

# AN EMPIRICAL STUDY ON RECONSTRUCTING SCIENTIFIC HISTORY TO FORECAST FUTURE TRENDS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

The advancement of scientific knowledge relies on synthesizing prior research to forecast future developments, a task that has become increasingly intricate. The emergence of large language models (LLMs) offers a transformative opportunity to automate and streamline this process, enabling faster and more accurate academic discovery. However, recent attempts either limit to producing surveys or focus overly on downstream tasks. To this end, we introduce a novel task that bridges two key challenges: the comprehensive synopsis of past research and the accurate prediction of emerging trends, dubbed *Dual Temporal Research Analysis*. This dual approach requires not only an understanding of historical knowledge but also the ability to predict future developments based on detected patterns. To evaluate, we present an evaluation benchmark encompassing 20 research topics and 210 key AI papers, based on the completeness of historical coverage and predictive reliability. We further draw inspirations from dual-system theory and propose a framework *HorizonAI* which utilizes a specialized temporal knowledge graph for papers, to capture and organize past research patterns (System 1), while leveraging LLMs for deeper analytical reasoning (System 2) to enhance both summarization and prediction. Our framework demonstrates a robust capacity to accurately summarize historical research trends and predict future developments, achieving significant improvements in both areas. For summarizing historical research, we achieve a 18.99% increase over AutoSurvey; for predicting future developments, we achieve a 7.71% increase over GPT-4o.

## 1 INTRODUCTION

For over 200,000 years, human intelligence has evolved, with knowledge-building processes underpinned by the dual imperatives of learning from the past and forecasting future directions (Sternberg, 2000). From the conceptual foundations of Ramon Llull’s “Tree of Knowledge” to Francis Bacon’s structured approach to human learning, both historical and contemporary scholars have emphasized the critical role of synthesizing past insights to drive future advancements. In recent years, modern frameworks addressing scientific discovery and knowledge structuring have further underscored this dual focus (Fire & Guestrin, 2019; Nagarajan et al., 2015).

The rapid growth of scientific publications presents an unprecedented challenge: researchers must now sift through vast amounts of literature to extract relevant historical insights and anticipate future trends (Fire & Guestrin, 2019). LLMs offer potential solutions by automating tasks such as retrieval, summarization, and analysis. However, most existing approaches either concentrate on retrospective literature reviews (Wang et al., 2024; Agarwal et al., 2024) or focus solely on generating novel research **by using simple concept-level link predictions lacking semantic relationships** (Krenn et al., 2023; Lu, 2021; Gu & Krenn; Tran & Xie, 2021). These narrow approaches neglect the essential integration of synthesizing past research with projecting future developments, a combination that is increasingly crucial for scientific discovery (Figure 1).

To address this gap, we propose *Dual Temporal Research Analysis (DTRA)*, a novel task that unifies the analysis of past research with the forecasting of future trends. In contrast to traditional methodologies, which focus on either historical synthesis or future speculation, our task bridges both by leveraging past knowledge to generate informed predictions. This twofold task is especially relevant in domains such as artificial intelligence (AI), where understanding prior research trajectories is essential for predicting emerging advancements.

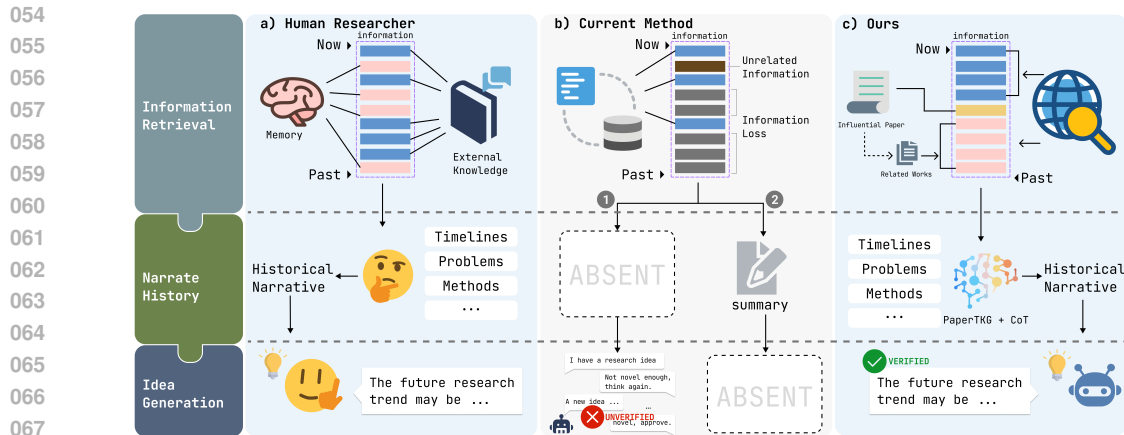


Figure 1: Comparison on dual temporal research analysis of a) human researchers b) current methods and c) our *HorizonAI*. Our framework resembles human researchers in the workflow while improving on thoroughness and logical reasoning, with both historical narrative and future prediction as output. In contrast, current methods focus only on either summarizing history (Wang et al., 2024; Edge et al., 2024) or generating future ideas (Baek et al., 2024; Si et al., 2024).

The *DTRA* consists of two interconnected phases: it involves consolidating and validating historical research trends, followed by the application of these insights to predict future developments. This approach mirrors the distinction between validation and experimentation, wherein past research serves as a foundation for verifying patterns and future predictions represent experimental, data-driven inferences (Chaiken, 1999; Posner, 2020).

Our framework *HorizonAI* (Figure 2) draws inspiration from Dual-System Theory (Chaiken, 1999), which posits that human cognition operates through two systems: System 1, which performs rapid, intuitive assessments, and System 2, which engages in deliberate, analytical reasoning. In the context of our framework, System 1 focuses on efficiently organizing historical data into structured formats such as temporal knowledge graphs (TKGs) (Cai et al., 2022), while System 2 conducts in-depth reasoning using Chain-of-Thought (CoT) (Wei et al., 2022; OpenAI, 2024b), to identify patterns and project future developments. Together, these systems enable a comprehensive analysis of both past and future research.

Given the novelty of the task, no established evaluation benchmarks or standardized methodologies currently exist. To address this, we introduce an evaluation benchmark *ResBench* that assesses the performance through historical completeness and predictive reliability. Extensive experiments demonstrate the superior performance of our proposed framework, *HorizonAI*, in tracing historical trends and making reliable future predictions. In comparison to existing baselines, it achieves higher predictive accuracy and generates more coherent, insightful content.

The main contributions of this paper are summarized as follows:

- **Integrating Historical Analysis and Future Forecasting:** We introduce *DTRA*, a task that combines the analysis of historical research with predictions about future trends. Unlike traditional methods that focus on either past research or future possibilities, this task incorporates both to generate informed projections, providing a more balanced perspective on scientific progress by ensuring that insights from the past inform future directions.
- **Cognitive-Inspired Framework:** Our approach *HorizonAI* is influenced by Dual-System Theory, suggesting that human cognition operates through both intuitive and analytical processes. In our framework, System 1 organizes past research into temporal knowledge graphs, while System 2 deliberately reasons to uncover patterns and anticipate future developments. This dual approach supports a more comprehensive, precise, and dynamic understanding of research trends.
- **Innovative Benchmark for Evaluation:** We propose a benchmark *ResBench* designed to evaluate *DTRA* based on historical coverage and predictive reliability. By incorporating datasets that span both historical and predictive dimensions, this benchmark provides the research community with a tool for systematically testing methods that summarize past research while forecasting future

developments. The integration of these two tasks enhances the performance of each, as insights from historical analysis inform predictions, leading to more accurate and contextually relevant forecasts.

## 2 DUAL-SYSTEM FRAMEWORK: *HorizonAI*

More formally, we define the *Dual Temporal Research Analysis (DTRA)* task as follows: Given the *input topic* and *source paper* ( $\mathcal{T}, \mathcal{P}$ ), the task is to narrate the *research history*  $H$  of the topic and generate *possible research ideas*  $F$ .

To achieve this task, we propose *HorizonAI* (as illustrated in Figure 2), a framework inspired by Dual-System Theory. We retrieve related papers to represent the history  $h_{s \sim t} = \{P_1, P_2, \dots, P_n\}$  during the time interval of  $s \sim t$  and structure it into a graph  $\mathcal{G}$  (i.e. our *PaperTKG*), then using strategy  $S$  to search the graph for timeline generation  $\tau_{s \sim t} = S(\mathcal{G})$ . The historical narrative  $H$  is generated by LLMs using temporal reasoning based on the timeline. We sample possible future predictions  $F, p(F|H) > threshold$  based on  $H$ .

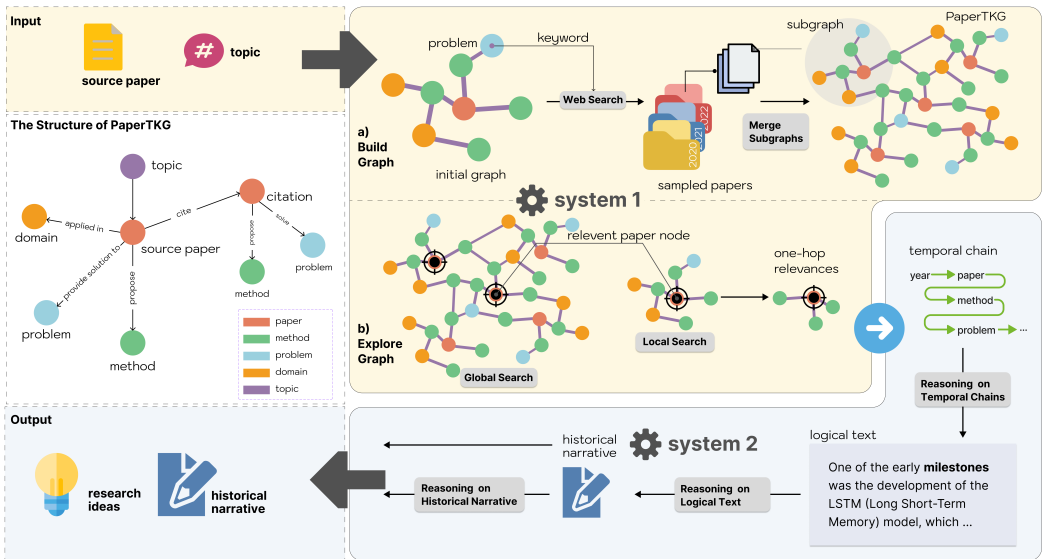


Figure 2: Diagram of *HorizonAI* framework. Given a topic and a source paper as input, *HorizonAI* goes through System 1 of a) structuring dynamically gathered historical information into *PaperTKG* and b) generating a timeline based on historical information and System 2 of reasoning on the timeline for historical narrative generation and future research trend prediction based on it.

### 2.1 *PaperTKG* CONSTRUCTION

To store historical information of research dynamically and structurally, we propose one specialized data structure—*PaperTKG*. It builds on the foundation of traditional TKGs by focusing specifically on academic papers. In *PaperTKG*, paper nodes are annotated with their timestamps and connected to entities such as methods, problems, domains, topics, and citations, as shown in the structure in Figure 2, storing all relevant and recent information to enhance reasoning in System 2 while reducing data processing costs.

The construction process of *PaperTKG* systematically progresses through three phases— building the initial graph, extending the graph by web search, and graph integration and refinement. The pseudo-code of our *PaperTKG* construction process can be found in Algorithm 1 and the prompts for it can be found in C.1.

**Paper2Graph** Converting papers to graphs (Paper2Graph) is a vital task in *PaperTKG* Construction. We use mainly the abstract and related work sections of a paper for that purpose. We derive the

**Algorithm 1** Temporal Knowledge Graph Construction

---

```

1: Input: Topic  $T$ , Source Paper  $S$ 
2: Output: Temporal Knowledge Graph  $G$ 
3: Phase 1: Build Initial Graph  $G_0 = \text{Paper2Graph}(T, S)$ 
4: Phase 2: Extend Graph  $G_0$ 
5: for each problem node  $p$  in  $G_0$  do
6:   for each year  $y$  from  $Y_{start}$  to  $Y_{end}$  do
7:     search for related papers with keywords  $\{T, p\}$ 
8:   end for
9: end for
10: for each new paper  $P_i, i = 1$  to  $N$  do
11:    $G_i = \text{Paper2Graph}(T, P_i)$ 
12: end for
13: Phase 3: Graph Integration
14: for each subgraph  $G_i, i = 1$  to  $N$  do
15:   for each problem node  $p$  in  $G_i$  do
16:     for each year  $y$  from  $Y_{start}$  to  $Y_{end}$  do
17:       search for related papers with keywords  $\{T, p\}$ 
18:     end for
19:   end for
20:   merge  $G_i$  into  $G_0$ 
21: end for
22: refine  $G$ : drop duplicates, discrete entities, complete missing entities
23: return  $G$ 

```

---

problem, application domain, and proposed method from the abstract while related works section provides insights into connections between existing methods, problems, and domains, annotated by the authors (Inevitable subjective bias is addressed by bulk sampling of papers - See Appendix A.3). Algorithm 2 details the pipeline.

**Algorithm 2** Paper2Graph: Entity and Relation Extraction

---

```

1: Input: topic  $T$ , paper  $P$ 
2: Output: subgraph  $\mathcal{G}$  with core concepts from  $P$  and its citations
3: Phase 0: Citation Matching
4: Match citations in related work to references
5: Create paper nodes and complete metadata via web search
6: Phase 1: Local Extraction
7: for each citation  $c$  do
8:   Extract method, problem, and domain related to  $c$  from context
9:   Establish entity relations
10: end for
11: Phase 2: Overall Connection
12: Infer relations between all entities
13: return  $\mathcal{G}$ 

```

---

**Graph Augmentation** We extend graphs built by Paper2Graph by incorporating subgraphs from papers cited and papers retrieved through a targeted web search. To create a concentrated historical dataset, we adopt a problem-centric sampling strategy, using problem nodes as search keywords rather than querying databases directly. Initially, problem nodes guide the first sampling round, yielding a fixed number of papers per year (also called uniform sampling, see Appendix A.2 for further explanation). Each sampled paper is converted to a subgraph using Paper2Graph, with their problem nodes driving the second and final sampling round. For a source paper with  $k_0$  problem nodes and  $n_0$  citations, the first round adds  $n_0 + L \cdot k_0 \cdot t$  subgraphs, sampling  $t$  papers annually over  $L$  years. Ultimately, we gather  $k_0 + \sum_{i=1}^{n_0+L \cdot k_0 \cdot t} k_i$  problems and  $1 + n_0 + \sum_{i=1}^{n_0+L \cdot k_0 \cdot t} n_i$  papers, ideally without duplication. To manage costs, we cap the total number of sampled papers.

## 2.2 QUERY GENERATION AND TEMPORAL REASONING

**Query Generation and Graph Exploration** Inspired by the local-to-global search strategy for summarization utilized by the GraphRAG (Edge et al., 2024), we design a global-to-local search strategy to narrow down the search range step by step without missing related nodes or relations. We first use global search to locate paper nodes related to the query, then apply local search to get detailed relations and neighbors of the paper node. **Global Search:** The query for global search is generated from the following three aspects: the application domain, the target problem to solve, and the method. We traverse all the paper nodes and select the ones related to our query by similarity, then we check the timestamps of these nodes to ensure that representative works from each year are included. **Local Search:** In the local search phase, we collect the one-hop relevances of the target paper nodes (i.e. the details of the paper) and structure them into a chain, finally, we get the timeline for the topic.

**Temporal Reasoning** Large Language Models (LLMs) are utilized to perform temporal reasoning (Yuan et al., 2024) using Chain-of-Thought (CoT) prompting (Wei et al., 2022). They are prompted to identify key research milestones, explain the methods and solutions, highlight connections between works, and emphasize the progression over time, step by step, to generate a coherent narrative of the research history. The prompt for converting the timeline to logical text can be found in Appendix C.2.

## 2.3 RESULT GENERATION

The output of our *HorizonAI* consists of two components: a historical narrative and a future prediction, with the latter being generated based on the former.

**Historical Narrative** We narrate the history from both local and holistic perspectives using temporal reasoning through CoT prompts (listed in Appendix C.2). For the local perspective, we structure the narrative by using the subtitles from the related work sections of the target papers as an outline, producing content resembling related work discussions. For the holistic perspective, we create outlines based on section titles from selected surveys to represent the overall development of the topic. Each section is then expanded with content following the outline, resulting in a survey-like narrative. **Future prediction** Existing approaches often emphasize the novelty of generated research ideas, overlooking that feasibility is a more critical factor than pure innovation. To ensure the ideas are grounded in practicality, we reference the local historical narrative and prompt LLMs to outline detailed roadmaps for realizing each idea. For each subdomain, multiple potential ideas are sampled, ensuring a balance between originality and implementability.

# 3 PROPOSED BENCHMARK: *ResBench*

## 3.1 DATA CONSTRUCTION

Table 1: Topics Used in Data Collection.

Index	Topic
1	In-context Learning
2	LLMs for Recommendation
3	LLMs-based Agents
4	Instruction Tuning for LLMs
5	LLMs for Information Retrieval
6	Safety in LLMs
7	Large Multi-Modal Language Models
8	LLMs for Software Engineering
9	LLM-Generated Texts Detection

Our dataset comprises papers from the arXiv<sup>1</sup> repository, specifically focusing on LLMs. **The dataset has 20 different topics, but considering the difficulty of manual verification, this article mainly evaluates 9 different topics (Table 1)**, covering the principles, techniques, and diverse applications of LLMs. For each topic, the dataset includes a source paper, a survey and at least 10 target papers. The source paper, an input for the task, serves as a starting point for collecting related literature to complete the historical information, while the surveys and target papers are used to evaluate the task outputs. In the surveys, every subsection’s content and title are included, alongside the corresponding references and notable research contributions. More details of the data composition are shown in Appendix B.

<sup>1</sup><https://arxiv.org/>

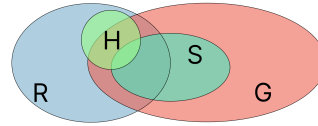
### 3.2 EVALUATION

The LLM evaluation consists of three main areas: historical completeness, predictive reliability, and text readability.

#### 3.2.1 GRAPH COMPLETENESS AND SEARCH EFFICIENCY

The completeness of graphs is evaluated by comparing the overlap between citations in target surveys and paper nodes in TKGs. Similarly, the search efficiency of graphs is assessed by measuring the overlap between paper nodes in retrieved during search and those in the surveys. The following sets involved in our evaluation are defined (their relationships are illustrated in Figure 3):

- $R$ : References in the surveys.
- $G$ : Paper nodes in the constructed graphs.
- $S$ : Paper nodes retrieved during graph search.
- $H$ : Key historical works in the surveys.



We define the following metrics to quantify the degree of overlap between these sets:

Figure 3: Venn Diagram. Note that  $S$  is a subset of  $G$  and  $H$  is a subset of  $R$ .

$$O_R = \frac{|G \cap R|}{|R|}, \quad O_H = \frac{|G \cap H|}{|H|}, \quad SE_R = \frac{|S \cap R|}{|G \cap R|}, \quad SE_H = \frac{|S \cap H|}{|G \cap H|}$$

Where  $O$  stands for *Overlap*, evaluating graph completeness upon construction and  $SE$  stands for *Search Efficiency*, evaluating graph search. More specifically:

- $O_R$ : Average proportion of target paper citations present in the generated graph. Used to evaluate the graph completeness.
- $O_H$ : Average proportion of key historical works present in the searched nodes. Used to evaluate the graph completeness, with a greater weight compared to  $O_R$ .
- $SE_R$ : The ratio between searched historical works and all the historical work referenced in the surveys and constructed in the graphs. Used to evaluate the search efficiency.
- $SE_H$ : The ratio between searched key historical works and all the key historical works referenced in the surveys and constructed in the graphs. Used to evaluate the search efficiency, with a greater weight compared to  $SE_R$ .

#### 3.2.2 PREDICTIVE RELIABILITY

Predictive reliability is evaluated through four perspectives: semantic similarity  $S_1$ , innovation and feasibility  $S_2$ , temporal consistency  $S_3$ , and contextual consistency  $S_4$ . All values are rated by LLM based on prompt instructions (detailed in Appendix C.3.1) on a scale of 1 to 5. The final rating is a weighted sum of these values:

$$\text{Final\_Score} = w_1 \cdot S_1 + w_2 \cdot S_2 + w_3 \cdot S_3 + w_4 \cdot S_4$$

Where  $w_1, w_2, w_3, w_4$  are weights for Semantic Similarity  $S_1$ , Innovation and Feasibility  $S_2$ , Temporal Consistency  $S_3$ , and Contextual Consistency  $S_4$ .

The explanation for the ranges of the final score is defined as:

- $\text{Final\_Score} \in [1,2)$ : The generated future directions show poor relevance to the target paper, with significant deficiencies in semantics, innovation, feasibility, or temporal consistency.
- $\text{Final\_Score} \in [2,3)$ : The generated future directions are somewhat relevant to the target paper but have several notable shortcomings.
- $\text{Final\_Score} \in [3,4)$ : The generated future directions are generally well-aligned with the target paper across multiple dimensions, though some improvements are still needed.
- $\text{Final\_Score} \in [4,5)$ : The generated future directions are highly relevant and excel in innovation, feasibility, temporal logic, and contextual consistency.

## 4 EXPERIMENTS

We use the GPT-4o (OpenAI, 2024a) for all the LLMs-involved processes (e.g. graph construction and reasoning) in our framework. We run experiments on the evaluation dataset on two subtasks for both our *HorizonAI* and baselines.

### 4.1 SUMMARIZING HISTORY - SURVEY COMPARISON

The performance of our history summarization subtask is assessed on the overlap degree of generated content and the target survey, whose result is illustrated in Table 7 where we also present the performance of graph completeness and search efficiency as a reference. As shown in Table 2, we compare the performance of our *HorizonAI* and AutoSurvey (Wang et al., 2024) in the history summarization subtask from three perspectives, namely total citation, key citation, and keyword. We conclude the performance of our framework as follows:

**Comprehensive historical representation** Under the conditions of limited information and the presence of bias in the writing of the target survey, the paper nodes in our *PaperTKG* have an average 39.35% overlap ratio with the citations in the human-written surveys and an even higher score of 46.35% regarding key citation overlap, demonstrating that our method of structuring history into *PaperTKG* to thoroughly and logically arrange scientific history is effective.

**Efficient Search Strategy** The average proportion of our searched works shared with the target survey among the total paper nodes of the graph reaches 69.92%, while on key references it is as high as 71.06%. The significantly high ratios of success indicate that using our global-to-local search strategy to fetch related paper nodes in the graph is powerful.

**Complete and Reliable History Summarization** A small proportion of searched past works are lost after the reasoning phase. The final generated content has an average of 24.25% citations in common with the target survey, compared to 5.26% of AutoSurvey. It reveals that our *HorizonAI* has a better ability to trace and summarize influential past works.

Table 2: Comparison of our *HorizonAI* and AutoSurvey (Wang et al., 2024) in history summarization subtask evaluated on nine topics. We use the overlap ratio of two aspects—total citation and key citation—to evaluate the performance of this task.

Evaluation Object	Citation Overlap(%)		Key Citation Overlap(%)	
	<i>HorizonAI</i> (ours)	AutoSurvey	<i>HorizonAI</i> (ours)	AutoSurvey
Topic 1	<b>42.86</b>	5.44	<b>53.01</b>	12.12
Topic 2	<b>27.32</b>	6.19	<b>38.93</b>	6.45
Topic 3	<b>2.27</b>	<b>2.27</b>	<b>37.87</b>	10.81
Topic 4	<b>24.24</b>	6.06	<b>47.98</b>	0.00
Topic 5	<b>27.93</b>	4.14	<b>46.19</b>	6.98
Topic 6	<b>5.00</b>	4.00	<b>10.00</b>	6.06
Topic 7	<b>35.57</b>	4.35	<b>50.00</b>	12.50
Topic 8	<b>19.75</b>	8.02	<b>24.24</b>	5.48
Topic 9	<b>33.33</b>	6.86	<b>25.40</b>	8.33
Average	<b>24.25</b>	5.26	<b>37.07</b>	7.64

### 4.2 PREDICTING FUTURE - RELATED WORKS COMPARISON

We use the subtitles of related work of the target paper as a guideline to generate the possible research idea, then we compare this generated idea with the actual one proposed by this paper in the abstract. The performance of the future prediction is evaluated on the comprehensive score of the content, covering content quality, relevance, innovation, and so on. Due to the existing works aimed at idea generation mainly focusing on novelty, they will naturally filter out previous works. In response to this situation, we use LLM and horizonAI without temporal logic reasoning as our baseline to evaluate how much the performance of HorizonAI will drop without adopting a workflow inspired by dual system theory. The result of this subtask, as is illustrated in Table 3, proved that with adequate historical narrative and temporal logic reasoning, LLMs can produce more reliable research ideas than the ones without.

Table 3: Comparison between our *HorizonAI*, LLMs and *HorizonAI* without temporal logic reasoning in future prediction. The final score calculation method is shown in Section 3. The full score is 5.

Evaluation Object	Baseline	Without Temporal Logic Reasoning	<i>HorizonAI</i> (ours)
Topic 1	3.77	2.30	3.91
Topic 2	3.42	1.10	3.98
Topic 3	3.88	1.25	3.85
Topic 4	3.68	2.20	3.78
Topic 5	3.50	2.05	3.76
Topic 6	3.69	1.15	4.01
Topic 7	3.44	2.00	3.87
Topic 8	3.84	2.28	3.88
Topic 9	3.23	1.90	3.91
Average	3.60	2.14	3.88

Table 4: Ablation for input quality. Problem type 1 stands for the interdisciplinary topics, and topics related to it are: A - Bias and Fairness in LLMs, B - LLMs in Medicine, C - Domain Specialization of LLMs, D - Challenges of LLMs in Education. Problem type 2 stands for the topic and source paper misalignment case, and the topic related to it is: E - Explainability for LLMs

Problem Type	Topic	Source Paper	Number of Paper Nodes	History Completeness (%)	Citation Overlap (%)
1	A	Gupta et al. (2023)	44	2.99	1.45
	B	Singhal et al. (2023)	996	6.19	2.44
	C	Li et al. (2022)	131	16.67	5.56
	D	Leinonen et al. (2023)	36	0.00	0.00
2	E	Gao et al. (2023)	620	2.68	1.52

### 4.3 ABLATION STUDY

**Effect of Input** We designed two possible problem types regarding the inputs to determine their influence on our method (as illustrated in Table 4). The first case involves interdisciplinary topics that require more relevant historical information compared to topics within the AI field. Additionally, obtaining related data from arXiv is relatively more challenging. Topics related to medicine, education, and society are selected for it. The results show that our method with a broad cross-domain topic as input suffers from graph augmentation failure, leading to an unwanted history completion performance. The second case involves a mismatch between the topic and source paper. In this case, a source paper with less relevance to the topic is given as input. This leads to the graph used for representing history expanding in the wrong direction, which explains the bad performance. In conclusion, our method is sensible to the inputs (i.e. the topic and the source paper), either a vague topic or mismatched inputs will lead to unwanted results.

**Effect of Graph Augmentation Strategy** We test the performance of history completeness under four different graph argumentation methods to determine the effect of web retrieval strategy on our framework. The complexity of collecting historical work increases sequentially from Method 1 to Method 4, with Method 4 being the one used in our framework. As shown in Figure 4, the performance variation trends of different search strategies across topics are consistent, and the more comprehensive the search method, the higher the citation overlap of the enhanced graph. Among them, Method 4 achieves the best performance across all topics. This result indicates that the graph argumentation method has a significant impact

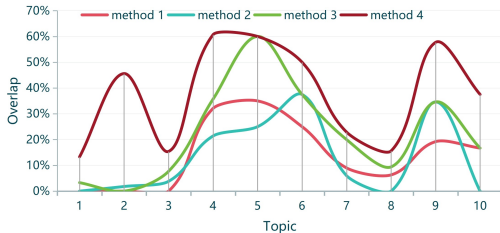


Figure 4: Diagram of search strategy performance. Method 1 is human efforts, Method 2 is greedy-search from the original problem in source paper, Method 3 uses similarity ranking to search from the original problem, while Method 4 (ours) updates on method 3 by searching on all central problems.



432 on completing historical data. More data does not necessarily mean better historical reproduction;  
433 rather, retrieving data from multiple dimensions is more beneficial for historical completion.  
434

## 435 5 RELATED WORKS

### 436 5.1 RETRIEVAL ON KNOWLEDGE GRAPHS (KGs)

437  
438 Recent search strategies using Knowledge Graphs (KGs) improve retrieval by leveraging structured  
439 relationships in Large Language Models (LLMs) to enhance inference and interpretability (Pan et al.,  
440 2024; Yang et al., 2024). Structuring LLM interactions with KGs refines retrieval performance,  
441 enabling effective responses to complex queries (Sun et al., 2023; Jiang et al., 2023; 2024). Unlike  
442 traditional RAG methods that rely on text embeddings, KGs serve as indices to enhance precision  
443 by navigating relevant subgraphs. Approaches like KAPING (Baek et al., 2023), G-Retriever (He  
444 et al., 2024), and Graph-ToolFormer (Zhang, 2023) enhance retrieval by using graph metrics to refine  
445 search results, while SURGE (Kang et al., 2023) and FABULA (Ranade & Joshi, 2023) leverage  
446 KGs for narrative generation grounded in factual subgraphs. Systems like ITRG (Feng et al., 2023)  
447 and IR-CoT (Trivedi et al., 2023) facilitate multi-hop question answering by tracing interconnected  
448 knowledge nodes, while Selfmem (Cheng et al., 2023) employs KGs for generation-augmented  
449 retrieval. GraphRAG (Edge et al., 2024) advances these approaches by introducing a unique local-  
450 to-global search strategy with a self-generated graph index, which inspired us for our global-to-local  
451 retrieval approach.  
452

### 453 5.2 TEMPORAL KNOWLEDGE GRAPHS (TKGs)

454  
455 Temporal Knowledge Graphs (TKGs) extend traditional Knowledge Graphs(KGs) by incorporating  
456 temporal information, enabling the representation of dynamically changing facts (Ji et al., 2021).  
457 Temporal Knowledge Graph Completion (TKGC) is a key task in TKGs, focusing on filling in miss-  
458 ing information and predicting future relationships (Wang et al., 2023; Xu et al., 2023; Zhang et al.,  
459 2023; Xiong et al., 2024b). Additionally, TKGs specialize in applications like event tracking and  
460 historical data analysis, providing a more nuanced framework for mapping extensive academic liter-  
461 ature. Specialized models such as Know-Evolve (Trivedi et al., 2017) and TA-TransE (García-Durán  
462 et al., 2018) have advanced temporal reasoning, while Xiong et al. (2024a) introduced a two-step  
463 framework for language-based temporal reasoning that translates narratives into TKGs. Despite  
464 these innovations, many existing models remain general and do not specifically address the orga-  
465 nizational needs of academic information. Our proposed *PaperTKG* thus serves as a tailored TKG  
466 structure designed explicitly for managing scholarly papers, enabling the tracking of paper evolu-  
467 tion, citation networks, and topic trends over time, thereby fulfilling the demand for a specialized  
468 TKG in academia.  
469

### 470 5.3 LLMs IN SCIENTIFIC DEVELOPMENT

471  
472 Large Language Models (LLMs) are recognized for their transformative potential in scientific re-  
473 search, owing to their ability to process and analyze vast datasets beyond human capacity. Recent  
474 studies, such as those by Baek et al. (2024), Yang et al. (2023), and Qi et al. (2023), focus on  
475 Literature-based Discovery (LBD) (Swanson, 1986), using LLMs to mine academic publications  
476 for correlations and generate research insights. Wang et al. (2024) explores the possibility of LLMs  
477 automatically generating survey papers, while other works (Elsevier, 2024; Agarwal et al., 2024)  
478 emphasize automated retrieval and summarization of existing literature, often neglecting the pre-  
479 diction of future research trends. In a pioneering effort, Li & Flanigan (2024) formalizes future  
480 language modeling, aiming to predict future textual data based on temporal histories. Additionally,  
481 several studies (Si et al., 2024; Baek et al., 2024; Zheng et al., 2024) develop LLM-based agents  
482 for research idea generation, a critical step in the early stages of scientific inquiry. AI Scientist (Lu  
483 et al., 2024) represents the first comprehensive system for fully automated scientific discovery using  
484 LLMs, generating novel research ideas independent of prior work, though it requires multiple iter-  
485 ations to yield viable outcomes. In contrast, we introduce *HorizonAI*, a dual-system approach that  
integrates both the summarization of past research and the prediction of future directions, offering  
superior performance in both tasks compared to existing models.

## 6 CONCLUSION

In this paper, we introduced *Dual Temporal Research Analysis (DTRA)*, a novel task that integrates the summarization of historical research with the prediction of future trends. Our framework, *HorizonAI*, draws inspiration from Dual-System Theory to organize past research using *PaperTKG* (a temporal knowledge graph for papers) and employs LLMs for in-depth reasoning to generate both historical narratives and future projections.

Through extensive evaluation on the *ResBench* benchmark, we demonstrated that bridging the tasks of historical analysis and future forecasting enhances the performance of both. Our results showed significant improvements in summarizing past works and generating accurate predictions compared to existing methods.

The integration of historical insights with predictive reasoning offers a balanced perspective on scientific progress, showing the potential of *HorizonAI* as a robust tool for supporting research across multiple domains. Future work will focus on expanding the dataset beyond AI-related topics and enhancing search capabilities to incorporate a wider range of academic databases.

## 7 LIMITATION AND FUTURE WORKS

### 7.1 LIMITATIONS

1. The dataset currently focuses on AI-related topics with surveys available in 2024, but it can be extended to a broader range of domains. This design was chosen to facilitate more precise evaluation and easier expert feedback, but future work should include diverse research fields to enhance generalizability.
2. Currently, the search and graph construction processes are time-consuming due to the reliance on third-party web APIs that often struggle with bulk access. This issue can be addressed by using specialized API keys or developing our own databases. Nevertheless, the current system still offers higher efficiency compared to manual research, and the constructed graphs can be reused for further analysis.
3. **The current assessment of future idea generation relies solely on LLMs in content analysis; however, while the accuracy of utilizing our algorithm is guaranteed, the inclusion of expert reviewers will provide additional insight into the feasibility and reliability of the ideas generated.**

### 7.2 FUTURE WORKS

In addressing previous limitations, we encourage extending the dataset beyond AI-related topics to include a broader range of research fields. This expansion would allow for a more comprehensive evaluation of the framework’s generalizability across different domains. Additionally, we plan to enhance our data collection by incorporating papers from other sources beyond Arxiv, such as peer-reviewed journals and other preprint servers, using advanced tools for PDF information extraction. Furthermore, integrating expert reviews into the evaluation process will provide more reliable insights into the feasibility and practical relevance of the generated future ideas, moving beyond sole reliance on LLM evaluations.

On the other hand, we will continue to explore ways to enhance research efficiency in the era of LLMs and AI. This area holds significant potential, and beyond generalization and future direction prediction, we aim to enable AI to contribute to the actual realization of future research topics. This will involve collaboration with researchers in experiment design and result analysis, integrating AI more deeply into the research process.

## REFERENCES

- 540  
541  
542 Shubham Agarwal, Issam H Laradji, Laurent Charlin, and Christopher Pal. Litllm: A toolkit for  
543 scientific literature review. *arXiv preprint arXiv:2402.01788*, 2024.
- 544  
545 Jinheon Baek, Alham Fikri Aji, and Amir Saffari. Knowledge-augmented language model prompt-  
546 ing for zero-shot knowledge graph question answering, 2023. URL [https://arxiv.org/  
547 abs/2306.04136](https://arxiv.org/abs/2306.04136).
- 548  
549 Jinheon Baek, Sujay Kumar Jauhar, Silviu Cucerzan, and Sung Ju Hwang. Researchagent: Iterative  
550 research idea generation over scientific literature with large language models. *arXiv preprint  
551 arXiv:2404.07738*, 2024.
- 552  
553 Borui Cai, Yong Xiang, Longxiang Gao, He Zhang, Yunfeng Li, and Jianxin Li. Temporal knowl-  
554 edge graph completion: A survey. *arXiv preprint arXiv:2201.08236*, 2022.
- 555  
556 Shelly Chaiken. Dual-process theories in social psychology. *Guilford Press google schola*, 2:206–  
557 214, 1999.
- 558  
559 Xin Cheng, Di Luo, Xiuying Chen, Lema Liu, Dongyan Zhao, and Rui Yan. Lift yourself up:  
560 Retrieval-augmented text generation with self memory, 2023. URL [https://arxiv.org/  
561 abs/2305.02437](https://arxiv.org/abs/2305.02437).
- 562  
563 Darren Edge, Ha Trinh, Newman Cheng, Joshua Bradley, Alex Chao, Apurva Mody, Steven Truitt,  
564 and Jonathan Larson. From local to global: A graph rag approach to query-focused summariza-  
565 tion. *arXiv preprint arXiv:2404.16130*, 2024.
- 566  
567 Elsevier. <https://www.elsevier.com/products/scopus/scopus-ai>. [https://www.elsevier.  
568 com/products/scopus/scopus-ai](https://www.elsevier.com/products/scopus/scopus-ai), 2024.
- 569  
570 Zhangyin Feng, Xiaocheng Feng, Dezhi Zhao, Maojin Yang, and Bing Qin. Retrieval-generation  
571 synergy augmented large language models, 2023. URL [https://arxiv.org/abs/2310.  
572 05149](https://arxiv.org/abs/2310.05149).
- 573  
574 Michael Fire and Carlos Guestrin. Over-optimization of academic publishing metrics: observing  
575 goodhart’s law in action. *GigaScience*, 8(6):giz053, 2019.
- 576  
577 Yunfan Gao, Tao Sheng, Youlin Xiang, Yun Xiong, Haofen Wang, and Jiawei Zhang. Chat-  
578 rec: Towards interactive and explainable llms-augmented recommender system. *arXiv preprint  
579 arXiv:2303.14524*, 2023.
- 580  
581 Alberto García-Durán, Sebastijan Dumančić, and Mathias Niepert. Learning sequence encoders  
582 for temporal knowledge graph completion, 2018. URL [https://arxiv.org/abs/1809.  
583 03202](https://arxiv.org/abs/1809.03202).
- 584  
585 Xuemei Gu and Mario Krenn. Impact4cast: Forecasting high-impact research topics via machine  
586 learning on evolving knowledge graphs. In *ICML 2024 AI for Science Workshop*.
- 587  
588 Shashank Gupta, Vaishnavi Shrivastava, Ameet Deshpande, Ashwin Kalyan, Peter Clark, Ashish  
589 Sabharwal, and Tushar Khot. Bias runs deep: Implicit reasoning biases in persona-assigned llms.  
590 *arXiv preprint arXiv:2311.04892*, 2023.
- 591  
592 Xiaoxin He, Yijun Tian, Yifei Sun, Nitesh V. Chawla, Thomas Laurent, Yann LeCun, Xavier Bres-  
593 son, and Bryan Hooi. G-retriever: Retrieval-augmented generation for textual graph understand-  
ing and question answering, 2024. URL <https://arxiv.org/abs/2402.07630>.
- 588  
589 Shaoxiong Ji, Shirui Pan, Erik Cambria, Pekka Marttinen, and S Yu Philip. A survey on knowledge  
590 graphs: Representation, acquisition, and applications. *IEEE transactions on neural networks and  
591 learning systems*, 33(2):494–514, 2021.
- 592  
593 Jinhao Jiang, Kun Zhou, Zican Dong, Keming Ye, Wayne Xin Zhao, and Ji-Rong Wen. Structgpt:  
A general framework for large language model to reason over structured data. *arXiv preprint  
arXiv:2305.09645*, 2023.

- 594 Jinhao Jiang, Kun Zhou, Wayne Xin Zhao, Yang Song, Chen Zhu, Hengshu Zhu, and Ji-Rong  
595 Wen. Kg-agent: An efficient autonomous agent framework for complex reasoning over knowl-  
596 edge graph. *arXiv preprint arXiv:2402.11163*, 2024.
- 597
- 598 Minki Kang, Jin Myung Kwak, Jinheon Baek, and Sung Ju Hwang. Knowledge graph-augmented  
599 language models for knowledge-grounded dialogue generation, 2023. URL <https://arxiv.org/abs/2305.18846>.
- 600
- 601 Rodney Michael Kinney, Chloe Anastasiades, Russell Authur, Iz Beltagy, Jonathan Bragg, Alexan-  
602 dra Buraczynski, Isabel Cachola, Stefan Candra, Yoganand Chandrasekhar, Arman Cohan, Miles  
603 Crawford, Doug Downey, Jason Dunkelberger, Oren Etzioni, Rob Evans, Sergey Feldman, Joseph  
604 Gorney, David W. Graham, F.Q. Hu, Regan Huff, Daniel King, Sebastian Kohlmeier, Bailey  
605 Kuehl, Michael Langan, Daniel Lin, Haokun Liu, Kyle Lo, Jaron Lochner, Kelsey MacMil-  
606 lan, Tyler C. Murray, Christopher Newell, Smita R Rao, Shaurya Rohatgi, Paul Sayre, Zejiang  
607 Shen, Amanpreet Singh, Luca Soldaini, Shivashankar Subramanian, A. Tanaka, Alex D Wade,  
608 Linda M. Wagner, Lucy Lu Wang, Christopher Wilhelm, Caroline Wu, Jiangjiang Yang, An-  
609 gele Zamarron, Madeleine van Zuylen, and Daniel S. Weld. The semantic scholar open data  
610 platform. *ArXiv*, abs/2301.10140, 2023. URL <https://api.semanticscholar.org/CorpusID:256194545>.
- 611
- 612 Mario Krenn, Lorenzo Buffoni, Bruno Coutinho, Sagi Eppel, Jacob Gates Foster, Andrew Grit-  
613 sevskiy, Harlin Lee, Yichao Lu, João P Moutinho, Nima Sanjabi, et al. Forecasting the future  
614 of artificial intelligence with machine learning-based link prediction in an exponentially growing  
615 knowledge network. *Nature Machine Intelligence*, 5(11):1326–1335, 2023.
- 616
- 617 Juho Leinonen, Paul Denny, Stephen MacNeil, Sami Sarsa, Seth Bernstein, Joanne Kim, Andrew  
618 Tran, and Arto Hellas. Comparing code explanations created by students and large language  
619 models. In *Proceedings of the 2023 Conference on Innovation and Technology in Computer  
620 Science Education V. 1*, pp. 124–130, 2023.
- 621
- 622 Changmao Li and Jeffrey Flanigan. Future language modeling from temporal document history.  
623 *arXiv preprint arXiv:2404.10297*, 2024.
- 624
- 625 Margaret Li, Suchin Gururangan, Tim Dettmers, Mike Lewis, Tim Althoff, Noah A Smith, and Luke  
626 Zettlemoyer. Branch-train-merge: Embarrassingly parallel training of expert language models.  
627 *arXiv preprint arXiv:2208.03306*, 2022.
- 628
- 629 Chris Lu, Cong Lu, Robert Tjarko Lange, Jakob Foerster, Jeff Clune, and David Ha. The ai scien-  
630 tist: Towards fully automated open-ended scientific discovery. *arXiv preprint arXiv:2408.06292*,  
631 2024.
- 632
- 633 Yichao Lu. Predicting research trends in artificial intelligence with gradient boosting decision trees  
634 and time-aware graph neural networks. In *2021 IEEE International Conference on Big Data (Big  
635 Data)*, pp. 5809–5814. IEEE, 2021.
- 636
- 637 Meenakshi Nagarajan, Angela D Wilkins, Benjamin J Bachman, Ilya B Novikov, Shenghua Bao,  
638 Peter J Haas, María E Terrón-Díaz, Sumit Bhatia, Anbu K Adikesavan, Jacques J Labrie, et al.  
639 Predicting future scientific discoveries based on a networked analysis of the past literature. In  
640 *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and  
641 Data Mining*, pp. 2019–2028, 2015.
- 642
- 643 OpenAI. Hello GPT-4o. <https://openai.com/index/hello-gpt-4o/>, 2024a.
- 644
- 645 OpenAI. Learning to reason with llms, 2024b. URL <https://openai.com/index/learning-to-reason-with-llms>.
- 646
- 647 Shirui Pan, Linhao Luo, Yufei Wang, Chen Chen, Jiapu Wang, and Xindong Wu. Unifying large  
648 language models and knowledge graphs: A roadmap. *IEEE Transactions on Knowledge and Data  
649 Engineering*, 2024.
- 650
- 651 Ingmar Posner. Robots thinking fast and slow: on dual process theory and metacognition in embod-  
652 ied ai. *None*, 2020.

- 648 Biqing Qi, Kaiyan Zhang, Haoxiang Li, Kai Tian, Sihang Zeng, Zhang-Ren Chen, and Bowen Zhou.  
649 Large language models are zero shot hypothesis proposers. *arXiv preprint arXiv:2311.05965*,  
650 2023.
- 651 Priyanka Ranade and Anupam Joshi. Fabula: Intelligence report generation using retrieval-  
652 augmented narrative construction. In *Proceedings of the International Conference on Advances*  
653 *in Social Networks Analysis and Mining*, volume 40 of *ASONAM '23*, pp. 603–610. ACM,  
654 November 2023. doi: 10.1145/3625007.3627505. URL [http://dx.doi.org/10.1145/](http://dx.doi.org/10.1145/3625007.3627505)  
655 [3625007.3627505](http://dx.doi.org/10.1145/3625007.3627505).
- 656 Chenglei Si, Diyi Yang, and Tatsunori Hashimoto. Can llms generate novel research ideas? a large-  
657 scale human study with 100+ nlp researchers. *arXiv preprint arXiv:2409.04109*, 2024.
- 658 Karan Singhal, Shekoofeh Azizi, Tao Tu, S Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan  
659 Scales, Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl, et al. Large language models encode  
660 clinical knowledge. *Nature*, 620(7972):172–180, 2023.
- 661 Robert J Sternberg. *Handbook of intelligence*. Cambridge University Press, 2000.
- 662 Jiashuo Sun, Chengjin Xu, Lumingyuan Tang, Saizhuo Wang, Chen Lin, Yeyun Gong, Heung-  
663 Yeung Shum, and Jian Guo. Think-on-graph: Deep and responsible reasoning of large language  
664 model with knowledge graph. *arXiv preprint arXiv:2307.07697*, 2023.
- 665 Don R Swanson. Undiscovered public knowledge. *The Library Quarterly*, 56(2):103–118, 1986.
- 666 Ngoc Mai Tran and Yangxinyu Xie. Improving random walk rankings with feature selection and  
667 imputation science4cast competition, team hash brown. In *2021 IEEE International Conference*  
668 *on Big Data (Big Data)*, pp. 5824–5827. IEEE, 2021.
- 669 Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. Interleaving re-  
670 trieval with chain-of-thought reasoning for knowledge-intensive multi-step questions, 2023. URL  
671 <https://arxiv.org/abs/2212.10509>.
- 672 Rakshit Trivedi, Hanjun Dai, Yichen Wang, and Le Song. Know-evolve: Deep temporal reasoning  
673 for dynamic knowledge graphs, 2017. URL <https://arxiv.org/abs/1705.05742>.
- 674 Jiapu Wang, Boyue Wang, Meikang Qiu, Shirui Pan, Bo Xiong, Heng Liu, Linhao Luo, Tengfei Liu,  
675 Yongli Hu, Baocai Yin, et al. A survey on temporal knowledge graph completion: Taxonomy,  
676 progress, and prospects. *arXiv preprint arXiv:2308.02457*, 2023.
- 677 Yidong Wang, Qi Guo, Wenjin Yao, Hongbo Zhang, Xin Zhang, Zhen Wu, Meishan Zhang, Xinyu  
678 Dai, Min Zhang, Qingsong Wen, et al. Autosurvey: Large language models can automatically  
679 write surveys. *arXiv preprint arXiv:2406.10252*, 2024.
- 680 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed Chi, Quoc V  
681 Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models.  
682 In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in*  
683 *Neural Information Processing Systems*, volume 35, pp. 24824–24837. Curran Associates, Inc.,  
684 2022. URL [https://proceedings.neurips.cc/paper\\_files/paper/2022/](https://proceedings.neurips.cc/paper_files/paper/2022/file/9d5609613524ecf4f15af0f7b31abca4-Paper-Conference.pdf)  
685 [file/9d5609613524ecf4f15af0f7b31abca4-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/9d5609613524ecf4f15af0f7b31abca4-Paper-Conference.pdf).
- 686 Siheng Xiong, Ali Payani, Ramana Kompella, and Faramarz Fekri. Large language models can learn  
687 temporal reasoning. *arXiv preprint arXiv:2401.06853*, 2024a.
- 688 Siheng Xiong, Yuan Yang, Faramarz Fekri, and James Clayton Kerce. Tilp: Differentiable learning  
689 of temporal logical rules on knowledge graphs. *arXiv preprint arXiv:2402.12309*, 2024b.
- 690 Wenjie Xu, Ben Liu, Miao Peng, Xu Jia, and Min Peng. Pre-trained language model with prompts  
691 for temporal knowledge graph completion. *arXiv preprint arXiv:2305.07912*, 2023.
- 692 Linyao Yang, Hongyang Chen, Zhao Li, Xiao Ding, and Xindong Wu. Give us the facts: Enhancing  
693 large language models with knowledge graphs for fact-aware language modeling. *IEEE Transac-*  
694 *tions on Knowledge and Data Engineering*, 2024.

702 Zonglin Yang, Xinya Du, Junxian Li, Jie Zheng, Soujanya Poria, and Erik Cambria. Large  
703 language models for automated open-domain scientific hypotheses discovery. *arXiv preprint*  
704 *arXiv:2309.02726*, 2023.

705  
706 Chenhan Yuan, Qianqian Xie, Jimin Huang, and Sophia Ananiadou. Back to the future: To-  
707 wards explainable temporal reasoning with large language models. In *Proceedings of the ACM*  
708 *Web Conference 2024, WWW '24*, pp. 1963–1974, New York, NY, USA, 2024. Association  
709 for Computing Machinery. ISBN 9798400701719. doi: 10.1145/3589334.3645376. URL  
710 <https://doi.org/10.1145/3589334.3645376>.

711 Jiawei Zhang. Graph-toolformer: To empower llms with graph reasoning ability via prompt aug-  
712 mented by chatgpt, 2023. URL <https://arxiv.org/abs/2304.11116>.

713  
714 Mengqi Zhang, Yuwei Xia, Qiang Liu, Shu Wu, and Liang Wang. Learning latent relations for tem-  
715 poral knowledge graph reasoning. In *Proceedings of the 61st Annual Meeting of the Association*  
716 *for Computational Linguistics (Volume 1: Long Papers)*, pp. 12617–12631, 2023.

717 Yuxiang Zheng, Shichao Sun, Lin Qiu, Dongyu Ru, Cheng Jiayang, Xuefeng Li, Jifan Lin, Bin-  
718 jie Wang, Yun Luo, Renjie Pan, et al. Openresearcher: Unleashing ai for accelerated scientific  
719 research. *arXiv preprint arXiv:2408.06941*, 2024.

720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755

756 A FURTHER CLARIFICATION ON STRATEGY

757  
758 A.1 HOW DO WE REPRESENT RESEARCH HISTORY?

759  
760 There are various ways to define the history of research, but academic papers remain one of the  
761 most common forms of scholarly communication. A typical method for scientists to understand the  
762 evolution of a field is by reviewing related papers. Thus, we define the history of a field or research  
763 question as the collection of papers associated with it:

764  
765 
$$H = \{P_1, P_2, \dots, P_n\}$$
 (1)

766 where  $P_i$  denotes the  $i$ -th paper in the collection, and  $n$  is the total number of relevant papers.

767 To expand this collection more effectively, we extract papers from the related works sections:

768  
769 
$$H' = \bigcup_{i=1}^n RW(P_i)$$
 (2)

770 where  $H'$  represents the extended history, and  $RW(P_i)$  is the set of cited works in  $P_i$ .

771 The progression of knowledge in a field can be represented as a trajectory, with papers ordered  
772 temporally:

773  
774 
$$T = \{P_{\sigma(1)}, P_{\sigma(2)}, \dots, P_{\sigma(n)} \mid t_{\sigma(1)} \leq t_{\sigma(2)} \leq \dots \leq t_{\sigma(n)}\}$$
 (3)

775 where  $\sigma$  is a permutation of indices ensuring the papers are arranged chronologically.

776 To represent this trajectory more efficiently, we use temporal knowledge graphs, detailed further in  
777 Section 3 of the Methods.

778  
779 A.2 HOW DO WE SAMPLE FROM HISTORICAL PAPERS?

780 Sampling is essential to ensure the accuracy and diversity of research findings, mitigating bias.  
781 Table 5 compares several sampling strategies: *Uniform Sampling*, *Proportional Sampling*, *Citation-*  
782 *based Sampling*, *Random Sampling*, and *Stratified Sampling*. We select *Uniform Sampling* due to  
783 its ability to maintain an even distribution across years, minimizing variance and offering ease of  
784 implementation.

785  
786 Table 5: Overview of Sampling Methods and Their Variance. The other methods have  $> 0$  variance  
787 influenced by factors such as publication volume, citation counts, and strata representation.  $P(p_i)$   
788 is the probability of selecting papers  $p_i$ ,  $N_j$  is the number of papers in a time period  $y_j$ ,  $N$  is the  
789 total number of papers,  $C_i$  is the citation count of paper  $p_i$ ,  $M$  is the total sample size, and  $m_j$  is the  
790 sample size for each period or stratum.

791  
792  
793  
794  
795  
796  
797  
798  
799

Sampling Method	Mathematical Definition	Pros	Cons	Variance in $m_j$
Uniform Sampling	$P(p_i) = \frac{1}{ P_{y_j} }$ , $\forall i \in y_j$	Balanced representation across time periods	May exclude influential papers from prolific years	0
Proportional Sampling	$P(p_i) = \frac{N_j}{N}$ , $m_j = P(p_i) \cdot M$	Reflects natural publication volume	Over-represents years with high publication counts	$> 0$
Citation-based Sampling	$P(p_i) = \frac{C_i}{\sum_{P \in P_{y_j}} C_P}$	Focuses on highly influential papers	Skews toward older papers; Ignores recent work	$> 0$
Random Sampling	$P(p_i) = \frac{1}{N}$	Simple, unbiased by time or citation	May miss important trends; Over-represents recent years	$> 0$
Stratified Sampling	$P(p_i) = \frac{N_j}{N}$ , $m_j = P(p_i) \cdot M$	Ensures representation across strata	Complex to implement; Over-represent dominant strata	$> 0$

800  
801  
802  
803  
804  
805  
806  
807  
808  
809

810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863

### A.3 ADDRESSING SUBJECTIVE BIAS IN PUBLICATIONS

Subjective bias is inherent in individual papers, as each presents knowledge from a particular viewpoint. Consequently, integrating biased papers into a study introduces this subjectivity. However, by aggregating enough diverse papers, we can mitigate individual biases and approach a more objective historical representation:

$$B(H) = \sum_{i=1}^n B(P_i) \tag{4}$$

To reduce bias, we aim to incorporate a sufficiently large set of papers. The bias in an expanded collection  $H'$  of papers can be approximated by:

$$B(H') \approx \frac{1}{|H'|} \sum_{i=1}^{|H'|} B(P_i) \tag{5}$$

As the size of  $H'$  increases, the overall bias approaches a more balanced representation of the field.



## 864 B DATA COLLECTION

865  
866 The data collection process focused on identifying relevant papers across three categories: **source**  
867 **papers**, **target papers**, and **surveys**. These papers were selected based on their influence, citation  
868 count, and relevance to the topic, ensuring a comprehensive overview of the field.

### 870 B.1 SOURCE PAPERS

871  
872 Source papers refer to highly cited and influential papers published before **2023**. These papers were  
873 carefully chosen based on their significant contributions to the field and their role in shaping founda-  
874 tional knowledge. Each source paper was analyzed for its related work sections, which included:

- 875 • The titles of cited references.
- 876 • Summaries of key points from these references.

877  
878 We selected these papers using Semantic Scholar’s influential sorting feature (Kinney et al., 2023),  
879 ensuring the source papers were ranked by citation count. Data was extracted from HTML and PDF  
880 formats, with arXiv<sup>2</sup> and ar5iv<sup>3</sup> providing the HTML versions for most papers. All data underwent  
881 manual verification to ensure accuracy.

### 883 B.2 TARGET PAPERS

884  
885 Target papers refer to newer, high-quality papers and surveys from **2024**, chosen for their cutting-  
886 edge insights. These papers were similarly ranked by **citation count** and were selected to reflect the  
887 most current trends and advancements in the field. In target papers, we also focused on their related  
888 work sections, capturing:

- 889 • The titles of cited references.
- 890 • Key points and summaries relevant to the topic.

891  
892 Target papers helped bridge the gap between historical research and the latest developments, pro-  
893 viding a forward-looking perspective.

### 895 B.3 SURVEYS

896  
897 Surveys were treated as a separate category, as they provide an overview of the field and summa-  
898 rize key developments. For surveys, we included every subsection’s content and title, alongside  
899 corresponding references. In addition, we focused on identifying **notable research contributions**,  
900 defined as:

- 901 • Articles or works that were frequently cited.
- 902 • Papers described with extensive detail by the survey authors, often corresponding to key  
903 subheadings.

904  
905 These typically represent key historical works and important research results such as the develop-  
906 ment of technologies like Transformers.

907  
908 Surveys were instrumental in identifying key historical works (*key history*) that had a lasting impact  
909 on the field. These works were defined by their influence and the significant number of references  
910 made to them within the surveys.

911 More details of the data composition are shown in Table 6.

---

912 <sup>2</sup><https://arxiv.org/>

913 <sup>3</sup><https://ar5iv.labs.arxiv.org/>

918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971

Table 6: Details of the data information.

Entity Type	Content	Example
topic		"in-context learning"
year_start		"2021"
year_end		"2024"
source_paper ↓ reference(full)	name, arxiv_id, isAPA, abstract, reference, related_work	{ "name": "Chain-of-Thought...", "arxiv_id": "2201.11903", "isAPA": true, "abstract": "We explore how generating...", "reference": [ reference1, reference2,... ], "related work": "7Related Work...", "date": "2022" } "Subhro Roy and Dan Roth. 2015. Solving general arithmetic word problems. EMNLP"
target_list ↓ target_paper ↓ subtitle reference	name, arxiv_id, subtitles, reference, related_work	[target_paper1, target_paper2,...]  { "name": "Long-context LLMs...", "arxiv_id": "2404.02060", "subtitles": [subtitle1, subtitle2,...], "reference": [reference1, reference2,... ], "related_work": "2Related Work..." } "Reinforcement Learning via Supervised Learning (RvS)"  "Eva: Exploring the limits of masked visual representation learning at scale"
survey ↓ reference subtitle(full) ↓ key_history	name, arxiv_id, subtitles, all_references  name, key_history, refer- ences_in_this_section reference_title, key_word	{ "name": "In-context Learning...", "arxiv_id": "2401.11624", "subtitles": [subtitle1, subtitle1,...], "all_references": [reference1, reference2,...] } "Eva: Exploring the limits of masked visual representation learning at scale"  { "name": "Few-shot...", "key_history": [key_history1, key_history2,...] , "references_in_this_section": [reference1, reference2,...] }  { "reference_title": "Attention is all you need", "key_word": "Transformer Models" }
topic_history ↓ reference	name, arxiv_id, reference	{ "name": "Long-context LLMs...", "arxiv_id": "2404.02060", "reference": [reference1, reference2,...] } "Eva: Exploring the limits of masked visual representation learning at scale"

## 972 C PROMPTS

### 973 C.1 PROMPTS FOR GRAPH CONSTRUCTION

#### 974 **Extract from Abstract**

```

975 EXTRACT_THEME = '''Please extract the key issue addressed, the proposed
976 method, and the application domain from the abstract of the paper
977 titled *{title}* and present the information in the following JSON
978 format.
979
980 {{
981   "problem": {{
982     "name": the key issue it addressed,
983     "description": a more detailed description of this key issue
984   }}
985   "method": {{
986     "name": the method it proposed,
987     "description": a more detailed description of this method
988   }}
989   "domain": {{
990     "name": the application domain,
991     "description": a more detailed description of this domain
992   }}
993 }}]
994 If any of the information is not available, please fill the corresponding
995 value with 'null'. Note that the descriptions should be extracted
996 from context, DO NOT simply use your prior knowledge to complete them
997 .
998 Abstract Content: {abstract}'''

```

#### 997 **Extract from related works**

```

999 LEVEL1 = '''Please extract the method and problem entities related to the
1000 citation '{citations}' from the excerpt of the paper titled '{title}
1001 ', and identify the relations between these entities and the
1002 citation. Please respond with the following JSON format.
1003 [
1004   {{
1005     "entity name": The name of the entity that has relation
1006     with the citation '{citations}'. DO NOT extract human
1007     names as entities,
1008     "entity type": The type of the entity, selected from '
1009     method', 'problem', and 'domain',
1010     "description": Description of the entity extracted from
1011     the context, null if not enough information,
1012     "relation": The relationship between the citation and the
1013     entity extracted from the context can be expressed
1014     using phrases such as 'applied in', 'proposed by',
1015     and others. Ensure that the relation is explicitly
1016     mentioned in the text and avoid inferring any
1017     relations based on prior knowledge. Do not use vague
1018     description like 'related to'
1019   }},
1020   ...
1021 ]
1022 '''
1023 LEVEL2 = '''Find out the relationships between these entities in the
1024 content. DO NOT add relations including entities that do not exist in
1025 the list. Please respond with the following format.
1026 [{{
1027   "entity1": The name of the entity1,
1028   "relation": The relationship between entity1 and entity2. Ensure
1029   that the relation is explicitly mentioned in the text and
1030   avoid inferring any relations based on prior knowledge. Do
1031   not use vague description like 'related to',

```

```

1026     "entity2": The name of the entity2
1027     }},
1028     ...]
1029     Entities: {entities}
1030     Content: {content}'''

```

1031

## 1032 C.2 PROMPTS FOR REASONING

1033

### 1034 **Generate Related Works**

1035

```

1036 generate_relatedwork_prompt = f"""
1037 Let's generate a high-quality "Related Work" section for a research paper
1038 by following a structured reasoning approach. We will use the
1039 following steps to ensure clarity and depth.
1040 **Step 1: Analyze the topic and the narrative's progression.**
1041 The topic is '{topic}', and the subtitle is '{subtitle}'. Here is the
1042 time-based progression of research developments:\n\n{cot_narrative}\n
1043 \n
1044 Analyze the key themes, shifts, and milestones in the narrative to
1045 extract the most relevant and impactful works that shaped the field
1046 over time.
1047 **Step 2: Identify key studies.**
1048 Based on the analysis, identify the most influential and representative
1049 studies that have contributed to the advancement of this field.
1050 Select works that either introduced foundational concepts, solved
1051 critical challenges, or advanced the field in significant ways.
1052 **Step 3: Structure the Related Work section.**
1053 Organize the selected studies in a way that emphasizes their contribution
1054 to the progression of the field. The section should naturally flow
1055 either chronologically or thematically, ensuring a balance between
1056 foundational works and recent innovations. You may highlight any gaps
1057 or ongoing debates in the literature to contextualize how these
1058 works relate to your research.
1059 Now, based on this reasoning, generate the "Related Work" section. Make
1060 sure it is flexible but retains a coherent narrative that aligns with
1061 academic standards. Incorporate key research areas and their
1062 evolution in the field, using a mix of foundational works and recent
1063 studies. The section should demonstrate a clear understanding of how
1064 these works interrelate and how they contribute to the current
1065 research landscape.
1066 Please provide the response in "Related Work" section only, structured as
1067 follows:
1068 """
1069 example="""
1070 **Related Work**
1071 The field of {main_topic} has evolved significantly over the past few
1072 decades, particularly in areas such as {key_areas_1}, {key_areas_2},
1073 and {key_areas_3}. Early works such as {Author1 et al., Year} laid
1074 the groundwork for {specific concept or technique}, introducing key
1075 methods that have since been built upon by later studies.
1076 For instance, *{Key Area 1}* has been a major focus, starting with
1077 foundational research by {Author2 et al., Year}, who proposed {a
1078 major contribution}. Building on this, subsequent studies like {
1079 Author3 et al., Year} have refined these approaches, introducing
1080 innovations such as {specific advancement} that have made a
1081 substantial impact in the field.
1082 In contrast, *{Key Area 2}* represents a more recent development, with
1083 groundbreaking contributions by {Author4 et al., Year}, who explored
1084 {a novel approach or finding}. This has opened new avenues for
1085 research, particularly in {specific application or challenge}, as
1086 evidenced by {Author5 et al., Year}, whose work has further expanded
1087 on these ideas.
1088 Additionally, the intersection of *{Key Area 3}* with {related field} has
1089 also gained attention in recent years. Notably, {Author6 et al.,

```

```

1080     Year} demonstrated {key contribution}, which has been instrumental in
1081     advancing the understanding of {specific problem or question}.
1082 While much progress has been made, there remain open questions,
1083     especially in {specific area of ongoing research}, where recent
1084     studies by {Author7 et al., Year} indicate that further exploration
1085     is needed to fully realize the potential of {key technique or
1086     approach}.
1087 By reviewing these works, we gain a comprehensive understanding of how
1088     the field has evolved and where it is headed, providing essential
1089     context for the contributions of our own research.
1090 ""
1091 generate_relatedwork_prompt+=example
1092
1093 Generate Future Idea
1094
1095 Now, let's think step by step to generate predictions for future
1096     research directions based on the provided time-based narrative and
1097     subtitle.
1098
1099 Step 1: Analyze the time-based narrative for future trends. The
1100     narrative '{cot_narrative}' shows how the research has evolved over
1101     time. Carefully examine this narrative to extract clues about current
1102     trends, technological bottlenecks, and potential gaps in research
1103     that could drive future developments.
1104
1105 Step 2: Consider challenges and opportunities. Based on the trends
1106     and patterns identified in the narrative, think about the challenges
1107     the field currently faces and the opportunities for future
1108     innovations. What are the key bottlenecks, and what cutting-edge
1109     technologies could overcome them?
1110
1111 Step 3: Predict future research directions. Based on the analysis,
1112     predict possible future directions for the topic '{topic}' and
1113     subtitle '{subtitle}'. These directions should be logically derived
1114     from the observed research trends and potential advancements in
1115     technology.
1116
1117 Step 4: Structure the future directions and technical roadmap.
1118     Organize the predicted future research directions into a well-
1119     structured roadmap, clearly outlining the steps researchers might
1120     take to advance in this area. Present the future directions in JSON
1121     format.
1122 Please provide the response in JSON format only, structured as follows:
1123 Example output format:
1124 ```json
1125 {{
1126   "Future_Directions": {{
1127     "1. Title of Future Direction": {{
1128       "Description": "Detailed description of the future research
1129         direction, derived from trends in the narrative.",
1130       "Technical_Roadmap": [
1131         "First step in technical roadmap, based on observed trends.",
1132         "Next steps, reflecting future possibilities derived from the
1133         narrative."
1134       ]
1135     }},
1136     "2. Another Future Direction": {{
1137       "Description": "Another future research direction logically derived
1138         from current research challenges and gaps.",
1139       "Technical_Roadmap": [
1140         "First step, addressing challenges seen in the narrative.",
1141         "Subsequent steps reflecting the roadmap towards technological
1142         advancements."
1143       ]
1144     }}
1145   }}
1146 }}

```

```

1134   }}
1135   }}
1136   ```

```

### 1138 C.3 PROMPTS FOR EVALUATION

#### 1139 C.3.1 PROMPT FOR FUTURE RELIABILITY EVALUATION

##### 1141 **$S_1$ : Semantic Similarity**

1143 The following is the abstract of a target paper:\n\n{target\_abstract}\n\n  
 1144 Below is a set of future research directions predicted by another  
 1145 paper:\n\n{future\_text}\n\n

1146 Your task is to carefully compare the abstract and the predicted  
 1147 future directions step by step using chain of thought reasoning.

1148 Step 1: Analyze the future directions and extract key research themes  
 1149 or topics.

1150 Step 2: Compare each theme with the abstract to identify any matching  
 1151 concepts.

1152 Step 3: Assign a score based on the extent of matching as follows:

- 1153 - **Score 0**: No matches at all; none of the key themes or topics  
 1154 from the future directions are present in the abstract.
- 1155 - **Score 1**: Very few matches; only 1 out of 5 key themes match  
 1156 with the abstract.
- 1157 - **Score 2**: Some matches; 2 out of 5 key themes match with the  
 1158 abstract.
- 1159 - **Score 3**: Moderate matches; 3 out of 5 key themes match with the  
 1160 abstract.
- 1161 - **Score 4**: Mostly matches; 4 out of 5 key themes match with the  
 1162 abstract.
- 1163 - **Score 5**: All matches; all 5 key themes from the future  
 1164 directions are present in the abstract.

1165 Please **only return the final score as a single number**.

##### 1166 **$S_2$ : Innovation and Feasibility**

1167 The following is a future research direction proposed by a research paper  
 1168 :

1169 "{future\_direction}"  
 1170

1171 Please evaluate the following aspects step by step:

1172 Step 1: Analyze the future direction to determine if it proposes a  
 1173 novel or innovative idea compared to current research in the  
 1174 field, specifically related to the topic: {topic} and subtitle: {  
 1175 subtitle}. Does it introduce new concepts, techniques, or  
 approaches that are not commonly explored?

1176 Step 2: Assess the technical feasibility of this future direction.  
 1177 Can it be realistically implemented with current technology, or  
 1178 does it require significant breakthroughs?

1179 Step 3: Rate the future direction on a scale from 1 to 5:

- 1180 - 1: No innovation and technically infeasible.
- 1181 - 2: Slight innovation but mostly infeasible.
- 1182 - 3: Moderately innovative and feasible with technical challenges.
- 1183 - 4: Innovative and feasible with current technology, with minor  
 1184 challenges.
- 1185 - 5: Highly innovative and technically feasible without significant  
 1186 challenges.

1187 Please **only return the final score as a single number**.

##### **$S_3$ : Temporal Consistency**

1188 The following is a future research direction proposed by a research paper  
1189 :

1190  
1191 "{future\_direction}"

1192 Please evaluate the following step by step:

1193 Step 1: Analyze the current state of research and technological  
1194 progress in the relevant field of {topic}, and related to the  
1195 subtitle: {subtitle}. Identify the key milestones and major  
1196 developments up to now.

1197 Step 2: Determine if this future direction logically builds upon  
1198 recent developments, or if it requires an unrealistic leap  
1199 forward in technology.

1200 Step 3: Assess whether the future direction aligns with the current  
1201 pace of technological development. If it seems unrealistic for  
1202 the near future, explain why.

1203 Step 4: Rate the temporal consistency of the future direction on a  
1204 scale from 1 to 5:

- 1205 - 1: Does not fit the timeline at all.
- 1206 - 2: Slightly inconsistent with the timeline.
- 1207 - 3: Moderately consistent with the timeline, with some gaps.
- 1208 - 4: Largely consistent with minor inconsistencies.
- 1209 - 5: Fully consistent with the timeline and logically follows current  
1210 research progress.

1211 Please **only** return the final score as a single number.

1212

#### 1213 **S<sub>4</sub>: Contextual Consistency**

1214

1215 The following is a future research direction proposed by a research paper  
1216 :

1217 "{future\_direction}"

1218 The target paper's abstract is as follows:

1219 "{target\_abstract}"

1220 Please evaluate the following aspects step by step:

1221 Step 1: Identify the key research challenges or limitations discussed  
1222 in the target paper's abstract, which relates to the topic: {  
1223 topic} and subtitle: {subtitle}. What are the primary issues the  
1224 target paper seeks to address?

1225 Step 2: Determine if the proposed future direction addresses any of  
1226 these challenges or builds upon the research presented in the  
1227 target paper.

1228 Step 3: Assess whether the proposed future direction logically  
1229 follows the research context or is disconnected from the  
1230 challenges identified in the target paper.

1231 Step 4: Rate the contextual consistency of the future direction on a  
1232 scale from 1 to 5:

- 1233 - 1: Completely disconnected from the research context.
- 1234 - 2: Slightly relevant, but mostly misaligned with the research  
1235 context.
- 1236 - 3: Moderately related to the research context, but missing key  
1237 connections to the identified challenges.
- 1238 - 4: Largely consistent, addressing most of the research challenges  
1239 with minor gaps.
- 1240 - 5: Fully aligned with the research context, addressing key  
1241 challenges comprehensively.

1242 Please **only** return the final score as a single number.

1243

### 1242 C.3.2 WEIGHT SETTINGS FOR FUTURE RELIABILITY EVALUATION

1243  
1244 **Weight Distribution:** In our future prediction evaluation task, we use the following weights:

1245 ( $w_1 = 0.4$ ) for historical completeness and prediction reliability (this includes manual evaluation),  
1246 ( $w_2, w_3, w_4 = 0.2$ ) each for the LLM-based assessments of prediction reliability, text generation  
1247 quality, and other factors. We assign the highest weight ( $w_1 = 0.4$ ) to historical completeness  
1248 because this portion involves manual evaluation, which we believe is more accurate and reliable,  
1249 particularly for complex tasks such as evaluating the completeness of the graph and the reliability  
1250 of future predictions. Manual evaluation offers higher credibility compared to the more automated  
1251 LLM assessments.

1252 On the other hand, the LLM-based evaluations are given lower weights because they rely on auto-  
1253 mated models, and their results can be more dynamic and fluctuate depending on the context, input,  
1254 and other factors. While LLM assessments are valuable for scalability, we acknowledge that their  
1255 results are not as stable or trustworthy as human assessments in these particular tasks.

1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295



## D RESULT SPECIFICS

### D.1 HISTORICAL SUMMARIZATION EVALUATION RESULTS

The specific values for our evaluation metrics of the nine topics can be viewed below:

Table 7: The performance of history completeness and history summarization. We utilize the overlap degree of paper nodes in the graph and citations in the survey to evaluate the completeness of history. The search phase is a preliminary process of history summarization, we show the search efficiency on both overall and key citations to indicate the reasoning performance. The history summarization performance is assessed on the overlap degree of generated content and the target survey. See the definition of the metrics in Section 3.

Evaluation Object	Graph Completeness(%)		Search Efficiency(%)		Generated Content Overlap Degree(%)
	$O_R$	$O_H$	$SE_R$	$SE_H$	
Topic 1	50.00	65.74	75.71	74.51	42.86
Topic 2	30.00	46.00	63.00	79.00	27.32
Topic 3	70.00	73.00	43.00	36.00	2.27
Topic 4	57.00	38.00	81.00	68.00	24.24
Topic 5	34.00	59.00	77.00	73.00	27.93
Topic 6	21.00	35.00	17.00	19.00	5.00
Topic 7	36.00	50.00	97.00	100.00	35.57
Topic 8	25.00	25.00	79.00	90.00	19.75
Topic 9	31.13	25.40	96.54	100.00	33.33
Average	39.35	46.35	69.92	71.06	24.25

### D.2 CASE STUDY ON FUTURE PREDICTION: TOPIC 2 - LLMs FOR RECOMMENDATION

#### Future Directions Generated by Our *HorizonAI*

The ratings are 4 for  $S_1$ , 4 for  $S_2$ , 4.17 for  $S_3$ , 3.75 for  $S_4$ , with an average of 3.98.

- Enhancing Generative Recommendations with LLMs** *Description:* Building on the paradigm-shifting survey on generative recommendations, focus on overcoming efficiency challenges and improving the ethical integration of LLMs in recommender systems. *Technical Roadmap:* 1. Develop optimized algorithms to reduce computational overhead in LLM-based generative recommendations. 2. Explore pruning and quantization techniques to make LLMs more efficient for real-time applications. 3. Conduct extensive studies on the ethical implications of LLM recommendations, formulating guidelines for ethical usage.
- Integrating User Mobility Patterns in Recommendations** *Description:* Leverage LLMs to better understand and predict user mobility patterns, enhancing recommendation relevance in location-based services. *Technical Roadmap:* 1. Create datasets that capture diverse mobility patterns and integrate them with user activity data. 2. Develop LLM-based models that can learn and predict user mobility patterns to provide context-aware POI recommendations. 3. Evaluate the effectiveness of mobility-pattern integrated recommendations in real-world scenarios and refine models accordingly.
- Cross-Domain Knowledge Transfer in Recommender Systems** *Description:* Address the limitations of single-domain pre-training by exploring techniques for cross-domain knowledge transfer using LLMs. *Technical Roadmap:* 1. Investigate methodologies for transferring knowledge between different domains to enhance recommendation capabilities. 2. Develop frameworks to assess the quality and applicability of transferred knowledge across domains. 3. Implement and test cross-domain transfer models in various recommendation contexts to determine best practices.
- Enhancing Behavioral Simulation for Recommendation Development** *Description:* Refine the agent-based simulation methods introduced by RecAgent to create even more sophisticated user behavior models. *Technical Roadmap:* 1. Expand user and recommender modules to capture a broader range of user behaviors and interac-

1350  
1351  
1352  
1353  
1354  
1355  
1356  
1357  
1358  
1359  
1360  
1361  
1362  
1363  
1364  
1365  
1366  
1367  
1368  
1369  
1370  
1371  
1372  
1373  
1374  
1375  
1376  
1377  
1378  
1379  
1380  
1381  
1382  
1383  
1384  
1385  
1386  
1387  
1388  
1389  
1390  
1391  
1392  
1393  
1394  
1395  
1396  
1397  
1398  
1399  
1400  
1401  
1402  
1403

tions. 2. Integrate reinforcement learning techniques to continuously improve simulated user-adaptive responses. 3. Validate simulated environments against real-world data to ensure the fidelity and applicability of simulations.

5. **Ethical and Fair Recommendation Systems** *Description*: Further research into developing frameworks for ensuring fairness, transparency, and ethical considerations in LLM-powered recommender systems. *Technical Roadmap*: 1. Develop metrics and benchmarks to evaluate fairness and transparency in recommendation systems. 2. Create auditing tools to detect and mitigate biases in LLM-generated recommendations. 3. Formulate policies and best practices to uphold ethical standards in deploying LLM-based recommenders.

6. **Real-Time Adaptability in LLM-Based Recommender Systems** *Description*: Focus on making LLM-based recommender systems more adaptive in real-time scenarios to provide timely and context-sensitive recommendations. *Technical Roadmap*: 1. Develop techniques for rapid model updates and fine-tuning based on real-time user feedback and environmental changes. 2. Implement scalable infrastructure to support real-time adaptability without compromising performance. 3. Test real-time adaptable systems in diverse recommendation applications and refine models for robustness.

#### Future Directions Generated by GPT-4o

The ratings are 3.6 for  $S_1$ , 3.3 for  $S_2$ , 3.8 for  $S_3$ , 2.98 for  $S_4$ , with an average of 3.42.

1. **Integrating LLMs with Geographic Information Systems (GIS) for Next-POI Recommendation** *Description*: Develop sophisticated models that combine the power of LLMs with GIS to provide highly accurate next-POI recommendations, leveraging spatial and temporal data. *Technical Roadmap*: 1. Develop methods to preprocess and integrate GIS data with LLM inputs, ensuring data coherence and compatibility. 2. Enhance LLM capabilities to understand geographical nuances and contexts by fine-tuning on spatial datasets. 3. Create hybrid models that combine LLM predictions with GIS analytics, implementing validation mechanisms specific to location-based predictions. 4. Deploy and iteratively improve the recommendation system using real-world data and feedback loops.
2. **Personalization in LLM-Based Recommender Systems** *Description*: Address personalization challenges by developing LLMs capable of deeply understanding individual user preferences and historical behavior for tailor-made recommendations. *Technical Roadmap*: 1. Aggregate extensive user interaction data and develop pre-training tasks that reflect typical user behaviors. 2. Investigate transfer learning methods to adapt general LLMs to specific user profile data, enhancing personalization accuracy. 3. Integrate reinforcement learning techniques to continuously adapt and improve recommendations based on real-time user feedback. 4. Evaluate and implement privacy-preserving mechanisms to ensure user data protection while personalizing recommendations.
3. **Predicting User Mobility Patterns Using LLMs** *Description*: Leverage LLMs to predict complex user mobility patterns by integrating various data sources like social media, transportation logs, and historical movements. *Technical Roadmap*: 1. Identify and unify diverse data sources (e.g., social media, GPS logs) to create comprehensive mobility datasets. 2. Design pre-training tasks focused on mobility pattern recognition to enhance LLM understanding of movement data. 3. Implement sequence-based LLM architectures to capture temporal dynamics and predict future movements accurately. 4. Develop evaluation frameworks to assess the predictive performance and refine models based on predictive accuracy and reliability metrics.
4. **Predicting User Mobility Patterns Using LLMs** *Description*: Leverage LLMs to predict complex user mobility patterns by integrating various data sources like social media, transportation logs, and historical movements. *Technical Roadmap*: 1. Identify and unify diverse data sources (e.g., social media, GPS logs) to create comprehensive mobility datasets. 2. Design pre-training tasks focused on mobility pattern recognition to enhance LLM understanding of movement data. 3. Implement sequence-based

1404  
1405  
1406  
1407  
1408  
1409  
1410  
1411  
1412  
1413  
1414  
1415  
1416  
1417  
1418  
1419  
1420  
1421  
1422  
1423  
1424  
1425  
1426  
1427  
1428  
1429  
1430  
1431  
1432  
1433  
1434  
1435  
1436  
1437  
1438  
1439  
1440  
1441  
1442  
1443  
1444  
1445  
1446  
1447  
1448  
1449  
1450  
1451  
1452  
1453  
1454  
1455  
1456  
1457

LLM architectures to capture temporal dynamics and predict future movements accurately. 4. Develop evaluation frameworks to assess the predictive performance and refine models based on predictive accuracy and reliability metrics.

5. **Cross-Domain Implications of LLM-Based Recommendation Systems** *Description*: Explore the applicability and implications of LLM-based recommender systems across various domains (e.g., retail, entertainment, health) to uncover new opportunities and challenges. *Technical Roadmap*: 1. Conduct domain-specific studies to understand the unique requirements and constraints of LLM applications in different sectors. 2. Develop adaptable LLM architectures that can efficiently switch contexts and deliver domain-specific recommendations. 3. Implement cross-domain transfer learning techniques to enhance LLM generalizability while preserving domain-specific nuances. 4. Continuously monitor and document the performance, ethical considerations, and user satisfaction across these diverse applications.