of SimKGC_{IB+PB+SN} [36], but this module is disabled when handling emerging entities. This means not all PLM-based models are suitable for the fully-inductive setting.

As the experiments show, our model SST-BERT has a strong fully-inductive ability to reason flexibly over complicated TKGs based on implicitly learned semantic rules. We combine the structured sentences and prior knowledge in PLMs, and the overall framework is similar in the transductive and fully-inductive settings, so there is no bias in the temporal relation prediction capability of SST-BERT in the two settings. In addition, SST-BERT has the most robust results among all baselines except SimKGC_{IB+PB+SN} due to its low performance. For example, the variances of the Hit@1 results of SST-BERT are 6.14, 2.42, 6.45 and 6.43 for the four datasets. In contrast, all the variances of BERTRL are greater than 10. Therefore, SST-BERT can be robustly adapted to the fully-inductive setting in the real world.

4.7 Ablation study

4.7.1 Effect of relation paths, historical descriptions and $t2v$.

As mentioned in Section 3.2, our generated structured sentences consist of relation paths and historical descriptions, these converted texts in the natural language form provide both structural and semantic knowledge. Moreover, relation paths serve as the source of fully-inductive ability and historical descriptions additionally offer background information. To emphasize the significance of the two key components of our model SST-BERT, we conduct an ablation study to remove relation paths (w/o paths) and remove historical descriptions (w/o history). By removing $t2v$ (w/o $t2v$), we also illustrate to what extent the explicitly encoded time contributes to SST-BERT. The results 6 show that removing any part will reduce the performance. Two kinds of descriptions are vital for SST-BERT

fully-inductive benchmarks	SST-BERT	w/o paths	w/o history	w/o t2v
ICEWS14 (v2)	0.588	0.536	0.491	0.552
ICEWS05-15 (v2)	0.798	0.724	0.680	0.743

Table 6: Ablation study of our model SST-BERT (MRR).

to recall the time-oriented knowledge stored in the PLMs, proving that temporal descriptions of entities are remarkable for the time-sensitive TKGc task. The role of t2v is to make SST-BERT precisely recognize the time of the target quadruples without being distracted by other factors in the structured sentences.

4.7.2 Effect of time masking pre-training task.

In this section, we explore the performance improvements from the pre-trained TempBERT with the growth of pre-training epochs. In Figure 2, we replace the time masking pre-training task with the original random masking pre-training task in BERT [5] to illustrate our contribution. First, both pre-training tasks improve the performance of SST-BERT compared with no pre-training. Furthermore, our proposed time masking outperforms the best results gotten by original random masking after 2 ~ 5 epochs due to the time-type targeted masking. Secondly, SST-BERT observes obvious growth in the early pre-training epochs, and the final performance gains for the two datasets are about 21.9% and 14.0% on average. Therefore, our proposed time masking is efficient and effective. In our experiments, we choose to pre-train 10 epochs, since we find the performance tends to be relatively stable after pre-training 10 epochs and the cost of pre-training and the performance achieved are well balanced.

4.8 Explainability and Case Study

Temporal logical rules are usually considered to be explainable. TLogic [16] utilizes temporal random walks to get temporal logical rules in symbolic form. When TLogic applies the found rules to answer questions, it must precisely match each relation and time-tamp in the rules, including relation types, relation order and time order. If there are no matching body groundings in the graph, then no answers will be predicted for the given questions. In contrast, SST-BERT leverages generated relation paths and historical descriptions to form rules in the semantic space, achieving explainability. While pre-training, our proposed time masking MLM task increases the ability of BERT to understand the meaning of timestamps. Furthermore, through training, we implicitly inject semantic rules into the parameters, which are more flexible than symbolic rules. As long as the semantics of the two edges are similar, SST-BERT can detect the relevance of rules rather than abandon them. Finally, for a prediction task $(s, r, ?, t_{begin}, t_{end})$ or $(s, r, ?, t)$, we aggregate all the sentences found and their probabilities for different candidate tail entities. The tail entity o with the highest probability is chosen as the answer and the relations paths and historical descriptions are regarded as the explanation of the occurrence of the target quadruple.

For a prediction question $(Iraq, Engage in diplomatic cooperation, ?, 2014/12/29)$, Figure 3 illustrates the reasoning evidence for the ground truth answer *Iran*. We can see that there are two countries, *China* and *Japan*, and one person, *Barack Obama*, engaged in the happening of the target quadruple; there are two countries, *France* and *Oman*, offering historical descriptions for the target entities.

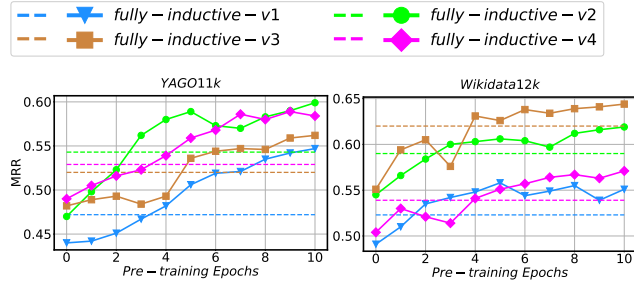


Figure 2: The dashed line represents the best MRR results of SST-BERT within 10 pre-training epochs with the original random masking pre-training task in BERT[5]. The solid line represents the MRR results of SST-BERT with our proposed time masking pre-training task from 0 to 10 epochs. The same colour represents two results are under the same benchmark.

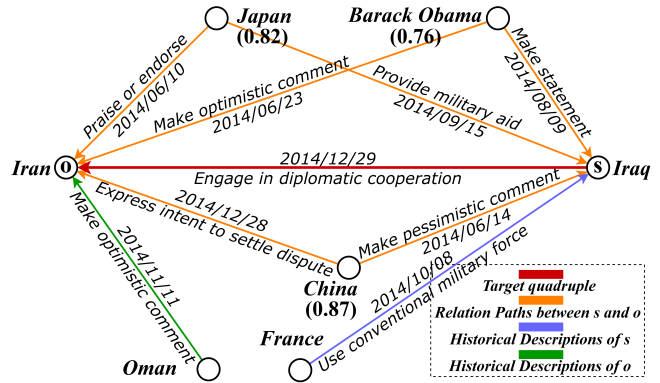


Figure 3: Explainability for a prediction case $(Iraq, Engage in diplomatic cooperation, ?, 2014-12-29)$ and its answer *Iran*. Numbers in brackets stand for the probabilities $P_{sentence\phi_s}$ of the structured sentences generated by SST-BERT.

First, the three relation paths and the two selected historical descriptions construct three structured sentences for the target quadruple, which have the probabilities $P_{sentence\phi_s}$ of 0.87, 0.82, 0.76, ranking top 3. The final score given by $f_{rt}(s, o)$ is 0.83, ranking first among candidate entities. Secondly, two historical facts, $(Japan, Provide military aid, Iraq, 2014/9/15)$ and $(France, Use conventional military force, Iraq, 2014/10/08)$, show the whole affair was military related. Finally, the edge $(China, Express intent to settle dispute, Iran, 2014/12/28)$ temporally and semantically closest to the target quadruple was the direct cause of the occurrence of *engaging in diplomatic cooperation* between *Iraq* and *Iran*.

5 CONCLUSION

In this paper, our model SST-BERT incorporates structured sentences with time-enhanced BERT as a comprehensively considered and time-sensitive solution to predict missing temporal relations and extends to the fully-inductive setting over TKGs. The rule-like relation paths in the natural language form enable SST-BERT to

reason in a flexible and explainable way. The historical descriptions enhance the target entities with easily accessible historical information inside TKGs and make SST-BERT external resource-independent. *TempBERT* pre-trained by our proposed *time masking* strategy in a specially generated TKG-aimed corpus rich in time tokens makes SST-BERT more sensitive to the time changing than BERT used in PLM-based baselines. We generate various benchmarks of the four datasets for the fully-inductive setting to evaluate SST-BERT comprehensively. Experiments show the outperformance of SST-BERT and the effectiveness of our proposed modules. Moreover, SST-BERT is robust to the different test sizes and different ratios of seen and unseen entities in the fully-inductive setting.

REFERENCES

- [1] Antoine Bordes, Nicolas Usunier, Alberto García-Durán, Jason Weston, and Oksana Yakhnenko. 2013. Translating Embeddings for Modeling Multi-relational Data. In *Advances in Neural Information Processing Systems 26: 27th Annual Conference on Neural Information Processing Systems 2013. Proceedings of a meeting held December 5-8, 2013, Lake Tahoe, Nevada, United States*, Christopher J. C. Burges, Léon Bottou, Zoubin Ghahramani, and Kilian Q. Weinberger (Eds.), 2787–2795. <https://proceedings.neurips.cc/paper/2013/hash/1cecc7a77928ca8133fa24680a88d2f9-Abstract.html>
- [2] Elizabeth Boschee, Jennifer Lautenschlager, Sean O'Brien, Steve Shellman, James Starz, and Michael Ward. 2015. ICEWS Event Coded Event Data.
- [3] Zhongwu Chen, Chengjin Xu, Fenglong Su, Zhen Huang, and You Dou. 2023. Meta-Learning Based Knowledge Extrapolation for Temporal Knowledge Graph. *arXiv e-prints*, Article arXiv:2302.05640 (Feb. 2023), arXiv:2302.05640 pages. <https://doi.org/10.48550/arXiv.2302.05640> arXiv:2302.05640 [cs.AI]
- [4] Shib Sankar Dasgupta, Swayambhu Nath Ray, and Partha Talukdar. 2018. HyTE: Hyperplane-based Temporally aware Knowledge Graph Embedding. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Brussels, Belgium, 2001–2011. <https://doi.org/10.18653/v1/D18-1225>
- [5] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. Association for Computational Linguistics, Minneapolis, Minnesota, 4171–4186. <https://doi.org/10.18653/v1/N19-1423>
- [6] Fredo Erxleben, Michael Günther, Markus Kröttsch, Julian Alfredo Mendez, and Denny Vrandečić. 2014. Introducing Wikidata to the Linked Data Web. In *International Workshop on the Semantic Web*.
- [7] Dieter A. Fensel, Umutkan Simsek, Kevin Angele, Elwin Huaman, Elias Kärle, Aleksandra Panasiuk, Ioan Toma, Jürgen Umbrich, and Alexander Wahler. 2020. Knowledge Graphs: Methodology, Tools and Selected Use Cases. *Knowledge Graphs* (2020).
- [8] Alberto García-Durán, Sebastijan Dumančić, and Mathias Niepert. 2018. Learning Sequence Encoders for Temporal Knowledge Graph Completion. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Brussels, Belgium, 4816–4821. <https://doi.org/10.18653/v1/D18-1516>
- [9] Rishab Goel, Seyed Mehran Kazemi, Marcus Brubaker, and Pascal Poupart. 2020. Diachronic Embedding for Temporal Knowledge Graph Completion. In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*. AAAI Press, 3988–3995. <https://aaai.org/ojs/index.php/AAAI/article/view/5815>
- [10] Zhen Han, Peng Chen, Yunpu Ma, and Volker Tresp. 2021. Explainable Subgraph Reasoning for Forecasting on Temporal Knowledge Graphs. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net. <https://openreview.net/forum?id=pGIHq1m7PU>
- [11] Shaoxiong Ji, Shirui Pan, Erik Cambria, Pekka Marttinen, and Philip S. Yu. 2022. A Survey on Knowledge Graphs: Representation, Acquisition, and Applications. *IEEE Transactions on Neural Networks and Learning Systems* 33, 2 (2022), 494–514. <https://doi.org/10.1109/TNNLS.2021.3070843>
- [12] Diederik P. Kingma and Jimmy Ba. 2015. Adam: A Method for Stochastic Optimization. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, Yoshua Bengio and Yann LeCun (Eds.). <http://arxiv.org/abs/1412.6980>
- [13] Timothée Lacroix, Guillaume Obozinski, and Nicolas Usunier. 2020. Tensor Decompositions for Temporal Knowledge Base Completion. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net. <https://openreview.net/forum?id=rke2P1BFwS>
- [14] Julien Leblay and Melisachew Wudage Chekol. 2018. Deriving Validity Time in Knowledge Graph. In *Companion Proceedings of the The Web Conference 2018 (Lyon, France) (WWW '18)*. International World Wide Web Conferences Steering Committee, Republic and Canton of Geneva, CHE, 1771–1776. <https://doi.org/10.1145/3184558.3191639>
- [15] S Liu, B Grau, I Horrocks, and EV Kostylev. 2021. INDIGO: GNN-based inductive knowledge graph completion using pair-wise encoding.
- [16] Yushan Liu, Yunpu Ma, Marcel Hildebrandt, Mitchell Joblin, and Volker Tresp. 2022. TLogic: Temporal Logical Rules for Explainable Link Forecasting on Temporal Knowledge Graphs. *AAAI* (2022).
- [17] Zhenghao Liu, Chenyan Xiong, Maosong Sun, and Zhiyuan Liu. 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)*. Association for Computational Linguistics, Melbourne, Australia, 2395–2405. <https://doi.org/10.18653/v1/P18-1223>
- [18] Xin Lv, Yankai Lin, Yixin Cao, Lei Hou, Juanzi Li, Zhiyuan Liu, Peng Li, and Jie Zhou. 2022. Do Pre-trained Models Benefit Knowledge Graph Completion? A Reliable Evaluation and a Reasonable Approach. In *Findings of the Association for Computational Linguistics: ACL 2022*. Association for Computational Linguistics, Dublin, Ireland, 3570–3581. <https://doi.org/10.18653/v1/2022.findings-acl.282>
- [19] Kaixin Ma, Hao Cheng, Xiaodong Liu, Eric Nyberg, and Jianfeng Gao. 2022. Open Domain Question Answering with A Unified Knowledge Interface. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, Dublin, Ireland, 1605–1620. <https://doi.org/10.18653/v1/2022.acl-long.113>
- [20] Sijie Mai, Shuangjia Zheng, Yuedong Yang, and Haifeng Hu. 2020. Communicative Message Passing for Inductive Relation Reasoning. In *AAAI Conference on Artificial Intelligence*.
- [21] Seyed Mehran Kazemi, Rishab Goel, Sepehr Eghbali, Janahan Ramanan, Jaspreet Sahota, Sanjay Thakur, Stella Wu, Cathal Smyth, Pascal Poupart, and Marcus Brubaker. 2019. Time2Vec: Learning a Vector Representation of Time. *ArXiv preprint abs/1907.05321* (2019). <https://arxiv.org/abs/1907.05321>
- [22] Minh Van Nguyen, Bonan Min, Franck Dernoncourt, and Thien Nguyen. 2022. Joint Extraction of Entities, Relations, and Events via Modeling Inter-Instance and Inter-Label Dependencies. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Association for Computational Linguistics, Seattle, United States, 4363–4374. <https://doi.org/10.18653/v1/2022.naacl-main.324>
- [23] Byungkook Oh, Seungmin Seo, Jimin Hwang, Dongho Lee, and Kyong-Ho Lee. 2022. Open-World Knowledge Graph Completion for Unseen Entities and Relations via Attentive Feature Aggregation. *Inf. Sci.* 586, C (2022), 468–484. <https://doi.org/10.1016/j.ins.2021.11.085>
- [24] Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and Alexander Miller. 2019. Language Models as Knowledge Bases?. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. Association for Computational Linguistics, Hong Kong, China, 2463–2473. <https://doi.org/10.18653/v1/D19-1250>
- [25] Guy D. Rosin, Ido Guy, and Kira Radinsky. 2022. Time Masking for Temporal Language Models. In *Proceedings of the Fifteenth ACM International Conference on Web Search and Data Mining (Virtual Event, AZ, USA) (WSDM '22)*. Association for Computing Machinery, New York, NY, USA, 833–841. <https://doi.org/10.1145/3488560.3498529>
- [26] Taylor Shin, Yasaman Razeghi, Robert L. Logan IV, Eric Wallace, and Sameer Singh. 2020. AutoPrompt: Eliciting Knowledge from Language Models with Automatically Generated Prompts. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational Linguistics, Online, 4222–4235. <https://doi.org/10.18653/v1/2020.emnlp-main.346>
- [27] Yusheng Su, Xu Han, Zhengyan Zhang, Yankai Lin, Peng Li, Zhiyuan Liu, Jie Zhou, and Maosong Sun. 2020. CokeBERT: Contextual Knowledge Selection and Embedding towards Enhanced Pre-Trained Language Models. *ArXiv preprint abs/2009.13964* (2020). <https://arxiv.org/abs/2009.13964>
- [28] Fabian M. Suchanek, Gjergji Kasneci, and Gerhard Weikum. 2007. Yago: a core of semantic knowledge. In *Proceedings of the 16th International Conference on World Wide Web, WWW 2007, Banff, Alberta, Canada, May 8-12, 2007*, Carey L. Williamson, Mary Ellen Zurko, Peter F. Patel-Schneider, and Prashant J. Shenoy (Eds.). ACM, 697–706. <https://doi.org/10.1145/1242572.1242667>
- [29] Tianxiang Sun, Yunfan Shao, Xipeng Qiu, Qipeng Guo, Yaru Hu, Xuanjing Huang, and Zheng Zhang. 2020. CoLAKE: Contextualized Language and Knowledge Embedding. In *Proceedings of the 28th International Conference on Computational Linguistics*. International Committee on Computational Linguistics, Barcelona, Spain (Online), 3660–3670. <https://doi.org/10.18653/v1/2020.coling-main.327>
- [30] Yu Sun, Shuohuan Wang, Yu-Kun Li, Shikun Feng, Hao Tian, Hua Wu, and Haifeng Wang. 2020. ERNIE 2.0: A Continual Pre-Training Framework for Language Understanding. In *The Thirty-Fourth AAAI Conference on Artificial Intelligence*.

- Intelligence, AAAI 2020, *The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020*, New York, NY, USA, February 7–12, 2020. AAAI Press, 8968–8975. <https://aaai.org/ojs/index.php/AAAI/article/view/6428>
- [31] Zhiqing Sun, Zhi-Hong Deng, Jian-Yun Nie, and Jian Tang. 2019. RotatE: Knowledge Graph Embedding by Relational Rotation in Complex Space. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6–9, 2019*. OpenReview.net. <https://openreview.net/forum?id=HkgEQnRqYQ>
- [32] Komal Teru, Etienne Denis, and Will Hamilton. 2020. Inductive Relation Prediction by Subgraph Reasoning. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13–18 July 2020, Virtual Event (Proceedings of Machine Learning Research, Vol. 119)*. PMLR, 9448–9457. <http://proceedings.mlr.press/v119/teru20a.html>
- [33] Komal Teru, Etienne Denis, and Will Hamilton. 2020. Inductive Relation Prediction by Subgraph Reasoning. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13–18 July 2020, Virtual Event (Proceedings of Machine Learning Research, Vol. 119)*. PMLR, 9448–9457. <http://proceedings.mlr.press/v119/teru20a.html>
- [34] Théo Trouillon, Johannes Welbl, Sebastian Riedel, Éric Gaussier, and Guillaume Bouchard. 2016. Complex Embeddings for Simple Link Prediction. In *Proceedings of the 33rd International Conference on Machine Learning, ICML 2016, New York City, NY, USA, June 19–24, 2016 (JMLR Workshop and Conference Proceedings, Vol. 48)*, Maria-Florina Balcan and Kilian Q. Weinberger (Eds.). JMLR.org, 2071–2080. <http://proceedings.mlr.press/v48/trouillon16.html>
- [35] Bo Wang, Tao Shen, Guodong Long, Tianyi Zhou, and Yi Chang. 2021. Structure-Augmented Text Representation Learning for Efficient Knowledge Graph Completion. *Proceedings of the Web Conference 2021* (2021).
- [36] Liang Wang, Wei Zhao, Zhuoyu Wei, and Jingming Liu. 2022. SimKGC: Simple Contrastive Knowledge Graph Completion with Pre-trained Language Models. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, Dublin, Ireland, 4281–4294. <https://doi.org/10.18653/v1/2022.acl-long.295>
- [37] Xintao Wang, Qi He, Jiaqing Liang, and Yanghua Xiao. 2022. Language Models as Knowledge Embeddings. *IJCAI abs/2206.12617* (2022).
- [38] Ruobing Xie, Zhiyuan Liu, Jia Jia, Huanbo Luan, and Maosong Sun. 2016. Representation Learning of Knowledge Graphs with Entity Descriptions. In *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence, February 12–17, 2016, Phoenix, Arizona, USA*, Dale Schuurmans and Michael P. Wellman (Eds.). AAAI Press, 2659–2665. <http://www.aaai.org/ocs/index.php/AAAI/AAAI16/paper/view/12216>
- [39] Chengjin Xu, Yung-Yu Chen, Mojtaba Nayyeri, and Jens Lehmann. 2021. Temporal Knowledge Graph Completion using a Linear Temporal Regularizer and Multivector Embeddings. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Association for Computational Linguistics, Online, 2569–2578. <https://doi.org/10.18653/v1/2021.naacl-main.202>
- [40] Chengjin Xu, Yung-Yu Chen, Mojtaba Nayyeri, and Jens Lehmann. 2021. Temporal Knowledge Graph Completion using a Linear Temporal Regularizer and Multivector Embeddings. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Association for Computational Linguistics, Online, 2569–2578. <https://doi.org/10.18653/v1/2021.naacl-main.202>
- [41] Chengjin Xu, Mojtaba Nayyeri, Fouad Alkhoury, Hamed Shariat Yazdi, and Jens Lehmann. 2020. TeRo: A Time-aware Knowledge Graph Embedding via Temporal Rotation. In *Proceedings of the 28th International Conference on Computational Linguistics*. International Committee on Computational Linguistics, Barcelona, Spain (Online), 1583–1593. <https://doi.org/10.18653/v1/2020.coling-main.139>
- [42] Chenjin Xu, Mojtaba Nayyeri, Fouad Alkhoury, Hamed Yazdi, and Jens Lehmann. 2020. Temporal knowledge graph completion based on time series gaussian embedding. In *International Semantic Web Conference*. Springer, 654–671.
- [43] Chengjin Xu, Mojtaba Nayyeri, Yung-Yu Chen, and Jens Lehmann. 2022. Geometric Algebra based Embeddings for Static and Temporal Knowledge Graph Completion. *CoRR abs/2202.09464* (2022). arXiv:2202.09464 <https://arxiv.org/abs/2202.09464>
- [44] Chengjin Xu, Fenglong Su, and Jens Lehmann. 2021. Time-aware Graph Neural Network for Entity Alignment between Temporal Knowledge Graphs. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Online and Punta Cana, Dominican Republic, 8999–9010. <https://doi.org/10.18653/v1/2021.emnlp-main.709>
- [45] Chengjin Xu, Fenglong Su, Bo Xiong, and Jens Lehmann. 2022. Time-Aware Entity Alignment Using Temporal Relational Attention (*WWW '22*). Association for Computing Machinery, New York, NY, USA, 788–797. <https://doi.org/10.1145/3485447.3511922>
- [46] Bishan Yang, Wen-tau Yih, Xiaodong He, Jianfeng Gao, and Li Deng. 2015. Embedding Entities and Relations for Learning and Inference in Knowledge Bases. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7–9, 2015, Conference Track Proceedings*, Yoshua Bengio and Yann LeCun (Eds.). <http://arxiv.org/abs/1412.6575>
- [47] Liang Yao, Chengsheng Mao, and Yuan Luo. 2019. KG-BERT: BERT for Knowledge Graph Completion. *ArXiv abs/1909.03193* (2019).
- [48] Hanwen Zha, Zhiyu Chen, and Xifeng Yan. 2022. Inductive Relation Prediction by BERT. In *AAAI Conference on Artificial Intelligence*.