

From Decomposition to Validation: A Multi-stage Framework for Temporal Knowledge Graph Question Answering with LLMs

Anonymous ACL submission

Abstract

Temporal Knowledge Graph Question Answering (TKGQA) requires reasoning over time-dependent facts. However, existing approaches suffer from two fundamental limitations: (1) general-purpose text retrievers often fail to properly align entities, relations, and temporal constraints in the question with relevant TKG facts; and (2) large language models (LLMs) remain unreliable in temporal reasoning, frequently exhibiting hallucinations. To address these challenges, we propose **TempReasoner**, a multi-stage framework for TKGQA. TempReasoner decomposes complex temporal questions into structured sub-questions, retrieves temporally relevant facts using a trained TKG-Retriever, filters noisy facts through LLM-generated alignment signals, and generates answers along with supporting evidence set. A dedicated validation module further refines uncertain predictions by issuing verification queries, improving both answer accuracy and reliability. Experiments on two widely used benchmarks, MultiTQ and CronQuestions, show that TempReasoner consistently delivers strong performance. In particular, TempReasoner achieves substantial improvements on multi-hop temporal questions, demonstrating its effectiveness in temporal reasoning and evidence grounding. Additional analysis further confirms the efficiency of its retrieval component and consistent model-agnostic generalization.

1 Introduction

Many real-world questions involve explicit or implicit temporal constraints, requiring accurate reasoning over time. However, current large language models (LLMs) remain unreliable when answering temporal questions. As illustrated in Figure 1(a), directly querying LLMs with temporal questions often leads to incorrect answers without background knowledge. To better support temporally grounded reasoning, the Temporal Knowledge Graph Ques-

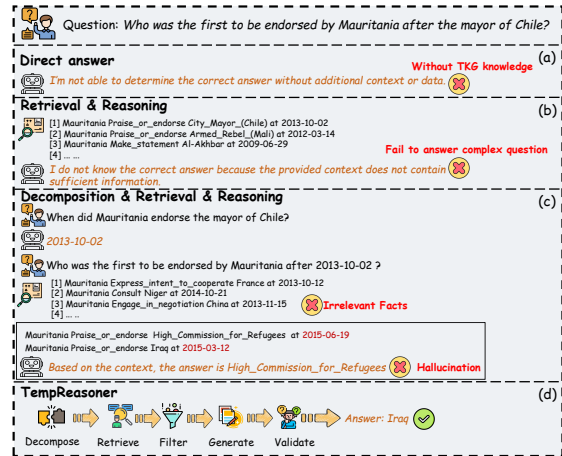


Figure 1: Illustration of the paradigms in the current LLM-based TKGQA task. (a) Direct LLM answering. (b) Retrieval & Reasoning. (c) Decomposition & Retrieval & Reasoning. (d) Our proposed TempReasoner

tion Answering (TKGQA) task has been introduced, which leverages temporal knowledge graphs (TKGs) as structured background knowledge to represent time-dependent facts and facilitate temporal reasoning (Chen et al., 2024c; Gao et al., 2024; Qian et al., 2024; Gong et al., 2025; QianyiHu et al., 2025).

In TKGQA, a straightforward **Retrieval & Reasoning** paradigm, as shown in Figure 1(b), is often insufficient for complex multi-hop temporal questions. To mitigate this limitation, several methods (Chen et al., 2024c; Gong et al., 2025) adopt a **Decomposition & Retrieval & Reasoning** strategy that decomposes complex questions into sub-questions, performs retrieval and reasoning for each sub-question, and aggregates the intermediate results to produce the final answer. Although this strategy alleviates the difficulty of complex temporal reasoning to some extent, it still faces several fundamental challenges:

Limitations of retrievers for TKGs. TKGs contain a large number of semantically similar entities and relations, making it easy for retrieval models

to return many superficially related but ultimately irrelevant facts. As illustrated in Figure 1(c), when using a general-purpose encoder to retrieve facts directly from a question: "Who was the first to be endorsed by Mauritania after 2013-10-02?", the relation should be *Praise_or_endorse*. However, the retrieved results often include semantically proximate but incorrect candidates: "Mauritania Engage_in_negotiation China at 2013-11-15".

Limitations of LLMs in generating answers for temporal reasoning. Although LLMs have demonstrated strong performance across a wide range of generation tasks, they still exhibit notable hallucinations in TKGQA. As illustrated in Figure 1(c), LLMs generate incorrect conclusions even when provided with relevant contextual facts (*Mauritania Praise_or_endorse Iraq at 2015-03-12* in this case). To further investigate this limitation, we conducted an experiment using GPT-4o-mini¹ and DeepSeek-V3². Both models were asked to answer *first* or *last* temporal questions³ based on 100 retrieved facts presented in random order. GPT-4o-mini achieved only 36% Hits@1 accuracy, while DeepSeek-V3 performed slightly better at 48%. These results reveal that current LLMs still struggle with reliable temporal reasoning even when contextual evidence is provided.

To overcome these challenges, we propose **TempReasoner**, a multi-stage framework for Temporal Knowledge Graph Question Answering. TempReasoner follows a staged processing pipeline in which each step is handled by a specialized module. The framework consists of five core modules: *Decompose*, *Retrieve*, *Filter*, *Generate*, and *Validate*, each responsible for a distinct stage in the reasoning workflow. This staged design allows the LLM to focus more effectively on each subtask, thereby reducing hallucinations that often arise when directly generating answers to complex temporal questions (Zhang and Zhang, 2025).

We summarize our contributions as follows:

- We introduce **TempReasoner**, a multi-stage temporal reasoning framework for TKGQA that unifies question decomposition, temporally aligned retrieval, alignment-based filtering, controlled answer generation, and answer

validation.

- We develop a temporal-sensitive and structure-aware **TKG-Retriever** trained with temporally informed hard negatives. This enables more accurate alignment between question semantics, temporal constraints, and TKG facts in the TKGQA task.
- Experiments on MultiTQ (Chen et al., 2023) and CronQuestions (Saxena et al., 2021) demonstrate that TempReasoner consistently outperforms strong baselines across multiple evaluation dimensions, with particularly notable gains on complex temporal questions.

2 Related Work

Existing approaches to temporal knowledge graph question answering (TKGQA) can be broadly grouped into **semantic similarity-based**, **semantic parsing-based**, and **retrieval-augmented generation** methods.

Semantic similarity-based methods embed questions and temporal facts into a shared representation space and perform reasoning via similarity scoring. Representative works include ExaQT (Jia et al., 2021), LGQA (Liu et al., 2023), TwiRGCN (Sharma et al., 2023), as well as language-model-based approaches such as TempoQR (Mavromatis et al., 2022) and MultiQA (Chen et al., 2023). While effective for simple queries, these methods often struggle with complex multi-hop temporal reasoning due to the lack of explicit stepwise inference.

Semantic parsing-based methods convert questions into executable logical forms defined over temporal operators, exemplified by TEQUILA (Jia et al., 2018), SF-TQA (Ding et al., 2022), and Prog-TQA (Chen et al., 2024b). Although interpretable, these approaches are sensitive to parsing errors and prone to cascading failures in complex scenarios.

Retrieval-augmented generation methods retrieve relevant temporal facts or subgraphs and rely on LLMs for reasoning and answer generation, including GenTKGQA (Gao et al., 2024), TimeR4 (Qian et al., 2024), ARI (Chen et al., 2024c), TempAgent (QianyiHu et al., 2025), and RTQA (Gong et al., 2025). Most existing methods couple retrieval and generation within a single reasoning loop, with limited mechanisms for explicitly

¹<https://platform.openai.com/docs/models/gpt-4o-mini>

²<https://github.com/deepseek-ai/DeepSeek-V3>

³For example: "Which country was the first/last that China wanted to cooperate with?"

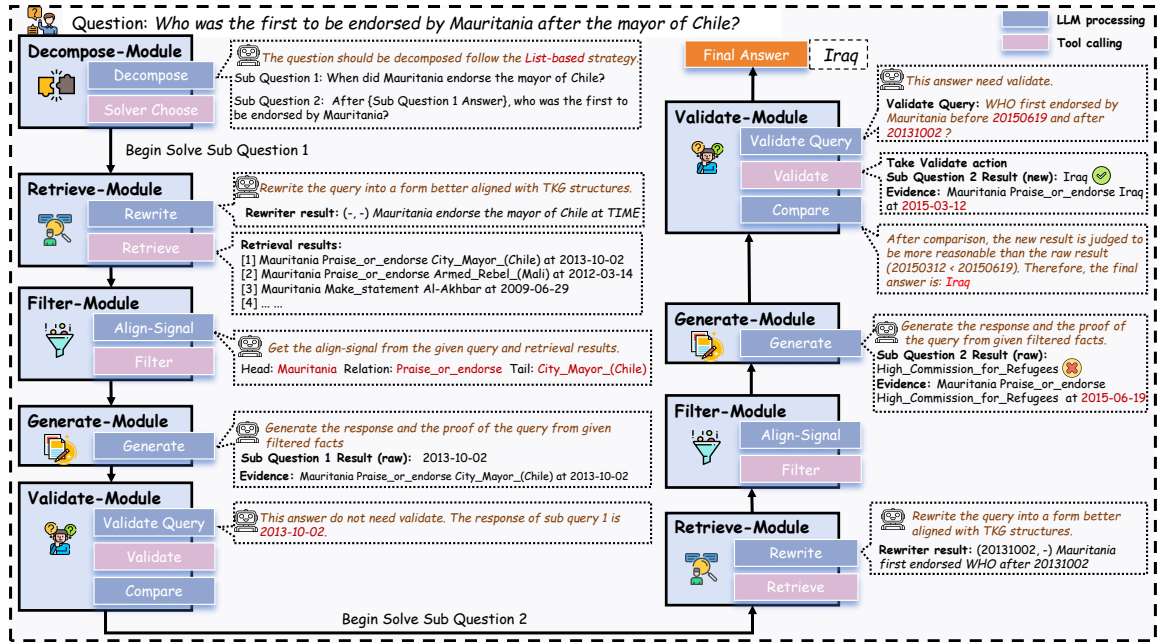


Figure 2: Overview of the *TempReasoner* framework. The figure illustrates this complete process on the example question “Who was the first to be endorsed by Mauritania after the mayor of Chile?”.

validating or correcting erroneous intermediate results. In contrast, our work explicitly decomposes temporal questions and introduces dedicated filtering and validation stages to improve reasoning reliability.

3 Method

3.1 Overview

We propose *TempReasoner*, a multi-stage framework that performs TKGQA through structured decomposition, retrieval, filtering, generation, and validation. As illustrated in Figure 2, *TempReasoner* treats the QA task as a multi-stage reasoning process to produce a reliable answer.

The *Decompose-Module* (Section 3.2) selects a List- or Tree-based strategy based on the question type and decomposes the question into sub-questions, whose results are later combined according to the selected strategy. The *Retrieve-Module* (Section 3.3) rewrites the question into a TKG-compatible query format and retrieves the top- K temporal facts using a trained retriever. The *Filter-Module* (Section 3.4) constructs an alignment signal to filter noisy or irrelevant evidence. The *Generate-Module* (Section 3.5) produces a candidate answer based on the filtered facts. Finally, the *Validate-Module* (Section 3.6) generates a validation query to assess the correctness of the candidate answer and selects the final answer based on

validation consistency, thereby improving answer reliability. A detailed illustration of the full algorithmic pipeline can be found in Appendix A.1 and we provide task preliminaries in Appendix A.2.

Each module in *TempReasoner* follows a unified design pattern based on LLM-driven reasoning and tool-calling. Details of the module architectures are given in Appendix A.3.

3.2 Question-structure-based Decomposition

The *Decompose-Module* is designed to decompose temporal questions into structured sub-questions. We adopt two representative decomposition strategies according to question structure: a List-based strategy and a Tree-based strategy. The List-based strategy handles sequential queries that require stepwise inference. For example, in the query “Who was the first to be endorsed by Mauritania after the mayor of Chile?”, one must first determine when Mauritania endorsed the mayor of Chile, and only then answer who was the first after that timestamp. The Tree-based strategy targets queries that require retrieving multiple partial results in parallel before combining them. For instance, the query “Which country did China wish to meet with in the same month as Japan?” requires retrieving whom China wished to meet and whom Japan wished to meet, and then identifying the entity that satisfies the same-month temporal constraint across both results. We formalize the decomposition process

Table 1: Example of sample construction in *TKG-Retriever*. The corresponding TKG fact is (*Agence_France-Presse, Praise_or_endorse, Iraq, 2014-10-01*). Three query templates are used: the raw-query template, the last-query template, and the first-query template. Five negative-sample templates are constructed: *inv_neg*, *sim_ent_neg*, *sim_rel_neg*, *temp_last_neg*, and *temp_first_neg*. All examples in the table place WHO in the subject position, and analogous templates can be constructed by placing WHO in the object position.

query-template	query example	neg-template	neg example
<i>raw-query</i>		<i>inv_neg</i>	
WHO {r} {o} at {t} ?	WHO Praise or endorse Iraq at 2014-10-01 ?	{o} {r} {s} at {t}	Iraq Praise_or_endorse Agence_France-Presse at 2014-10-01
		<i>sim_ent_neg</i>	
		{s _{sim} } {r} {o} at {t}	Government_(France) Praise_or_endorse Iraq at 2014-10-01
		<i>sim_rel_neg</i>	
		{s} {r _{sim} } {o} at {t}	Agence_France-Presse Make_optimistic_comment Iraq at 2014-10-01
<i>last-query</i>		<i>temp_last_neg</i>	
WHO last {r} {o} before {t+1} ?	WHO last Praise or endorse Iraq before 2014-10-02 ?	{s _{sim} } {r} {o} at {< t}	Government_(France) Praise_or_endorse Iraq at 2014-09-20
<i>first-query</i>		<i>temp_first_neg</i>	
WHO first {r} {o} after {t-1} ?	WHO first Praise or endorse Iraq after 2014-9-30 ?	{s _{sim} } {r} {o} at {> t}	Civil_Service_(France) Praise_or_endorse Iraq 2014-10-07

as:

$$(q_1, q_2, \dots) = \text{DecomposeModule}(Q, \mathcal{P}, S),$$

$$S \in \{\text{List-based}, \text{Tree-based}\}, \quad (1)$$

where Q denotes the original input question, \mathcal{P} represents the instruction prompt used to guide the LLM-based decomposition process, S indicates the decomposition strategy, and q_i is the i -th sub-question produced by the *Decompose-Module*.

3.3 TKG-aware Retrieval

The *Retrieve-Module* treats the question as a query and retrieves relevant facts from the TKG. General-purpose retrievers such as BGE-M3 (Chen et al., 2024a) often misalign entities or relations and fail to recognize temporal constraints embedded in queries. To mitigate these limitations, we design a retrieval pipeline that explicitly accounts for the structural and temporal characteristics of TKGs.

Structure-aligned rewrite. Each question is rewritten into a format that is aligned with the canonical quadruple structure used in temporal knowledge graphs. The rewritten query consists of (1) a temporal-constraint component that extracts explicit time bounds and restricts retrieval accordingly, and (2) a textual component where TIME and WHO denote unspecified timestamps and entities. Temporal operators such as *first/last* and *before/after* are preserved to guide later retrieval. For example, the question “Who was the first that China wanted to negotiate with after 20141010?” becomes “(20141010, -) China first want_to_negotiate_with WHO after 20141010,” where (20141010, -) restricts retrieval to facts occurring strictly after the specified date. This step facilitates better alignment between the question and the TKG storage format, thereby reducing semantic drift.

TKG-Retriever. To further enhance retrieval accuracy, we train a *TKG-Retriever* that incorporates semantic similarity and temporal order. We construct harder negative pairs that enforce fine-grained semantic and temporal discrimination instead of random negatives such as (Qian et al., 2024). For each entity e and relation r , we gather sets of semantically similar entities e_{sim} and relations r_{sim} . We introduce a residual quantization method for efficiently identifying semantically similar entities; details are provided in Appendix A.4. Given a fact (s, r, o, t) , its query is treated as a positive instance, and multiple hard negatives are constructed, including: (i) *inv_neg*: swapping subject and object, (ii) *sim_ent_neg*: replacing the subject with an entity from e_{sim} , (iii) *sim_rel_neg*: replacing the relation with one from r_{sim} .

To encode temporal sensitivity, we extend the raw-query template to *last-query* and *first-query* variants by injecting operators (*last/before* or *first/after*) and shifting timestamps to $t+1$ or $t-1$ ⁴. Corresponding temporal negatives (*temp_last_neg*, *temp_first_neg*) substitute the subject with a similar entity and assign timestamps $< t$ or $> t$. These negatives encourage the retriever to jointly model entity identity, relation semantics, and chronological ordering.

For each fact, we construct N_{neg} hard negatives, and treat other samples within the batch as in-batch negatives. The *TKG-Retriever* is trained using an InfoNCE loss \mathcal{L} :

$$Z = \exp(\text{sim}(q, p)/\tau) + \exp(\text{sim}(q, n_{\text{hard}})/\tau)$$

$$+ \sum_{b \in \mathcal{N}_{\text{batch}}} \exp(\text{sim}(q, b)/\tau),$$

$$\mathcal{L} = -\log \frac{\exp(\text{sim}(q, p)/\tau)}{Z}, \quad (2)$$

⁴e.g., 2014-10-01 + 1 \rightarrow 2014-10-02

where q denotes the embedding of the rewritten query, p is the embedding of its corresponding positive TKG fact, n_{hard} denotes one of the constructed hard negative samples, and $\mathcal{N}_{\text{batch}}$ denotes the set of in-batch negative samples within the same training batch. The function $\text{sim}(\cdot, \cdot)$ measures the similarity between two embeddings (e.g., cosine similarity), and τ is the temperature coefficient.

For each rewritten query, we apply the trained *TKG-Retriever* to obtain facts aligned with both its semantic content and temporal constraints. The retrieved evidence is then passed to the *Filter-Module* for further noise reduction.

3.4 Alignment-based Filtering

The *Filter-Module* is designed to mitigate noise in the facts retrieved by the retriever. Prior approaches either directly feed raw retrieved facts into the LLM (Gong et al., 2025) or apply a reranking step (Qian et al., 2024). However, these methods often lack explicit mechanisms to reliably filter irrelevant evidence, which may still mislead the LLM during reasoning.

To address this issue, the *Filter-Module* constructs an *alignment signal* from the retrieved evidence and applies it to filter noisy facts, producing a denoised evidence set. As illustrated in Figure 2, consider the sub-question #1: “When did Mauritania endorse the mayor of Chile?”. In this case, the question implies a target relation corresponding to *Praise_or_endorse*, with *Mauritania* as the subject. However, the raw retrieved set contains semantically related but irrelevant facts, such as (*Mauritania*, *Make_statement*, *Al-Akhbar*, 2009-06-29).

The *Filter-Module* generates an alignment signal based on the top- n retrieved facts and uses this signal to exclude such noisy evidence. One example of a generated alignment signal for this sub-question is shown below:

Alignment.Head: Mauritania
 Alignment.Relation: Praise_or_endorse
 Alignment.Tail: ?

After filtering, only the most relevant fact (*Mauritania*, *Praise_or_endorse*, *City_Mayor_(Chile)*, 2013-10-02) is retained, providing cleaner and more reliable input for subsequent answer generation. This filtering step reduces the impact of irrelevant evidence and helps alleviate error propagation and hallucinations in later stages. Additional imple-

mentation details of the *Filter-Module* are provided in Appendix A.5.

3.5 Controlled Answer Generation

The *Generate-Module* is responsible for producing a candidate answer for each sub-question after retrieval and filtering. Given a sub-question q_i produced by the *Decompose-Module* and a filtered fact set \mathcal{F}_i returned by the *Filter-Module*, the *Generate-Module* generates an answer together with a supporting *evidence set*, which consists of a subset of facts from \mathcal{F}_i that the model identifies as evidence for the generated answer.

We formalize the behavior of the *Generate-Module* as the following mapping:

$$(a_i, \hat{f}_i) = \text{GenerateModule}(q_i, \mathcal{F}_i, \Theta_{\text{LLM}}), \quad (3)$$

where q_i denotes the i -th sub-question, $\mathcal{F}_i = \{f_1, f_2, \dots, f_m\}$ is the filtered fact set aligned with q_i , and Θ_{LLM} represents the parameters of the LLM used for answer generation. The output a_i is a raw candidate answer for the sub-question, and $\hat{f}_i \subseteq \mathcal{F}_i$ denotes the set of supporting facts selected by the model as evidence for a_i .

Although the *Generate-Module* produces an answer candidate, it does not guarantee correctness or optimality. Instead, answer reliability is achieved through the upstream constraints imposed by retrieval and filtering, and is further refined by the downstream *Validate-Module*, which explicitly verifies and revises uncertain predictions.

3.6 Answer Validation and Self-Correction

The *Validate-Module* assesses whether a candidate answer a_i generated for a sub-question q_i requires further verification. Given q_i , its generated answer a_i , and the associated supporting facts \hat{f}_i , the module first examines whether a_i is sufficiently supported by \hat{f}_i and satisfies the temporal and relational constraints specified in q_i . If the retrieved evidence is consistent with the query constraints, no additional validation is performed.

When the support for a_i is insufficient or ambiguous, the *Validate-Module* constructs a validation query that explicitly re-examines the temporal and relational conditions implied by the candidate answer. Importantly, this validation query is designed to *challenge* the current answer by imposing stricter or counterfactual temporal constraints, rather than reinforcing the original prediction.

As illustrated in Figure 2, for the sub-question “Mauritania first endorsed WHO after 20131002”,

the module formulates a validation query based on the raw answer: “*WHO was first endorsed by Mauritania after 20131002 and before 20150619?*”. This query is processed through the same Retrieve–Filter–Generate pipeline to obtain an alternative candidate answer \hat{a}_i .

The final answer is selected as:

$$a_i^* = \text{Select}(a_i, \hat{a}_i), \quad (4)$$

where the selection function prefers the answer that is more strongly supported by retrieved evidence and better satisfies the imposed temporal constraints.

Overall, the *Validate-Module* serves as a second-pass temporal verifier that explicitly re-evaluates uncertain predictions under stricter constraints, enabling self-correction and improving answer reliability without assuming the initial generation to be correct.

4 Experiments

4.1 Datasets and Baselines

We evaluate TempReasoner on two widely used TKGQA benchmarks: MultiTQ (Chen et al., 2023) and CronQuestions (Saxena et al., 2021). Dataset statistics are provided in Appendix A.6.

We select a set of representative baselines for comparison, which can be broadly categorized into three groups: (1) **PLM-based methods**, including BERT(Kenton and Toutanova, 2019), ALBERT(Lan et al.), and DistilBERT(Sanh, 2019). (2) **Embedding-based methods**, such as EmbedKGQA(Saxena et al., 2020), CronKGQA(Saxena et al., 2021), and MultiQA(Chen et al., 2023). (3) **LLM-based methods**, which can be further divided into two subcategories: models that directly rely on LLMs to produce answers, including LLaMA2 and ChatGPT and models that incorporate knowledge graph information into LLM reasoning, including KG-RAG, CoT(Wei et al., 2022), ReAct(Yao et al., 2022), ARI(Chen et al., 2024c), RTQA(Gong et al., 2025), and TempAgent(QianyiHu et al., 2025). For several baselines, we report results following (Chen et al., 2024c) and (Gong et al., 2025).

4.2 Implementation Details

For training the TKG-Retriever, we fine-tune BGE-M3 (Chen et al., 2024a) as the base encoder. The model is trained on two NVIDIA A800 GPUs (80GB memory each), after which we build the

retrieval index using FAISS Library (Johnson et al., 2019) for efficient similarity search. For each quadruple, we construct $N_{\text{neg}} = 10$ negatives: five with WHO in the subject position and five in the object position.

During evaluation, we employ the DeepSeek-V3 model via the DeepSeek API⁵ as the backbone LLM for all agents. The retrieval context size is set to 100 facts, while the *Filter-Module* uses the top 20 retrieved facts to construct the filter alignment signals. We follow the data preprocessing as (Chen et al., 2024c). Model performance is evaluated using Hits@1, which measures the percentage of questions whose top-ranked prediction matches the ground-truth answer.

4.3 Main Results

Model	Overall	Question Type		Answer Type	
		Multiple	Single	Entity	Time
BERT	0.083	0.061	0.092	0.101	0.040
DistilBERT	0.083	0.074	0.087	0.102	0.037
ALBERT	0.108	0.086	0.116	0.139	0.032
EmbedKGQA	0.206	0.134	0.235	0.290	0.001
CronKGQA	0.279	0.134	0.337	0.328	0.156
MultiQA	0.293	0.159	0.347	0.349	0.157
LLaMA2	0.185	0.101	0.220	0.239	0.055
ChatGPT	0.102	0.077	0.147	0.137	0.020
ARI	0.380	0.210	0.680	0.394	0.344
TimeR4	0.728	0.335	<u>0.887</u>	0.639	0.945
TempAgent	0.702	0.316	0.857	0.624	0.870
RTQA	<u>0.765</u>	<u>0.424</u>	0.902	<u>0.692</u>	<u>0.942</u>
TempReasoner	0.805	0.631	0.875	0.765	0.899

Table 2: Performance of Hits@1 on the MultiTQ dataset. **Bold** numbers denote the best performance, and underlined numbers denote the second-best performance in each column.

We present the comparative results of TempReasoner against baseline models on the MultiTQ and CronQuestions datasets in Table 2 and Table 3 respectively. Specifically, TempReasoner achieves consistently strong results across both MultiTQ and CronQuestions.

On **MultiTQ**, TempReasoner obtains an Overall Hits@1 of 0.805, outperforming the strongest prior baseline RTQA. When examining question types, TempReasoner achieves 0.631 on Multiple questions, with the latter representing a large improvement of 20.7 percentage points over the previous best model. Multi-hop temporal queries often require chaining multiple temporal conditions or

⁵DeepSeek-V3-250324

Model	Overall	Question Type		Answer Type	
		Simple	Complex	Entity	Time
BERT	0.243	0.249	0.239	0.277	0.179
ALBERT	0.248	0.255	0.235	0.279	0.177
EmbedKGQA	0.288	0.290	0.286	0.411	0.056
CronKGQA	0.647	0.987	0.392	0.699	0.549
ChatGPT	0.249	0.250	0.247	0.246	0.253
KG-RAG	0.490	0.460	0.518	0.470	0.520
CoT	0.640	0.690	0.610	0.620	0.660
ReAct	0.685	0.835	0.525	0.650	0.755
ARI	0.707	0.860	0.570	0.660	0.800
TempAgent	<u>0.842</u>	<u>0.895</u>	<u>0.640</u>	<u>0.805</u>	0.921
TempReasoner	0.850	0.919	0.799	0.846	<u>0.857</u>

Table 3: Performance of Hits@1 on the CronQuestions. **Bold** numbers denote the best performance, and underlined numbers denote the second-best performance in each column.

selecting entities under strict ordering constraints, and TempReasoner’s improvements in this category highlight its enhanced reasoning capability. In terms of answer types, TempReasoner reaches 0.765 on Entity questions and 0.899 on Time questions. The improvement of 7.3 percentage points on Entity answers is notable, as most multi-hop temporal questions require identifying the correct entity.

On **CronQuestions**, TempReasoner also achieves strong overall performance, with an Overall Hits@1 of 0.850, outperforming all baseline models. TempReasoner performs consistently well across these settings, reaching 0.919 on Simple questions and 0.799 on Complex questions. In particular, its performance on Complex questions shows a clear advantage over prior methods, exceeding the next-best model (TempAgent) by nearly 16 percentage points. This result indicates that TempReasoner is more effective at handling multi-hop reasoning under temporal constraints. For answer types, TempReasoner achieves Hits@1 scores of 0.846 on Entity questions and 0.857 on Time questions, again demonstrating strong improvements over the baselines.

In addition, TempReasoner performs slightly worse than some baselines, particularly for questions with time as the final answer. This performance gap is mainly attributed to occasional errors in temporal expression generation by the LLM. Even when the retrieved evidence is correct, the model may produce natural language expressions (e.g., *January*) instead of the required normalized

time format (e.g., *2014-01*), which leads to mismatches under strict evaluation.

4.4 Ablation Study

We further conduct ablation studies to evaluate the contribution of each module in our framework. Since the MultiTQ test set is extremely large, we randomly sample 1,000 questions as the ablation benchmark and rerun our base model. The results are summarized in Table 4.

Effect of the Decompose-Module. The *w/o Dec.* variant represents a setting in which questions are processed in their original form, without decomposition. The performance drops by 17.9% on the overall score and exhibits a substantial 43.2% decrease on complex *Multiple* questions. This demonstrates that the decomposition strategies employed by the *Decompose-Module* are crucial for solving multi-hop temporal reasoning problems.

Effect of the Retrieve-Module and Filter-Module. The *w/o Ret.* variant retrieves context using the BGE-M3 encoder with FAISS instead of proposed *Retrieve-Module + Filter-Module*. This variant shows consistent degradation across all metrics, with an overall drop of 30.1%, highlighting that our structure-aligned retrieval combined with temporal filtering is significantly more effective at extracting relevant temporal facts from the TKG.

Effect of the Validate-Module. The *w/o Val.* variant represents a configuration in which the LLM-generated answers are used directly without verification. This leads to an 8.3% decrease on the overall score. Notably, compared with the moderate drop on *equal* questions (-4.3%), performance decreases much more sharply on stricter temporal constraints such as *after_first* (-20.9%) and *before_last* (-13%). These results indicate that the *Validate-Module* effectively identifies and corrects erroneous LLM outputs, especially under fine-grained temporal reasoning conditions.

5 Further Analysis

To better understand the behavior and generalization properties of TempReasoner, we conduct a series of further analyses from both the model and retrieval perspectives. In addition, we provide a case study of representative questions in Appendix A.7.

	Overall	answer type		qtype						qlabel	
		entity	time	after_first	before/after	before_last	first/last	equal	eq_multi	multiple	single
TempReasoner	0.823	0.793	0.893	0.583	0.850	0.678	0.913	0.909	0.709	0.646	0.894
<i>w/o Dec.</i>	0.644	0.548	0.863	0.261	0.670	0.104	<u>0.903</u>	0.853	0.345	0.214	0.815
<i>w/o Ret.</i>	0.522	0.447	0.694	0.304	0.450	0.391	0.636	0.678	0.200	0.319	0.628
<i>w/o Val.</i>	<u>0.740</u>	<u>0.679</u>	<u>0.879</u>	<u>0.374</u>	<u>0.735</u>	<u>0.548</u>	0.887	<u>0.866</u>	<u>0.673</u>	<u>0.502</u>	<u>0.835</u>

Table 4: Ablation study on the MultiTQ dataset. **Bold** numbers denote the best performance; underlined numbers denote the second-best performance in each column.

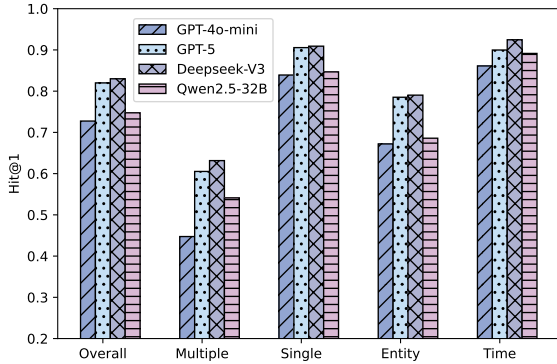


Figure 3: Performance comparison across different LLMs on four question types.

5.1 Performance of different LLMs

To further examine the adaptability of TempReasoner, we evaluated our framework using four LLMs ranging across different model scales: DeepSeek-V3, GPT-4o-mini, GPT-5, and Qwen2.5-32B. The results are shown in Figure 3 indicate that models with stronger reasoning capabilities generally achieve higher performance across all task types. Nevertheless, TempReasoner still maintains strong performance even when paired with comparatively weaker models such as GPT-4o-mini and Qwen2.5-32B. This suggests that our framework is adaptable and can generalize across different LLM backbones.

5.2 Performance of different Retrievers

To evaluate the effectiveness of the proposed TKG-Retriever, we compare it with two widely used text retrievers, BGE-M3 and Sentence-BERT (Reimers and Gurevych, 2019). All retrievers encode TKG facts and build FAISS indexes over their vector representations, and candidate facts are retrieved from the corresponding index for each query. Retrieval effectiveness is measured by the average rank of the ground-truth fact in the retrieved list, where a lower rank indicates better retrieval quality.

As shown in Figure 4, the TKG-Retriever con-

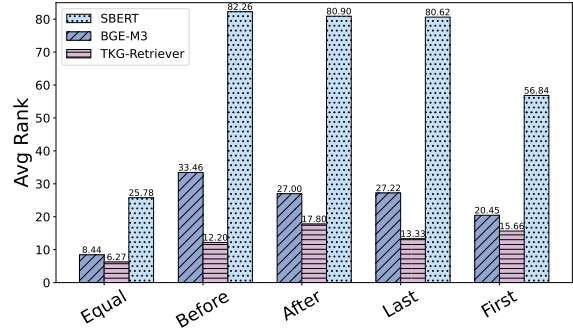


Figure 4: Comparison of retrieval performance between different retrievers.

sistently outperforms both BGE-M3 and Sentence-BERT across all evaluation settings. The performance gap is particularly pronounced for queries involving temporal constraints such as *before/after* and *last/first*, where sensitivity to temporal ordering is crucial. These results demonstrate that the proposed TKG-Retriever achieves more reliable temporal fact retrieval on TKGs.

6 Conclusion

We presented **TempReasoner**, a multi-stage framework for TKGQA that integrates question decomposition, TKG-aware retrieval, alignment-based filtering, controlled answer generation, and verification-based validation. By explicitly structuring the temporal reasoning process, TempReasoner addresses the limitations of general-purpose retrievers and mitigates temporal hallucination issues in LLMs. Experiments on MultiTQ and CronQuestions demonstrate consistent improvements across all evaluation settings, with particularly strong gains on complex multi-hop temporal queries. Additional analyses further demonstrate the effectiveness of the proposed TKG-Retriever and the framework’s ability to generalize across different LLM backbones. Overall, TempReasoner provides a principled approach for combining structured temporal reasoning with LLMs for TKGQA.

590 Limitations

591 While TempReasoner demonstrates strong perfor-
592 mance on complex temporal question answering, it
593 still has several limitations.

594 First, the proposed framework primarily operates
595 at the level of individual temporal facts (quadru-
596 ples) during retrieval and reasoning. Although the
597 multi-stage design effectively mitigates noise and
598 temporal misalignment, it does not explicitly ex-
599 ploit the global topological structure of temporal
600 knowledge graphs, such as multi-hop connectivity
601 patterns or path-level dependencies. As a result,
602 certain forms of structural reasoning that rely on ex-
603 plicit graph traversal or subgraph-level aggregation
604 are not directly modeled.

605 Second, TempReasoner relies on LLMs to per-
606 form key intermediate steps, including question de-
607 composition and alignment signal induction. While
608 the downstream filtering and validation modules
609 substantially reduce error propagation, imperfect
610 decomposition or signal generation may still affect
611 performance in some cases, particularly for ques-
612 tions with highly implicit temporal constraints.

613 In many real-world applications, temporally
614 grounded knowledge is primarily stored in un-
615 structured documents rather than explicit tempo-
616 ral knowledge graph form, and such data often
617 contains substantial noise. An important direction
618 for future work is to efficiently construct temporal
619 knowledge graphs from noisy temporal documents
620 and perform retrieval and reasoning over the result-
621 ing structured graphs.

622 Ethics Statement

623 All steps and data described in our paper follow the
624 ACL Ethics Policy⁶.

625 References

626 Jianlv Chen, Shitao Xiao, Peitian Zhang, Kun Luo, Defu
627 Lian, and Zheng Liu. 2024a. Bge m3-embedding:
628 Multi-lingual, multi-functionality, multi-granularity
629 text embeddings through self-knowledge distillation.
630 *arXiv preprint arXiv:2402.03216*.

631 Zhuo Chen, Zhao Zhang, Zixuan Li, Fei Wang, Yutao
632 Zeng, Xiaolong Jin, and Yongjun Xu. 2024b. Self-
633 improvement programming for temporal knowledge
634 graph question answering. In *Proceedings of the*

*2024 Joint International Conference on Computa- 635
tional Linguistics, Language Resources and Evalua- 636
tion (LREC-COLING 2024)*, pages 14579–14594. 637

Ziyang Chen, Dongfang Li, Xiang Zhao, Baotian Hu, 638
and Min Zhang. 2024c. Temporal knowledge ques- 639
tion answering via abstract reasoning induction. In 640
*Proceedings of the 62nd Annual Meeting of the As- 641
sociation for Computational Linguistics (Volume 1: 642
Long Papers)*, pages 4872–4889. 643

Ziyang Chen, Jinzhi Liao, and Xiang Zhao. 2023. Multi- 644
granularity temporal question answering over knowl- 645
edge graphs. In *Proceedings of the 61st Annual Meet- 646
ing of the Association for Computational Linguistics 647
(Volume 1: Long Papers)*, pages 11378–11392. 648

Wentao Ding, Hao Chen, Huayu Li, and Yuzhong Qu. 649
2022. Semantic framework based query generation 650
for temporal question answering over knowledge 651
graphs. In *Proceedings of the 2022 Conference on 652
Empirical Methods in Natural Language Processing*, 653
pages 1867–1877. 654

Yifu Gao, Linbo Qiao, Zhigang Kan, Zhihua Wen, 655
Yongquan He, and Dongsheng Li. 2024. Two-stage 656
generative question answering on temporal knowl- 657
edge graph using large language models. In *Findings 658
of the Association for Computational Linguistics ACL 659
2024*, pages 6719–6734. 660

Zhaoyan Gong, Juan Li, Zhiqiang Liu, Lei Liang, Hua- 661
jun Chen, and Wen Zhang. 2025. Rtqa: Recur- 662
sive thinking for complex temporal knowledge graph 663
question answering with large language models. In 664
*Proceedings of the 2025 Conference on Empirical 665
Methods in Natural Language Processing*, pages 666
9864–9881. 667

Zhen Jia, Abdalghani Abujabal, Rishiraj Saha Roy, Jan- 668
nik Strötgen, and Gerhard Weikum. 2018. Tequila: 669
Temporal question answering over knowledge bases. 670
In *Proceedings of the 27th ACM international confer- 671
ence on information and knowledge management*, 672
pages 1807–1810. 673

Zhen Jia, Soumajit Pramanik, Rishiraj Saha Roy, and 674
Gerhard Weikum. 2021. Complex temporal question 675
answering on knowledge graphs. In *Proceedings of 676
the 30th ACM international conference on informa- 677
tion & knowledge management*, pages 792–802. 678

Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2019. 679
Billion-scale similarity search with gpus. *IEEE 680
Transactions on Big Data*, 7(3):535–547. 681

Jacob Devlin Ming-Wei Chang Kenton and Lee Kristina 682
Toutanova. 2019. Bert: Pre-training of deep bidirec- 683
tional transformers for language understanding. In 684
Proceedings of NAACL-HLT, pages 4171–4186. 685

Zhenzhong Lan, Mingda Chen, Sebastian Goodman, 686
Kevin Gimpel, Piyush Sharma, and Radu Soricut. 687
Albert: A lite bert for self-supervised learning of lan- 688
guage representations. In *International Conference 689
on Learning Representations*. 690

⁶<https://www.aclweb.org/portal/content/acl-code-ethics>

691	Yonghao Liu, Di Liang, Mengyu Li, Fausto Giunchiglia,	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten	746
692	Ximing Li, Sirui Wang, Wei Wu, Lan Huang, Xi-	Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou,	747
693	aoyue Feng, and Renchu Guan. 2023. Local and	et al. 2022. Chain-of-thought prompting elicits rea-	748
694	global: Temporal question answering via information	soning in large language models. <i>Advances in neural</i>	749
695	fusion. In <i>IJCAI</i> , pages 5141–5149.	<i>information processing systems</i> , 35:24824–24837.	750
696	Costas Mavromatis, Prasanna Lakkur Subramanyam,	Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak	751
697	Vassilis N Ioannidis, Adesoji Adeshina, Phillip R	Shafraan, Karthik R Narasimhan, and Yuan Cao. 2022.	752
698	Howard, Tetiana Grinberg, Nagib Hakim, and George	React: Synergizing reasoning and acting in language	753
699	Karypis. 2022. Tempoqr: temporal question reason-	models. In <i>The eleventh international conference on</i>	754
700	ing over knowledge graphs. In <i>Proceedings of the</i>	<i>learning representations</i> .	755
701	<i>AAAI conference on artificial intelligence</i> , volume 36,	Wan Zhang and Jing Zhang. 2025. Hallucination mitiga-	756
702	pages 5825–5833.	tion for retrieval-augmented large language models:	757
703	Xinying Qian, Ying Zhang, Yu Zhao, Baohang Zhou,	a review. <i>Mathematics</i> , 13(5):856.	758
704	Xuhui Sui, Li Zhang, and Kehui Song. 2024. Timer4:		
705	Time-aware retrieval-augmented large language mod-		
706	els for temporal knowledge graph question answering.		
707	In <i>Proceedings of the 2024 Conference on Empiri-</i>		
708	<i>cal Methods in Natural Language Processing</i> , pages		
709	6942–6952.		
710	QianyiHu QianyiHu, Xinhui Tu, Guo Cong, and Shun-		
711	ping Zhang. 2025. Time-aware react agent for tem-		
712	poral knowledge graph question answering. In <i>Find-</i>		
713	<i>ings of the Association for Computational Linguistics:</i>		
714	<i>NAACL 2025</i> , pages 6013–6024.		
715	Nils Reimers and Iryna Gurevych. 2019. Sentence-bert:		
716	Sentence embeddings using siamese bert-networks.		
717	In <i>Proceedings of the 2019 Conference on Empirical</i>		
718	<i>Methods in Natural Language Processing and the 9th</i>		
719	<i>International Joint Conference on Natural Language</i>		
720	<i>Processing (EMNLP-IJCNLP)</i> , pages 3982–3992.		
721	V Sanh. 2019. Distilbert, a distilled version of bert:		
722	smaller, faster, cheaper and lighter. In <i>Proceedings</i>		
723	<i>of Thirty-third Conference on Neural Information</i>		
724	<i>Processing Systems (NIPS2019)</i> .		
725	Apoorv Saxena, Soumen Chakrabarti, and Partha Taluk-		
726	dar. 2021. Question answering over temporal knowl-		
727	edge graphs. In <i>Proceedings of the 59th Annual</i>		
728	<i>Meeting of the Association for Computational Lin-</i>		
729	<i>guistics and the 11th International Joint Conference</i>		
730	<i>on Natural Language Processing (Volume 1: Long</i>		
731	<i>Papers)</i> , pages 6663–6676.		
732	Apoorv Saxena, Aditay Tripathi, and Partha Talukdar.		
733	2020. Improving multi-hop question answering over		
734	knowledge graphs using knowledge base embeddings.		
735	In <i>Proceedings of the 58th annual meeting of the as-</i>		
736	<i>sociation for computational linguistics</i> , pages 4498–		
737	4507.		
738	Aditya Sharma, Apoorv Saxena, Chitranshu Gupta,		
739	Mehran Kazemi, Partha Talukdar, and Soumen		
740	Chakrabarti. 2023. Twirgcn: Temporally weighted		
741	graph convolution for question answering over tem-		
742	poral knowledge graphs. In <i>Proceedings of the 17th</i>		
743	<i>Conference of the European Chapter of the Asso-</i>		
744	<i>ciation for Computational Linguistics</i> , pages 2049–		
745	2060.		

A Appendix

A.1 Algorithm of TempReasoner

Algorithm 1: TempReasoner: A Multi-Stage Framework for TKGQA

Input: Temporal Knowledge Graph \mathcal{G} , query Q
Output: Final answer A

- 1 **Step 1: Query Decomposition;**
- 2 Obtain instruction prompt \mathcal{P} ;
- 3 $(q_1, q_2, \dots, q_n), S = \text{Decompose-Module}(Q, \mathcal{P});$
// $S \in \{\text{List-based}, \text{Tree-based}\}$
- 4 **Step 2: Sub-query Solving;**
- 5 **for** $i = 1$ **to** n **do**
 - // (a) Temporal-aware retrieval
 - 6 $\mathcal{R}_i = \text{Retrieve-Module}(q_i, \mathcal{G});$
 - // (b) Structure-guided filtering
 - 7 $\mathcal{F}_i = \text{Filter-Module}(q_i, \mathcal{R}_i);$
 - // (c) LLM-based answer generation
 - 8 $(a_i, \hat{f}_i) = \text{Generate-Module}(q_i, \mathcal{F}_i);$
// $a_i = \text{answer}, \hat{f}_i = \text{supporting}$
// evidence set
 - // (d) Optional validation
 - 9 **if** $\text{IsReliable}(a_i, \hat{f}_i)$ **then**
 - 10 | $\hat{a}_i \leftarrow a_i;$
 - 11 **else**
 - 12 | $\hat{q}_i \leftarrow \text{ConstructValidateQuery}(Q, a_i, \hat{f}_i);$
 - 13 | $\hat{a}_i = \text{Validate-Module}(\hat{q}_i);$
- 14 **Step 3: Answer Aggregation;**
- 15 **if** S is List-based **then**
- 16 | $A = \text{SequentialCombine}(\hat{a}_1, \dots, \hat{a}_n);$
- 17 **else**
- 18 | $A = \text{ParallelMerge}(\{\hat{a}_i\}_{i=1}^n);$
- 19 **return** $A;$

A.2 Preliminary

A Temporal Knowledge Graph (TKG) is defined as $\mathcal{G} = \{(s, r, o, t)\} \subseteq \mathcal{E} \times \mathcal{R} \times \mathcal{E} \times \mathcal{T}$, where \mathcal{E} denotes the set of entities, \mathcal{R} the set of relations, and \mathcal{T} the set of timestamps. Each quadruple (s, r, o, t) represents a temporal fact, indicating that entity s is connected to entity o by relation r at time t .

Temporal Knowledge Graph Question Answering (TKGQA) aims to answer a temporal query Q by retrieving relevant facts from the TKG. Queries may range from simple one-hop questions, such as “When did Indonesia want to make a visit to China?”, to complex multi-hop questions, such as “Which country was the first that China wanted to cooperate with after Indonesia?”. The answers may involve either a timestamp or an entity.

A.3 Module Architecture

We construct all modules following a general pattern. Each module consists of two components: an

LLM processing component and a tool-calling component. The LLM processing component includes the task instruction, few-shot exemplars, input formatting, and output specification. The tool-calling component allows the LLM to invoke external tools depending on the functional requirement of the module. For example, the *Retrieve-Module* calls the retriever to obtain relevant temporal facts, while the *Filter-Module* invokes the filtering tool to refine the retrieved facts.

A.4 Semantic Similarity Dictionary via RQ-KMeans

To efficiently obtain semantically similar entities and relations in the TKG, we construct a similarity dictionary using a Residual Quantization K-Means (RQ-KMeans) framework. The procedure consists of three steps.

(1) Encoder-based representation. We first encode every entity and relation using a pretrained encoder $\mathcal{E}(\cdot)$. For an entity or relation x , its embedding is:

$$\mathbf{h}_x = \mathcal{E}(x) \in \mathbb{R}^d. \quad (5)$$

(2) Multi-level residual quantization. We apply RQ-KMeans to build a hierarchical codebook of size:

$$\ell_1 \times \ell_2 \times \ell_3,$$

where ℓ_i is the number of cluster centers at level i . For each level, residual quantization is performed as:

$$\begin{aligned} \mathbf{r}_x^{(1)} &= \mathbf{h}_x, \\ c_x^{(1)} &= \arg \min_{c \in \mathcal{C}_1} \|\mathbf{r}_x^{(1)} - \mathbf{c}\|_2, \\ \mathbf{r}_x^{(2)} &= \mathbf{r}_x^{(1)} - \mathbf{c}_x^{(1)}, \\ c_x^{(2)} &= \arg \min_{c \in \mathcal{C}_2} \|\mathbf{r}_x^{(2)} - \mathbf{c}\|_2, \\ \mathbf{r}_x^{(3)} &= \mathbf{r}_x^{(2)} - \mathbf{c}_x^{(2)}, \\ c_x^{(3)} &= \arg \min_{c \in \mathcal{C}_3} \|\mathbf{r}_x^{(3)} - \mathbf{c}\|_2, \end{aligned} \quad (6)$$

where \mathcal{C}_i denotes the i -th level K-Means codebook. Thus each entity or relation x is assigned a *semantic code*:

$$\text{sid}(x) = (c_x^{(1)}, c_x^{(2)}, c_x^{(3)}). \quad (7)$$

(3) Prefix-based semantic similarity. Two entities or relations x and y are considered semantically similar if their semantic codes share a prefix:

$$\text{prefix}_k(\text{sid}(x)) = \text{prefix}_k(\text{sid}(y)), \quad k \in \{1, 2, 3\}. \quad (8)$$

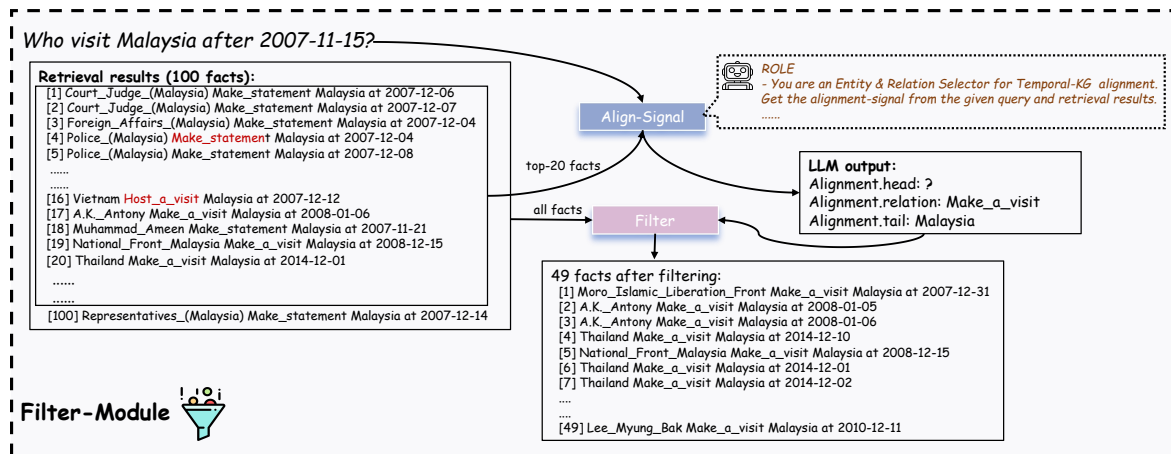


Figure 5: Illustration of *Filter-Module*.

A larger shared prefix (e.g., first two or all three levels) implies stronger semantic closeness.

This RQ-KMeans-based construction enables efficient lookup of semantically similar entities and relations, providing a scalable and accurate mechanism for supporting the hard-negative construction used by our TKG-Retriever.

A.5 Details of Filter-Module

Figure 5 provides a concrete example of how the *Filter-Module* operates for the question “Who visited Malaysia after 2007-11-15?”. The retriever first returns the top 100 facts whose textual surface matches portions of the query. However, many of these results correspond to semantically unrelated relations such as *Make_statement* or entities that only contain the token “Malaysia” in their names (e.g., *Court_Judge_(Malaysia)*, *Police_(Malaysia)*), rather than actual “visiting” events. This illustrates a common failure mode of TKG retrieval: high lexical similarity but low semantic alignment.

To refine this noisy candidate set, the *Filter-Module* first extracts the top- k retrieved facts and presents them to the LLM, asking it to infer an alignment signal. This signal captures the essential structural constraints that any valid supporting fact should satisfy. In this example, the LLM determines that the question requires facts matching the pattern: (*head*, *Make_a_visit*, *Malaysia*, *time*), where the head entity is unknown (null), the relation must correspond to the action “visit” (identified as *Make_a_visit*), and the tail must be *Malaysia*. This signal represents the LLM’s distilled understanding of which entities and relations are semantically compatible with the question in-

tent.

Once the alignment signal is obtained, the module applies it to *all* retrieved facts—including those beyond the top- k —to perform structured filtering. Any fact whose entity or relation does not match the inferred signal is removed. As shown in the figure, facts involving relations such as *Make_statement* or tail entities unrelated to Malaysia are eliminated.

After filtering, the candidate set is reduced from 100 facts to 49 high-quality facts, all of which correspond to genuine “visit Malaysia” events occurring after the specified date. This denoised evidence set preserves only semantically aligned temporal facts, enabling subsequent modules to reason over precise and relevant information rather than being distracted by irrelevant or misleading evidence.

This structured filtering process is crucial for mitigating semantic drift in retrieval and significantly improves the quality of evidence provided to downstream reasoning modules.

A.6 Details of Dataset

MultiTQ. MultiTQ is a large-scale temporal question answering dataset that incorporates multi-granularity temporal information. It provides a comprehensive evaluation protocol across several dimensions: Question Type (Multiple vs. Single), Answer Type (Entity vs. Time), and Time Granularity (year, month, and day). The detailed statistics of MultiTQ are presented in Table 5.

CronQuestions. CronQuestions is a large-scale temporal QA benchmark consisting of 410K unique question-answer pairs with annotated entities and timestamps. Questions in CronQuestions can be grouped into two major categories: simple

Type		Train #	Valid #	Test #
Single	Equal	135,890	18,983	17,311
	Before/After	75,340	11,655	11,073
	First/Last	72,252	11,097	10,480
Multiple	Equal Multi	16,893	3,213	3,207
	After First	43,305	6,499	6,266
	Before Last	43,107	6,532	6,247
Total		386,787	587,979	54,584

Table 5: Dataset Statistics of MultiTQ.

temporal reasoning (e.g., Simple Entity and Simple Time) and complex temporal reasoning (e.g., Before/After, First/Last, and Time-Join queries), depending on the temporal constraints involved. The detailed statistics of MultiTQ are presented in Table 6.

Category	Train #	Valid #	Test #
Simple Entity	90,651	7,745	7,812
Simple Time	61,471	5,197	5,046
Before/After	23,869	1,982	2,151
First/Last	118,556	11,198	11,159
Time Join	55,453	3,878	3,832
Simple Reasoning	152,122	12,942	12,858
Complex Reasoning	197,878	17,058	17,142
Entity Answer	225,672	19,362	19,524
Time Answer	124,328	10,638	10,476
Total	350,000	30,000	30,000

Table 6: Dataset Statistics of CronQuestions.

A.7 Case Study

Step-by-step Visualization As illustrated in Figure 6, we take the temporal question “Who was the first to visit Iraq after the Belgian Ministry?” as an example and present the key intermediate outputs produced by each module. This step-by-step visualization clearly demonstrates how TempReasoner, through the coordinated interaction of its multi-stage modules, progressively identifies the relevant evidence and ultimately arrives at the correct answer.

Specifically, The first sub-question aims to identify the time at which the Belgian Ministry visited Iraq. The Retrieve-Module rewrites the query into a TKG-compatible form and retrieves multiple candidate facts, many of which involve semantically related but irrelevant entities. The Filter-Module then performs entity–relation alignment, retaining only facts consistent with the inferred alignment constraints (i.e., *Ministry_(Belgium)*,

Make_a_visit, and *Iraq*). Based on the filtered evidence, the Generate-Module produces the answer *2009-06-04* together with its supporting fact. Since the retrieved evidence fully satisfies the temporal and relational constraints of the sub-question, the Validate-Module determines that no further verification is required and accepts the result as a reliable temporal anchor.

Using this inferred timestamp, the second sub-question is reformulated to identify entities that visited Iraq after *2009-06-04*. The retriever again collects candidate facts that satisfy the temporal constraint, including multiple visits occurring at different times. Through alignment-based filtering, the framework removes facts with incompatible relations or roles, leaving only valid visit events. The Generate-Module then selects *Chevron* as the earliest entity that visited Iraq after the anchor date, supported by the corresponding TKG fact.

At this stage, the Validate-Module constructs a stricter validation query by bounding the temporal interval between the inferred anchor date and the timestamp of the candidate answer. This query is re-evaluated through the same Retrieve–Filter–Generate pipeline to check whether an alternative entity satisfies the constraints more optimally. If a better-supported or earlier valid result were found during this verification step, the framework would replace the original answer accordingly. In this example, no such alternative is identified, and the initial answer *Chevron* is therefore selected as the final result. This case study demonstrates how TempReasoner incrementally refines evidence, enforces temporal consistency, and applies self-correction to ensure reliable temporal reasoning.

