

---

# GENTKG: Generative Forecasting on Temporal Knowledge Graph

---

Ruotong Liao<sup>1,2</sup>, Xu Jia<sup>3</sup>, Yunpu Ma<sup>1,4</sup>, Volker Tresp<sup>1,2</sup>  
<sup>1</sup>LMU Munich   <sup>2</sup>Munich Center for Machine Learning (MCML)  
<sup>3</sup>Technical University of Munich   <sup>4</sup>Siemens AG  
ruotong.liao@outlook.com,   volker.tresp@lmu.de

## Abstract

The rapid advancements in large language models (LLMs) have ignited interest in the temporal knowledge graph (tKG) domain, where conventional carefully designed embedding-based and rule-based models dominate. The question remains open of whether pre-trained LLMs can understand structured temporal relational data and replace them as the foundation model for temporal relational forecasting. Therefore, we bring temporal knowledge forecasting into the generative setting. However, challenges occur in the huge chasms between complex temporal graph data structure and sequential natural expressions LLMs can handle, and between the enormous data sizes of tKGs and heavy computation costs of finetuning LLMs. To address these challenges, we propose a novel retrieval augmented generation framework named GENTKG combining a temporal logical rule-based retrieval strategy and lightweight few-shot parameter-efficient instruction tuning to solve the above challenges. Extensive experiments have shown that GENTKG outperforms conventional methods of temporal relational forecasting under low computation resources with extremely limited training data as few as 16 samples. GENTKG also highlights remarkable cross-domain and in-domain generalizability with outperforming performance on unseen datasets without re-training. Our work reveals the huge potential of LLMs in the tKG domain and opens a new frontier for generative forecasting on tKGs.

## 1 Introduction

Forecasting the future lies in the intrinsic nature of humans to take controllability over the futural uncertainty ever since the existence of ancient fortunetellers who predict the future with insights into historical event rules. As the wave of Artificial General Intelligence (AGI) led by Large Language Models (LLMs) [4] showcases a persistent craving for World Models[20] that can model the complex information evolving in the real world, master the implicit rules and give predictions of what might happen next based on the historical observations[21], we term this challenge for LLMs as *Generative Forecasting*. We find Temporal Knowledge Graph (tKG) is a natural instance for investigating such a challenge attributed to the evolving world knowledge it contains and the task performed on it, namely *temporal knowledge graph forecasting*. In short sentence, tKGs can be viewed as streaming data sources where events come at the temporal axis hourly, daily, or yearly, etc., and tKG forecasting task aims to forecast future events at timestamp  $t + 1$  based on past historical events until  $t$ . Specifically, tKG originates from Knowledge Graph (KG) [22] which structures knowledge fact in the real world in the form of triples  $(e_s, r, e_o)$ , such as *(Paris, the capital of, France)*, where  $e_s, e_o$  represent the subject and object entity respectively, and  $r$  represents the observed predicate between the two entities. As world knowledge evolves constantly over time such as the inaugurated presidents of the USA, the Temporal Knowledge Graph (tKG) was introduced by [29] to indicate the temporal effectiveness of the world events by extending a timestamp  $t$  to form quadruples  $(e_s, r, e_o, t)$ . For example, *(Donald*

*Trump, the president of, the USA, 2021*) is followed by (*Joe Biden, the president of, the USA, 2023*). The tKG forecasting task aims to answer queries  $(e_s, r, ?, t)$  that predict the missing object given history events before  $t$ .

Conventional embedding-based graph representation learning methods [7, 8, 26, 31] require carefully designed models that embed indexed quadruples into hidden latent space and hence lose the semantic aspects of events in TKGs. Besides, it requires separate training for individual datasets and is hard to handle dataset modification and adaptation. In stark contrast, the rule-based methods [18] focus on mining temporal logic rules within the tKG graph structure in a symbolic way with limited scalability to different datasets without using neural networks. With the huge advancements emerging with numerous powerful large language models (LLMs) [30], early explorers have tried to utilize the emergent in-context learning ability of LLMs [6] to perform TKG forecasting tasks by flattening temporal knowledge facts to sequential texts but failed to compete with the above conventional methods [15]. It becomes questionable whether pre-trained LLMs can understand structured temporal relational data and replace both the embedding-based and rule-based approach as the foundation model for temporal relational forecasting.

To address the above issue, we bring temporal knowledge forecasting into the *generative forecasting* setting and deliberately prioritize the most influential factors in these two modalities: the temporal and structural characteristics of tKGs and the flexible natural language processing abilities of Large Language Models (LLMs). However, two challenges stand in the middle how to integrate them organically. The first challenge is the *modality challenge* between data structures. As tKG are complex temporal multi-relational graph data with tens of thousands of quadruples, it is hard to adapt to sequential natural language expressions that LLMs can process. The second challenge is the *computation challenge* with the enormous costs of fine-tuning LLMs especially with tens of thousands of quadruples. However, open-source and consumable training remains a core demand for advancing forefront research.

To solve the above two challenges, we propose **GENTKG**, a novel retrieval-augmented generation framework that solves the tKG forecasting task in the *generative forecasting* setting, outperforming both rule-based and embedding-based methods. Besides, GENTKG serves as an instantiation that sheds light on the promising *generative forecasting* ability of LLMs. For the first *modality challenge* between structured temporal graph data and sequential natural languages, we solve it in the retrieval phase. We utilize a temporal logical rule-based retrieval strategy that mines the temporal logic rules of the tKGs and then leverages these learned temporal logical rules to retrieve most temporally and logically relevant historical events to the give query. These events are then formalized as sequential natural languages in the temporal order and serve as specialized prompts to LLM. Although the prompts are in the form of sequential natural languages, they inherit structural information in the tKG implicitly since the extraction process are highly dependent on learned structural rules. The retrieval phase ends up with prompts that enable LLM to comprehend temporal relational data, regardless of the backbone LLM being used. We also deliver improvement over recent pure in-context learning method by introducing the retrieval phase on different LLMs only. Therefore, our retrieval strategy enables the input window of LLM to serve as the implicit and decouplable interface for communicating temporal and structural relational data to LLM.

For the second *computation challenge* between huge data and high computation costs, we solve it in the generation phase. We propose a few-shot parameter-efficient instruction-tuning framework for aligning LLMs with temporal relational forecasting task and reforming it into an autoregressive generation task. We further decompose the second *computation challenge* in two subtasks from the perspective of model and data respectively. The first subtask is to deal with the enormous computation costs and hardware requirements in training LLMs. We solve this challenge by a parameter-efficient fine-tuning (PEFT) adaptation method, specifically Low-rank Adaptation (LoRA)[13]. The second subtask is to deal with the enormous size of training data in tKGs. We deliberately think out of the box by bypassing learning the data as conventional methods and instead, letting the LLM learn the generative forecasting task on tKG. In other words, we reform data-centric model learning to task-centric LLM alignment that aligns LLMs with tKG forecasting task through instruction tuning, where we have specially designed task instruction, retrieved history as input, and generative predictions as output. Besides, we introduce few-shot tuning that further reduces training data to only 1024 prompt and response pairs which can be as few as 0.27% of original tens of thousands of training data with exceeding performance. Under extreme cases, we could further reduce to as few as 16 samples which is 0.0042% of original data while maintaining comparable performance to conventional methods.

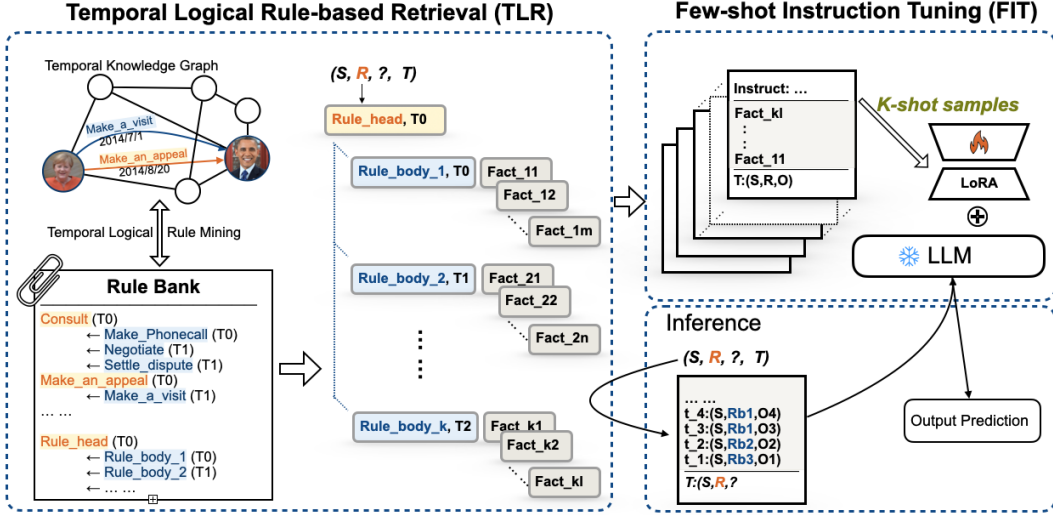


Figure 1: Framework of GENTKG. GENTKG first retrieves relevant facts based on a temporal logical rule-based retrieval strategy, then samples  $K$  prompts for few-shot parameter-efficient instruction-tuning of LLM that aligns LLM to the task of generative temporal knowledge graph forecasting.

Our approach offers a foundational framework for future explorations in generative forecasting on temporal knowledge graphs. Our contributions can be summarized as follows:

1. **Opening a frontier of generative forecasting on tKG.** To the best of our knowledge, we are the first to introduce instruction-tuned generative LLM to the tKG domain. Our framework GENTKG proposes a novel retrieval augmented generation paradigm for tKG forecasting, regardless of the backbone LLM.
2. **Drastically low computation costs with exceeding performance.** With only 16-shots parameter-efficient instruction tuning, we can already reach comparable results to conventional methods. With 1024-shots tuning, we can outperform existing rule-based, embedding-based, and the recent in-context-learning method.
3. **Task reformulation from data learning to task alignment.** We bypass designing specific models to learn specific tKG datasets. Instead, we novelly reform the data-centric learning to task-centric LLM alignment that aligns LLMs to generative forecasting task on tKG.
4. **Generalizability across datasets without retraining.** With one-time training on a single dataset, our GENTKG has showcased remarkably both cross-domain and in-domain generalizability with exceeding performance on multiple datasets without retraining.

## 2 Related Works

**Temporal Knowledge Graphs** Temporal knowledge graphs (tKGs) are multi-relational, directed graphs with labeled timestamped edges between entities (nodes). Let  $\mathcal{E}$  and  $\mathcal{P}$  represent a finite set of entities and predicates. A quadruple  $(e_s, r, e_o, t)$  represents a timestamped and labeled edge between a subject entity  $e_s \in \mathcal{E}$  and an object entity  $e_o \in \mathcal{E}$  at a timestamp  $t \in \mathcal{T}$ . Let  $\mathcal{F}$  represent the set of all true quadruples, i.e., real events in the world, the temporal knowledge graph forecasting is the task of predicting missing object entity at timestamp  $t$ , i.e.  $(e_s, r, ?, t)$  based on a set of observed facts  $\mathcal{O}$  before  $t$ , which is a subset of  $\mathcal{F}$ . Current methods can be categorized into two streams. On the one hand, embedding-based models learn representations of the quadruples with carefully designed embedding models[8, 7, 26, 11, 5]. On the other hand, the rule-based methods mine the temporal logical rules extracted and extract candidates directly on the temporal knowledge graphs[18].

**Investigating TKG with Language Models** The semantic part stored in the temporal knowledge graphs is heavily overlooked in either embedding-based or rule-based temporal knowledge graph methods. Early explorers had tryouts in introducing language models in the TKG domain, some

fused pre-trained language representations to the temporal knowledge embeddings [10], and some flattened explicit temporal events with the emergent in-context learning ability of large language models however not comparable with conventional performance [15]. Other researchers had tryouts in combing KG with LLM, utilizing the knowledge-aware prompting method [1, 25, 27, 32], however, cannot be transferred to the tKG domain due to their ignorance of temporal characteristics.

### 3 Generative Forecasting on Temporal Knowledge Graphs

In this section, we explain our GENTKG framework following its two-phase methodology: Retrieve-then-Generate, in two sections. Specifically, in Section 3.1, we explain our retrieval strategy, which commences by formulating a temporal-logic rule-based retrieval approach to capture historical facts that exhibit high temporal relevance and logical coherence. We name the retrieval strategy as TLR. Additionally, in Section 3.2, we delve into the details of our few-shot parameter-efficient instruction-finetuning which we termed as FIT, an essential component that aligns Large Language Models (LLMs) with the task of generative forecasting on temporal knowledge graphs.

#### 3.1 Temporal Logic Rule-based Retrieval

Our retrieval strategy is inspired by the phenomenon that a pair of entities can have many interactions at different timestamps such as a president visiting the same country multiple times. Another intuition behind this is that some relations tend to be temporally and logically sequential, for example in ICEWS14 we can see (*Angela Merkel, discuss by telephone, Barack Obama, 2014/07/22*) and (*Angela Merkel, consult, Barack Obama, 2014/08/09*). Therefore, we borrow a partial idea of TLogic[18] that mines the temporal logic rules hidden in the tKG structure. Notably, we opt to choose rules with a length equal to one that complies with the input length constraints of the LLMs. Then we propose a temporal logic rule-based retrieval strategy that retrieves the most temporally related and logically supportive history events for the given query based on these learned rules. To help understand our retrieval strategy, two definitions and the algorithm are given in the following.

**Definition I (Temporal Random Walk)** A non-increasing temporal random walk  $W$  starting from subject entity  $e_s \in \mathcal{E}$  to object entity  $e_o \in \mathcal{E}$  in the tKG  $\mathcal{G}$  is defined as a cycle of edges  $((e_s, r_1, e_o, t_2), (e_s, r_2, e_o, t_1))$  with  $t_2 \geq t_1$  where  $(e_s, r_i, e_o, t_i) \in \mathcal{G}$  and  $i \in 1, 2$ . The time constraints ensure that the edges are traversed only backward in time.

**Definition II (Temporal Logical Rule)** A cyclic temporal logical rule  $R$  is defined as  $(E_1, r_h, E_2, T_2) \leftarrow (E_1, r_b, E_2, T_1)$  with  $T_2 \geq T_1$ , where  $E_i$  and  $T_i$  for  $i \in 1, 2$  are replaceable variables that represent entities and timestamps. The left-hand side of  $R$  is called the rule head, with  $r_h$  being the head relation, while the right-hand side is called the rule body, with  $r_b$  being the body relation. A rule head can be supported by multiple rule bodies denoting different rules as  $\mathcal{TR}$ . A  $\mathcal{TR}$  implies that if the rule body holds then the rule head is true for a future timestamp  $T_2$ . The confidence of a rule  $\text{conf}(\mathcal{TR})$  is defined as dividing the rule support by the body support, where the support is the number of quadruples satisfying rule bodies or rule heads with time constraints.

**Rule Learning** Let  $r_h$  be a fixed relation, for which we want to learn rules. We sample an edge  $(e_1, r_h, e_2, t)$ , which will serve as the rule head, uniformly from all edges with relation  $r_h$ . Then the temporal random walker samples iteratively candidate edges adjacent to the current object  $\mathcal{C}(e_2, t) := \{(e_2, r, e_1, \hat{t}) \mid (e_2, r, e_1, \hat{t}) \in \mathcal{G}, \hat{t} < t\}$ , where  $\hat{t}$  is the timestamp associated with the next transition edge. Besides, we use an exponentially weighted transition distribution that prioritizes temporally closer edges during sampling which is defined as

$$\mathbb{P}(u; e_2, t) = \frac{\exp(t_u - t)}{\sum_{\hat{u} \in \mathcal{C}(e_2, t)} \exp(t_{\hat{u}} - t)} \quad (1)$$

where  $t_u$  denotes the timestamp of edge  $u$ . After a fixed sampling we can collect a set of temporal walks satisfying the rule  $(E_1, r_h, E_2, T_2) \leftarrow (E_1, r_b, E_2, T_1)$ . We then estimate the confidence of a rule following the definition II. Specifically, we count the number of body supports that there exists a relation head  $r_h$  connects  $e_1$  and  $e_2$  with the latest timestamp than all body timestamps.

Table 1: Instruction Prompt Design

PART	CONTENT
Task Instruction	You must be able to correctly predict the next {object_label} from a given text consisting of multiple quadruplets in the form of "{time}:{subject},{relation},{object_label}.{object}" and the query in the form of "{time}:{subject},{relation}," in the end. You must generate only the single number for {object_label} without any explanation.
Task Input	93: [Abdulrahman, Make_statement, 8092.Government_(Nigeria)] 113: [Abdulrahman, Make_statement, 8092.Government_(Nigeria)] 162: [Abdulrahman, Praise_or_endorse, 15546.Muslim_(Nigeria)] 197: [Abdulrahman, Consult, 8488.Governor_(Nigeria)] 197: [Abdulrahman, Make_statement, 8092.Government_(Nigeria)] 228: [Abdulrahman, Praise_or_endorse, 15414.Muhammadu_Buhari] 270: [Abdulrahman, Make_an_appeal_or_request, 3835.Citizen_(Nigeria)] 270: [Abdulrahman, Praise_or_endorse,
Task Output	3835.Citizen_(Nigeria)]

**Temporal Logic Rule-based History Retrieval** After gaining learned temporal logical rule sets, we order them according to the associated confidence scores. For a given forecast query  $(e_s, r, ?, t)$  we retrieve a candidate subgraph  $\mathcal{G}_s(e_s, r, t)$  from the TKG  $\mathcal{G}$  containing temporally and logically relevant histories for the given query, with respect to the subject entity, relation, and timestamp. Since the query subject entity is fixed, there are two key factors in the retrieval algorithm, i.e. time window and rule grounding. First, we define the time window as  $TW = [t_-, t]$  with  $t_- := t - w$ , where the  $w \in \mathbb{N}^+$  represents the time window length backward starting from the query timestamp. The maximum length of  $w$  is  $\min\{t_{max}, t\}$  with  $t_{max}$  denoting the maximum timestamp of the datasets. Second, the query relation is fixed as a rule head. Within each  $TW$ , we apply the learned rules  $\mathcal{TR}$  and select top  $k$  various rule bodies  $r_{b_1}, r_{b_2}, \dots, r_{b_k}$  regarding to  $r$  in descending confidence and add historical events  $(e_s, r_b, e_o, t - w)$  to  $\mathcal{G}_s(e_s, r, t)$  for the given query. The size of  $\mathcal{G}_s(e_s, r, t)$  can be adjusted dynamically with respect to  $w$  and  $k$ . We stop the retrieval until a maximum history length  $N$  is reached. For instance, we retrieve history events iteratively with the top confident rule bodies for each time window backtrace step until a maximum history length of 50 is reached. At the end of the retrieval phase, we reorder all history events in descending order for each query.

### 3.2 Aligning LLMs to generative forecasting on TKG

In our proposed GENTKG framework, we contribute to transforming the conventional tKG forecasting task into an alignment task that aligns LLM with generative forecasting on tKGs. In this subsection, we present the second generation phase in our GENTKG framework utilizing few-shot parameter-efficient instruction tuning in the settings of low GPU resource consumption. Specifically, we first describe the instruction prompt design in 3.2.1, then describe the parameter-efficient instruction tuning for training our generative model in 3.2.2. Then, we explain the few-shot tuning strategy that efficiently aligns the LLM with temporal relational forecasting with as few as 1024 samples in 3.2.3 and explore the lower-bound of samples for few-shot tuning. Finally, we explore the possibility of the inductive capability of generative forecasting on tKGs in 3.2.4.

#### 3.2.1 Instruction Prompt Design

Instruction Tuning is a crucial technique that fine-tunes LLMs with human-curated instruction and response pairs as the training data, empowering LLMs with instruction-following capability. [33] The construction of an instruction pair sample is usually composed of three parts, i.e. task instruction, task input, and task output. Task instruction clarifies the definition of the task for LLMs to comprehend and gives explicit solutions for LLMs to follow and execute. Task input in natural languages is input data serving as context for LLMs and task output is the decoding results based on the input prompt, which consists of task instruction and task input. In our proposed GENTKG framework, we adapt the temporal knowledge graph forecasting task to the instruction task for LLMs with individual adaptation for the three parts partially following the setting in [15]. As a demonstration, we design the instruction as depicted in table 1. Except for the designed task instruction, the task input is modeled as ordered historical events retrieved from the previous phase for a given query  $(e_s, r, e_o, t)$  as described in 3.1. Each fact is filled in the template of " $t : [e_s, r, n_{e_o}.e_o]$ ". It is worth noting that

we map each candidate object  $e_o$  with a numerical index  $n_{e_o}$  to avoid unfair tokenization for different entities [28]. In this way, we use the probabilities of index tokens generated by the LLMs to get ranked scores of output entities in an indirect way. The query  $(e_s, r, e_o, t)$  is expressed in a similar but partial way as “ $t : [e_s, r,$ “ for LLM to complete as generative predictions.

### 3.2.2 Lightweight Parameter-efficient Instruction Tuning

Direct fine-tuning of the entire model is computationally demanding and time-consuming. To address these computational challenges, we adopt the Low-Rank Adaptation (LoRA) technique [13]. LoRA involves the freezing of pre-trained model parameters  $\theta_0$  while introducing trainable additional parameters  $\theta_1$  that can be decomposed into low-rank matrices  $\Delta\theta_0 = \mathbf{BA}$ ,  $\mathbf{B} \in \mathbb{R}^{d \times r}$ ,  $\mathbf{A} \in \mathbb{R}^{r \times k}$ ,  $r \ll \min(d, k)$ . By optimizing these rank decomposition matrices, we can efficiently incorporate supplementary information into the fixed original model.

At present, there are large amounts of LLMs released, such as GPT series [14, 23], T5 series [24], CHinchilla [12], and LLaMA [28]. Among these, proprietary models can only be accessed by APIs such as ChatGPT with limited adaptation and alignment possibilities that hinder the research purpose. To facilitate the research in generative forecasting on temporal knowledge graphs, we carefully opt for the open-sourcing LLMs, i.e. GPT-NeoX-20B [2] and LLaMA2-7B [28], which is the third-party reproduction of GPT-3 and open-sourcing public model respectively. Due to hardware limitations, we leave GPT-NeoX-20B frozen to investigate the effectiveness of our retrieval phase through its in-context learning ability. We perform the whole GENTKG framework on LLaMA2-7B for consumable adaptation purposes.

### 3.2.3 Efficient Alignment with Few-shot Tuning

Our framework contributes a remarkably efficient and effective few-shot training strategy. The hypothesis has been proven that alignment can be a simple process where the LLMs learn the style or format for responding to prompts and expose the knowledge and capabilities that were already acquired during pretraining [33]. Therefore, considering the volume of temporal knowledge graphs that usually possess tens of thousands of training data, we propose a  $K$ -shot training paradigm where only an extremely limited number of  $K$  samples are uniformly sampled from the temporal-ordered training set for language model adaptations. In our case, we select only 1024 samples which takes up as few as 0.27% of the original GDELT dataset sizes that conventional methods usually fully trained on. We further prove that our method can acquire temporal relational forecasting capability rapidly with severely limited training data (0.0027%) with an extreme 16-shot training setting while maintaining comparable performances to conventional method.

### 3.2.4 Inductive Setting

Our framework GENTKG also delivers remarkable performance gains in inductive settings. **Cross-domain generalizability.** GENTKG trained on one dataset can be inferred on other datasets. It only requires learning the temporal-logical rule-based retrieval strategy for the new datasets in the first phase to ensure proper prompts with relevant histories. However, it doesn’t require retraining LLM in the second phase. Still, high-performance gains are maintained and even comparable to the original setting.

**In-domain generalizability.** GENTKG maintains high-performance gains on the same dataset even the training data is limited to a small partition of original training data, such as 5%. This characteristic exceeds conventional methods which always suffer drastic performance drops once the critical value of the forecasting timestamp between the train and evaluation set changes.

We attribute this ability to the reason that the LLM is not fine-tuned specifically to a single dataset. Instead, the LLM is aligned to the temporal relational forecasting task itself rather than the learning of the tKG data.

## 4 Experimental Setup

In this section, we describe the experimental setup of our GENTKG framework. First, we describe four datasets then evaluation protocols, and finally the experimental design.

**Datasets** We use four benchmark datasets to evaluate GENTKG: 1) ICEWS14 [3] 2) ICEWS18 [3] 3) GDEL T [16] 4) YAGO [19]. The two versions of the Integrated Crisis Early Warning System (ICEWS) both consist of timestamped political events, e.g., (Angela Merkel, visit, India, 2015-03-25). The GDEL T and YAGO datasets are extracted from the subsets of GDEL T and YAGO knowledge bases containing facts and time information. Detailed information is shown in table 3 in the Appendix.

**Evaluation** Since GENTKG generates entity predictions directly, we use raw Hits@1/3/10 metric to evaluate the model performance on extrapolated link prediction. Hits@1/3/10 denotes the proportion of the actual missing entities ranked within the top 1/3/10.

**Baselines** Since our framework GENTKG is the first method to introduce retrieval-augmented instruction-tuned generative models into the temporal knowledge forecasting domain, it is necessary to include three types of existing methods as baselines. The first are embedding-based methods, represented by RE-GCN[17], xERTE [8], TANGO [9], and Timetraveler[26]. The rule-based method is TLogic[18] and the third type is the LLM-based ICL method with frozen parameters[15].

**Experiment Design** In order to comprehensively analyze our method compared to different conventional methods, there are three research questions to be answered. RQ1: How is the overall performance of the proposed GENTKG framework compared with the existing conventional embedding-based, rule-based TKG methods and LLM-based ICL method? RQ2: How well does the GENTKG framework enable the LLM to acquire cross-domain and in-domain generalizability with different datasets? RQ3: How do the components of the proposed GENTKG affect its effectiveness?

## 5 Experimental Results

Table 2: Temporal knowledge forecasting results: Hits@1/3/10. The best results among each metric are highlighted in bold and the second bests are underlined. Notably, the TLR denotes only using the retrieval phase in GENTKG.

Datasets		ICEWS14			ICEWS18			GDEL T			YAGO		
Method Type	Model	Hits@1	Hits@3	Hits@10	Hits@1	Hits@3	Hits@10	Hits@1	Hits@3	Hits@10	Hits@1	Hits@3	Hits@10
Embedding-based	RE-GCN	0.313	0.473	<b>0.626</b>	0.223	0.367	<b>0.525</b>	0.084	0.171	0.299	0.468	0.607	0.729
	xERTE	0.330	0.454	0.570	<u>0.209</u>	0.335	<u>0.462</u>	0.112	0.191	0.294	<b>0.769</b>	<u>0.787</u>	0.794
	TANGO	0.272	0.408	0.550	0.191	0.318	<u>0.462</u>	0.094	0.189	0.322	0.566	0.651	0.718
	Timetraveler	0.319	0.454	0.575	0.212	0.325	0.439	0.112	0.186	0.285	0.604	0.770	0.831
Rule-based	TLogic	0.332	0.476	0.602	0.204	0.336	0.480	0.113	0.212	<u>0.351</u>	0.638	0.650	0.660
In-Context Learning	GPT-NeoX-20B	0.319	0.439	0.538	0.179	0.297	0.41	0.098	0.165	0.253	0.669	<u>0.787</u>	<u>0.841</u>
	Llama2-7B	0.252	0.427	0.504	0.128	0.272	0.323	0.06	0.164	0.246	0.662	0.760	0.818
GENTKG	GPT-NeoX-20B +TLR	<u>0.35</u>	<u>0.485</u>	<u>0.593</u>	0.205	0.338	0.462	<u>0.156</u>	<u>0.241</u>	<u>0.349</u>	0.681	<b>0.807</b>	<b>0.861</b>
	Llama2-7B +GENTKG	<b>0.372</b>	<b>0.488</b>	0.563	0.217	<b>0.372</b>	0.443	<b>0.185</b>	<b>0.278</b>	<b>0.358</b>	0.715	0.767	0.786
	Llama2-7B - inductive	-	-	-	<b>0.234</b>	0.333	0.409	0.142	0.213	0.281	0.638	0.697	0.762

### 5.1 Performance Comparison

Experiment results can be seen on Table 2. To answer the first question RQ1, our results achieve state-of-the-art performance, surpassing all three types of existing conventional including embedding-based models, rule-based method, and LLM-based in-context learning method across four datasets regarding metric Hit@1 and Hit@3 while maintaining comparable results regarding Hits@10. Our method demonstrates the promising trend for retrieval-augmented large language models to serve as the foundation model for temporal relational forecasting, opening a new frontier in the TKG domain. More detailed results and analyses are presented in the following. We refer to GENTKG utilizing LLaMA2-7B as instantiation unless otherwise specified.

**Compared to embedding-based models.** For all datasets, GENTKG outperforms its best embedding-based model xERTE on ICEWS14, ICEWS18, GDEL T, and Timetraveler on YAGO. Specifically, the highest performance gain is observed on GDEL T with 58% higher on Hits@1. It is natural to conclude that GENTKG can outperform embedding-based methods.

**Compared to the rule-based model.** Compared to the rule-based model TLogic, GENTKG outperforms TLogic on Hits@1 and Hits@3 while maintaining comparable performance regarding Hits@10 on GDEL T and YAGO. The slight drops in Hits@10 on ICEWS14 and ICEWS18 are because TLogic is carefully designed on these datasets while our method has more generalizability and demonstrated better performance regarding accuracy than recall.

**Compared to in-context-learning method.** We analyze the performance of GENTKG on different

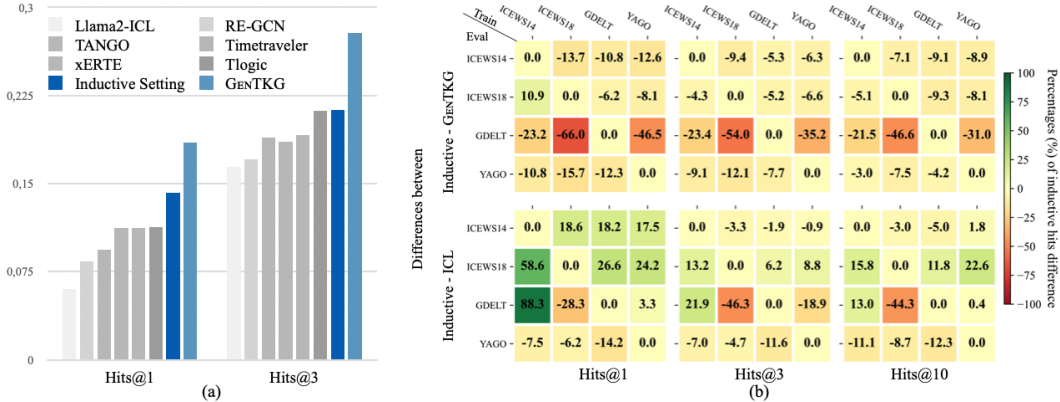


Figure 2: Cross-Domain Inductive Setting. (a) All models including GENTKG are trained and evaluated on the GDELT dataset, except that the inductive setting of GENTKG is trained on ICEWS14 and evaluated on GDELT. (b) Cross-checking. We cross-check the trained LLaMA2 in GENTKG on different training datasets and evaluation datasets. The performance drop compared to the original training setting takes up only small percentages. Even higher performance than ICL can be observed.

Language Model instantiations, i.e. GPT-NeoX-20B and LLaMA2-7B respectively. For GPT-NeoX-20B, we apply only the first retrieval phase of GENTKG due to hardware limitations. However, an average 10% performance increase is observed for all three metrics on all datasets even with pure retrieval-augmented in-context learning. For LLaMA2-7B, the performance gain of Hits@1 has increased by 47% on ICEWS14 and 69% on ICEWS18, even outperforming GPT-NeoX-20B which has two times more parameters, indicating the potential for greater performance of our proposed GENTKG framework if applied to larger language models.

## 5.2 Cross-domain Generalization

To answer the second question of GENTKG’s performance in the inductive setting, the empirical results indicate that the GENTKG framework manifests a substantial capability for cross-dataset generalization. Specifically, once the large language model has been aligned to the tKG forecasting task in the second phase on any dataset, the large language model can be applied directly to any other dataset. Therefore, on a new dataset, GENTKG only requires dataset-specific temporal-logical rule-based retrieval to formulate proper prompts in the first phase, and can directly infer the predictions without retraining in the second phase. As shown in Figure 2, when the model is trained exclusively on the ICEWS14 dataset, it still delivers comparable performance metrics on GDELT to conventional methods with a minor performance drop compared to the full fine-tuned version. We further demonstrate similar inductive results by cross-checking the training and evaluation datasets as shown in Figure 2. Although the model is trained exclusively on one dataset, it still delivers comparable metrics on disparate datasets, closely approximating the outcomes of models that were trained specifically on those identical datasets. This notable characteristic implies that the GENTKG framework is effectively capturing the underlying task-related features, as opposed to merely designed or trained specifically on the identical dataset, a limitation commonly shared in traditional approaches.

## 5.3 In-domain Generalization

Apart from generalization across different datasets, how well does GENTKG generalize to different training partitions within the same datasets? To investigate such a problem, we carefully designed various partitions of training data ranging in {5%, 10%, 20%, 30%, 50%, 75%, 100%}. All models trained on different training partitions are evaluated on the same evaluation set. Experiments have shown that conventional methods suffer from insufficient training data while GENTKG remains exceeding performance even with as few as 5% training data. This further proves that GENTKG successfully transforms conventional data-centric learning to the task-centric alignment of LLMs and overcomes the prediction instability under the incoming data stream in the forecasting setting.



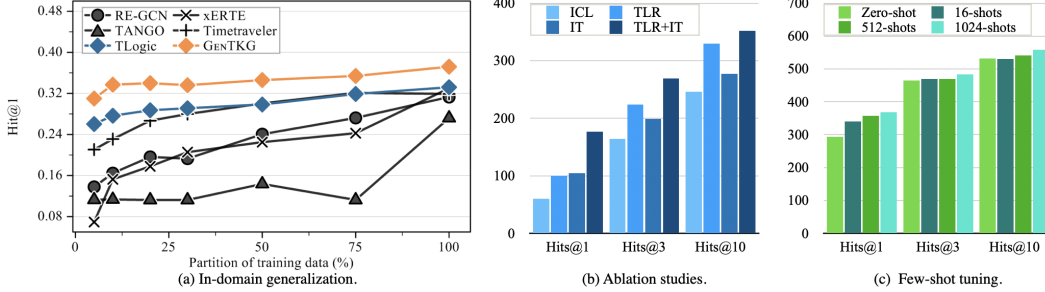


Figure 3: (a) In-domain generalizability of GENTKG exceeds conventional methods on all different partitions of training data on ICEWS14. (2) Ablation study of TLR and FIT phase in GENTKG framework. (c) Increasing few-shot training parameter  $K$  improves performance.

#### 5.4 Ablation studies

To comprehensively evaluate the contribution of each constituent element within our proposed GENTKG framework, we undertake the ablation studies on the ICEWS14 dataset. We assess the performance metrics across three distinct variants of the GENTKG framework: TLR, FIT, and TLR+FIT configurations. Here, TLR represents the variant that exclusively employs temporal logical rule-based retrieval, FIT denotes the variant solely implementing few-shot parameter-efficient instruction tuning, and TLR+FIT encapsulates the integration of all components within our proposed methodology GENTKG. We demonstrate the results in Figure 3. We observe that both phases in our proposed GENTKG framework contribute to distinct performance improvements. The whole pipeline enables GENTKG the ability to outperform existing methods.

#### 5.5 Few-shot Tuning

To delve further into the impact of sample size within the few-shot tuning paradigm, we conducted a series of experiments on the ICEWS14 dataset employing a range of shot sizes  $K$  from the set [16, 512, 1024]. For each configuration, we employed uniform sampling on the temporally-ordered training dataset. Empirical results indicate a consistent trend of performance improvement correlating with an increase in the number of training samples, as visualized in Figure 3. Remarkably, our findings suggest that the GENTKG framework is capable of achieving performance metrics closely approximating those of the default  $K = 1024$  setting, even when as few as 16 shots are used for tuning. This notable finding unlocks significant potential for GENTKG in the context of aligning LLMs with temporal relational forecasting tasks from the perspective of efficient alignment.

## 6 Limitations

Although our method significantly outperforms the existing approaches in various types, our approach is limited by the input context window of large language models. Specifically, for LLaMA2, the input context window is 4096 tokens with an average upper length limit of 50 history facts. We leave this to future work.

## 7 Conclusion

In this paper, we first raise the question of whether pre-trained LLMs can understand structured temporal relational data and replace existing tKG models as the foundation model for temporal relational forecasting tasks. We propose a GENTKG framework that can efficiently align LLM with temporal relational tasks through two stages: temporal logical rule-based retrieval and few-shot parameter-efficient fine-tuning. Extensive experimental results demonstrate that GENTKG framework outperforms conventional embedding-based, rule-based and ICL methods. Even though, our proposed framework is training-light and capable of alignment through consumable computation resources with extremely few training data, and exhibits strong cross-domain and in-domain transferability.

## References

- [1] Jinheon Baek, Alham Fikri Aji, and Amir Saffari. Knowledge-augmented language model prompting for zero-shot knowledge graph question answering. In *Proceedings of the 1st Workshop on Natural Language Reasoning and Structured Explanations (NLRSE)*, pages 78–106, Toronto, Canada, June 2023. Association for Computational Linguistics.
- [2] Sid Black, Stella Biderman, Eric Hallahan, Quentin Anthony, Leo Gao, Laurence Golding, Horace He, Connor Leahy, Kyle McDonell, Jason Phang, Michael Pieler, USVSN Sai Prashanth, Shivanshu Purohit, Laria Reynolds, Jonathan Tow, Ben Wang, and Samuel Weinbach. GPT-NeoX-20B: An open-source autoregressive language model. In *Proceedings of the ACL Workshop on Challenges & Perspectives in Creating Large Language Models*, 2022.
- [3] Elizabeth Boschee, Jennifer Lautenschlager, Sean O’Brien, Steve Shellman, James Starz, and Michael Ward. Icews coded event data. *Harvard Dataverse*, 12, 2015.
- [4] Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. Sparks of artificial general intelligence: Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*, 2023.
- [5] Zifeng Ding, Jingpei Wu, Bailan He, Yunpu Ma, Zhen Han, and Volker Tresp. Few-shot inductive learning on temporal knowledge graphs using concept-aware information. *arXiv preprint arXiv:2211.08169*, 2022.
- [6] Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiyong Wu, Baobao Chang, Xu Sun, Jingjing Xu, and Zhifang Sui. A survey for in-context learning. *arXiv preprint arXiv:2301.00234*, 2022.
- [7] Rishab Goel, Seyed Mehran Kazemi, Marcus Brubaker, and Pascal Poupart. Diachronic embedding for temporal knowledge graph completion. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 3988–3995, 2020.
- [8] Zhen Han, Peng Chen, Yunpu Ma, and Volker Tresp. Explainable subgraph reasoning for forecasting on temporal knowledge graphs. In *International Conference on Learning Representations*, 2020.
- [9] Zhen Han, Zifeng Ding, Yunpu Ma, Yujia Gu, and Volker Tresp. Learning neural ordinary equations for forecasting future links on temporal knowledge graphs. In *Proceedings of the 2021 conference on empirical methods in natural language processing*, pages 8352–8364, 2021.
- [10] Zhen Han, Ruotong Liao, Beiyan Liu, Yao Zhang, Zifeng Ding, Jindong Gu, Heinz Koeppel, Hinrich Schuetze, and Volker Tresp. Enhanced temporal knowledge embeddings with contextualized language representations. 2022.
- [11] Zhen Han, Yunpu Ma, Peng Chen, and Volker Tresp. Dyernie: Dynamic evolution of riemannian manifold embeddings for temporal knowledge graph completion. *arXiv preprint arXiv:2011.03984*, 2020.
- [12] Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al. Training compute-optimal large language models. *arXiv preprint arXiv:2203.15556*, 2022.
- [13] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.
- [14] Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large language models are zero-shot reasoners. *Advances in neural information processing systems*, 35:22199–22213, 2022.
- [15] Dong-Ho Lee, Kian Ahrabian, Woojeong Jin, Fred Morstatter, and Jay Pujara. Temporal knowledge graph forecasting without knowledge using in-context learning. *arXiv preprint arXiv:2305.10613*, 2023.

- [16] Kalem Leetaru and Philip A Schrod. Gdelt: Global data on events, location, and tone, 1979–2012. In *ISA annual convention*, volume 2, pages 1–49. Citeseer, 2013.
- [17] Zixuan Li, Xiaolong Jin, Wei Li, Saiping Guan, Jiafeng Guo, Huawei Shen, Yuanzhuo Wang, and Xueqi Cheng. Temporal knowledge graph reasoning based on evolutionary representation learning. 2021.
- [18] Yushan Liu, Yunpu Ma, Marcel Hildebrandt, Mitchell Joblin, and Volker Tresp. Tlogic: Temporal logical rules for explainable link forecasting on temporal knowledge graphs. In *Proceedings of the AAAI conference on artificial intelligence*, volume 36, pages 4120–4127, 2022.
- [19] Farzaneh Mahdisoltani, Joanna Biega, and Fabian M Suchanek. Yago3: A knowledge base from multilingual wikipedias. In *CIDR*, 2013.
- [20] Yutaka Matsuo, Yann LeCun, Maneesh Sahani, Doina Precup, David Silver, Masashi Sugiyama, Eiji Uchibe, and Jun Morimoto. Deep learning, reinforcement learning, and world models. *Neural Networks*, 152:267–275, 2022.
- [21] Grégoire Mialon, Roberto Dessì, Maria Lomeli, Christoforos Nalmpantis, Ram Pasunuru, Roberta Raileanu, Baptiste Rozière, Timo Schick, Jane Dwivedi-Yu, Asli Celikyilmaz, et al. Augmented language models: a survey. *arXiv preprint arXiv:2302.07842*, 2023.
- [22] Maximilian Nickel, Kevin Murphy, Volker Tresp, and Evgeniy Gabrilovich. A review of relational machine learning for knowledge graphs. *Proceedings of the IEEE*, 104(1):11–33, 2015.
- [23] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- [24] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551, 2020.
- [25] Md Rashad Al Hasan Rony, Ricardo Usbeck, and Jens Lehmann. Dialokg: Knowledge-structure aware task-oriented dialogue generation. *arXiv preprint arXiv:2204.09149*, 2022.
- [26] Haohai Sun, Jialun Zhong, Yunpu Ma, Zhen Han, and Kun He. Timetraveler: Reinforcement learning for temporal knowledge graph forecasting. *arXiv preprint arXiv:2109.04101*, 2021.
- [27] Jiashuo Sun, Chengjin Xu, Lumingyuan Tang, Saizhuo Wang, Chen Lin, Yeyun Gong, Heung-Yeung Shum, and Jian Guo. Think-on-graph: Deep and responsible reasoning of large language model with knowledge graph. *arXiv preprint arXiv:2307.07697*, 2023.
- [28] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- [29] Volker Tresp, Cristóbal Esteban, Yinchong Yang, Stephan Baier, and Denis Krompaß. Learning with memory embeddings. *arXiv preprint arXiv:1511.07972*, 2015.
- [30] Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. Emergent abilities of large language models. *arXiv preprint arXiv:2206.07682*, 2022.
- [31] Jinzhu Yang, Wei Zhou, Lingwei Wei, Junyu Lin, Jizhong Han, and Songlin Hu. Re-gcn: relation enhanced graph convolutional network for entity alignment in heterogeneous knowledge graphs. In *Database Systems for Advanced Applications: 25th International Conference, DASFAA 2020, Jeju, South Korea, September 24–27, 2020, Proceedings, Part II 25*, pages 432–447. Springer, 2020.
- [32] Miao Zhang, Rufeng Dai, Ming Dong, and Tingting He. Drlk: dynamic hierarchical reasoning with language model and knowledge graph for question answering. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 5123–5133, 2022.

- [33] Chunting Zhou, Pengfei Liu, Puxin Xu, Srinu Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, et al. Lima: Less is more for alignment. *arXiv preprint arXiv:2305.11206*, 2023.

Models	GPT-NeoX-20B	GPT-NeoX-20B+TLR	LLaMA2-7B	LLaMA2-7B+TLR	Llama2-7B-FT+TLR	
Datasets	inference time	inference time	inference time	inference time	inference time	finetune time
ICEWS14	0.15 min	0.14 min	0.14 min	0.14 min	0.67 min	126 h
ICEWS18	0.25 min	0.23 min	0.24 min	0.22 min	1.11 min	186 h
GDELT	0.30 min	0.29 min	0.28 min	0.28 min	1.34 min	172 h
YAGO	0.08 min	0.07 min	0.07 min	0.07 min	0.34 min	120 h

Figure 4: Computation Costs

## Appendix

### 7.1 Dataset Statistics.

Table 3: Dataset statistics.

Datasets	#train	#valid	#test	#entity	#relations	time gap
ICEWS14	74854	8514	7371	7128	230	1 day
ICEWS18	373018	45995	49545	23033	256	1 day
GDELT	79319	9957	9715	5850	238	15 mins
YAGO	220393	28948	22765	10778	23	1 year

### 7.2 Computation Costs.

In our experiments, we utilized A100 GPUs for both the finetuning and inference processes. Figure 4 provides statistics on the average inference time per sample and the GPU hours required for each finetuning experiment across different models and datasets.