

TEST OF TIME: A BENCHMARK FOR EVALUATING LLMs ON TEMPORAL REASONING

Anonymous authors

Paper under double-blind review

ABSTRACT

Large language models (LLMs) have showcased remarkable reasoning capabilities, yet they remain susceptible to errors, particularly in temporal reasoning tasks involving complex temporal logic. Existing research has explored LLM performance on temporal reasoning using diverse datasets and benchmarks. However, these studies often rely on real-world data that LLMs may have encountered during pre-training or employ anonymization techniques that can inadvertently introduce factual inconsistencies. In this work, we address these limitations by introducing novel synthetic datasets specifically designed to assess LLM temporal reasoning abilities in various scenarios. The diversity of question types across these datasets enables systematic investigation into the impact of the problem structure, size, question type, fact order, and other factors on LLM performance. Our findings provide valuable insights into the strengths and weaknesses of current LLMs in temporal reasoning tasks. To foster further research in this area, we will open-source the datasets and evaluation framework used in our experiments.

1 INTRODUCTION

Recent breakthroughs in large language model (LLM) research and applications have been significant (Vaswani et al., 2017; Devlin et al., 2019; Raffel et al., 2020; Brown et al., 2020; Touvron et al., 2023; Achiam et al., 2023; Team et al., 2023; Reid et al., 2024). These models, capable of generating new content, have fascinated the AI community, leading to the release of numerous LLMs trained on diverse tasks and data types (Zhao et al., 2023). All of these advancements have led to a growing consensus that LLMs are a pivotal advancement on the path to artificial general intelligence (AGI) (Bubeck et al., 2023). Benchmarking reasoning capabilities in LLMs is therefore a problem of pressing interest to the field (Huang & Chang, 2023).

In this work, we focus on temporal reasoning, an essential task for intelligent systems across many domains. Temporal reasoning is focused on understanding reasoning between events in time. Despite this area’s importance, existing temporal reasoning benchmarks do not effectively measure the full scope of temporal reasoning relationships. Instead, they typically rely on question-answering tasks based on Knowledge Graph (KG)-style temporal facts about well-known entities.

This overemphasis on KG-style temporal facts limits the scope of research and creates several issues. First, it neglects the diverse temporal structure and reasoning tasks found in real-world applications beyond KGs. Second, the results on such data often reflect a model’s ability to exploit prior knowledge rather than genuine temporal reasoning, making findings less relevant to domains where models lack this knowledge (see Figure 1 as an example.). In addition, previous research has shown that shortcuts or heuristics can often answer questions on these datasets without explicit temporal reasoning (Chen et al., 2022; Tan et al., 2023). Finally, the simple temporal structure of knowledge graphs overlooks the extensive time arithmetic skills required in real-world temporal questions.

Our Contributions: To address these limitations, we develop tasks specifically designed to assess temporal reasoning in a more comprehensive and controlled manner. Our benchmark, Test of Time, `TO_T`, is centered around the observation that temporal reasoning often involves two primary skills: 1) understanding the semantics and logic of time, and 2) the ability to carry out temporal arithmetic. `TO_T` has two tasks to cover each essential skill of temporal reasoning, which enable measuring and improving model performances along these two axes independently. `TO_T-Semantic`, a synthetic task, focuses on temporal semantics and logic; it allows for a flexible exploration of diverse graph

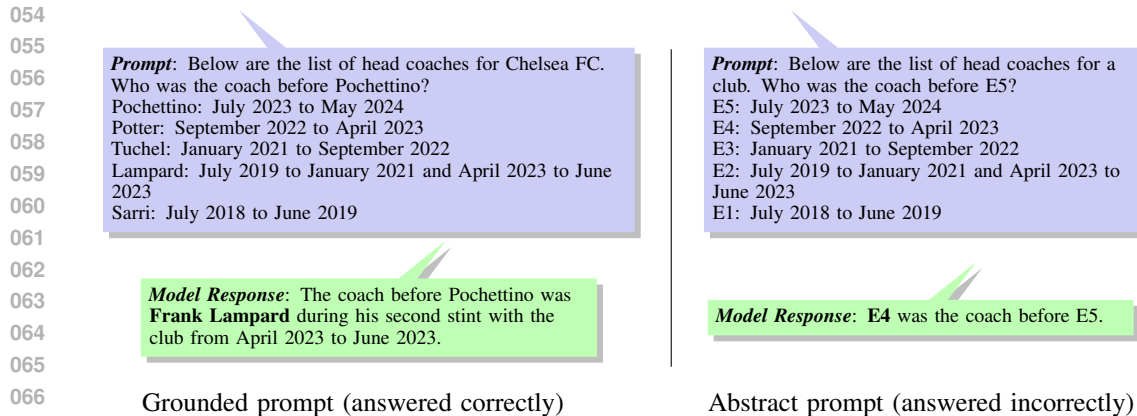


Figure 1: Comparison of the same temporal query using real (left) and anonymized (right) entity names. Gemini Advanced correctly answered the query with real names but failed with anonymized names, suggesting that LLMs might rely on their parametric knowledge to solve temporal tasks.

structures and reasoning task complexity, isolating and evaluating reasoning abilities independent of prior knowledge. `TOT-Arithmetic`, a crowd-sourced task, assesses the ability to perform calculations involving time points and durations. Our experimental results with `TOT` provide valuable insights into the strengths and weaknesses of current LLMs in temporal reasoning tasks.

2 RELATED WORK

Reasoning. The ability to draw valid conclusions from explicitly provided knowledge has been a fundamental goal for AI since its early days (McCarthy, 1959; Hewitt, 1969). In the past few years, several LLM-based techniques have been developed which have advanced the general automated reasoning capabilities of the state-of-the-art models (Wei et al., 2022; Yao et al., 2023), or their capabilities in specific directions including mathematical reasoning (Lewkowycz et al., 2022; Ahn et al., 2024), logical reasoning (Creswell et al., 2022; Kazemi et al., 2023b), multi-modal reasoning (Wang et al., 2024), commonsense reasoning (Zellers et al., 2019), and more. Advancing reasoning may explicitly or implicitly translate to improvements in several downstream NLP applications.

Temporal reasoning. Temporal reasoning has recently gained substantial attention (e.g., Vashishtha et al., 2020; Nylund et al., 2023; Hu et al., 2023; Gurnee & Tegmark, 2023; Liu et al., 2023; Xiong et al., 2024; Beniwal et al., 2024; Jia et al., 2024). Much research focuses on enhancing LLMs’ understanding of temporal concepts, primarily through pre-training and fine-tuning strategies to improve their temporal reasoning capabilities (e.g., Ning et al., 2019; Zhou et al., 2020; Yang et al., 2023; Xiong et al., 2024; Jia et al., 2024).

Benchmark creation is another active area, with many benchmarks centered on knowledge graphs (e.g., Jia et al., 2018; Neelam et al., 2021; Jia et al., 2021; Wang & Zhao, 2023; Chu et al., 2023; Su et al., 2024). While TempTabQA (Gupta et al., 2023) offers crowd-sourced questions based on Wikipedia infoboxes, the process is resource-intensive and prone to issues like LLM overuse by workers. The questions in Wang & Zhao (2023) are all multiple-choice, and do not require reasoning through a many temporal facts from a knowledge graph. The questions in Chu et al. (2023) are collected from ten existing real-world datasets, one of which requires reasoning through temporal facts provided in the context. In contrast, ToT goes beyond such datasets by providing controllable, comprehensive temporal relationship collections via synthetic graph generation. The questions in Timo Su et al. (2024) are grouped into two categories: math-time and pure-time. ToT-Arithmetic covers more domains in the math-time category and more focus on reasoning in the pure-time category. Xiong et al. (2024) recently proposed TGQA, a data set derived from the YAGO11k knowledge graph (Dasgupta et al., 2018). To prevent data leakage, TGQA changes each entity name to a name generated by GPT3.5 that is guaranteed to (i) align with the entity’s type and (ii) not be otherwise present in YAGO11k. This strategy has several weaknesses. First, it can introduce spurious entity name

Table 1: Comparison of T_{OT} against related benchmarks.

Benchmark	Semantics	Arithmetic	Real-World	Synthetic	Hermetic	Implicit
TimeSensitiveQA (Chen et al., 2021)	✓	✗	✓	✗	✗	✗
StreamingQA (Liska et al., 2022)	✓	✗	✓	✗	✗	✗
TempLama (Dhingra et al., 2022)	✓	✗	✓	✗	✗	✗
TEMPTABQA (Gupta et al., 2023)	✓	✗	✓	✗	✗	✓
TEMPREASON (Tan et al., 2023)	✓	✓	✓	✗	✗	✓
TIQ (Jia et al., 2024)	✓	✗	✓	✗	✗	✓
TempUN (Beniwal et al., 2024)	✓	✗	✓	✗	✗	✗
TGQA (Xiong et al., 2024)	✓	✗	✓	✗	✗	✗
T_{OT} (ours)	✓	✓	✓	✓	✓	✓

correlations (LLMs could even potentially guess the original entities due to their adjacent relations). Second, it can generate factually incorrect or anti-commonsensical claims, for instance, if an entity’s generated replacement name is a *real* name that is nonetheless not in YAGO11k. Finally, relying on GPT for copying facts introduces the potential for hallucinations to contaminate the dataset.

Synthetic datasets. A new trend in probing various LLMs capabilities, especially in the case of reasoning, is through synthetic data that allows for a more systematic evaluation. Previous work has developed synthetic datasets for probing and improving various kinds of reasoning including logical reasoning (Tafjord et al., 2021; Kazemi et al., 2023c; Saparov et al., 2023) and mathematical reasoning (Kazemi et al., 2023a; Srivastava et al., 2024). Most similar to our work, Fatemi et al. (2024) develop a synthetic probe for measuring the graph-based reasoning abilities of LLMs (Sanford et al., 2024; Perozzi et al., 2024). Our work extends this concept to the case of temporal reasoning with graph-like facts.

Present work. In this work, we introduce T_{OT} , a novel benchmark for temporal reasoning generated synthetically. Unlike many existing benchmarks that rely on knowledge graphs, T_{OT} aims to encompass a wider variety of graph structures. Our synthetic generation approach offers precise control over the type of data produced. Importantly, when evaluating LLMs against T_{OT} , they cannot exploit their latent knowledge for shortcuts; instead, they must genuinely reason with the presented facts. This design promotes a more rigorous assessment of temporal reasoning capabilities in LLMs. Table 1 provides a comprehensive comparison of T_{OT} with existing benchmarks across six key dimensions: **1- Semantics:** whether the benchmark has semantic-type questions, **2- Arithmetic:** whether the benchmark has arithmetic-type questions, **3- Real-world:** whether the benchmark has questions generated from real-world data, **4- Synthetic:** whether the benchmark has questions generated from synthetic data, **5- Hermetic:** whether the benchmark is sealed off from potential LLM training data, and **6- Implicit:** whether the benchmark includes implicit questions. Our analysis reveals that T_{OT} is unique in incorporating all these question types while effectively mitigating training data leakage. Notably, TEMPREASON (Tan et al., 2023) only covers one category of the arithmetic operations as defined in Section 3.2.

3 T_{OT} : A BENCHMARK FOR EVALUATING LLMs ON TEMPORAL REASONING

We propose that effective temporal reasoning hinges on two distinct skills: understanding the semantics and logic of time, and performing accurate temporal arithmetic. To measure and improve model performance along these independent axes, we create a dedicated task for each skill. By decoupling the evaluation of temporal semantics from arithmetic, we aim to provide a more nuanced analysis of LLM capabilities, pinpointing strengths and weaknesses in each aspect. Experiments on these tasks enable us to independently benchmark LLM performance on both dimensions.

3.1 ToT-SEMANTIC: A SYNTHETIC DATASET

The first task we introduce, ToT-Semantic, consists of synthetic problems designed to highlight temporal semantics and logic in reasoning. This task is unique because it allows us to freely experiment with a wide range of temporal dependencies and manipulate the difficulty of the reasoning problem. This allows us to isolate and analyze the core reasoning capabilities of an LLM, separating them from any reliance on pre-existing parametric knowledge. To create the ToT-Semantic task, we take the following steps (summarized in Figure 2):

Step 1: Generate a (random) graph. We begin by generating random structures that we will then use to create temporal questions. To ensure we generate a diverse set of random structures for this purpose, we turn to the literature on graph structure generation. From it, we employ several existing algorithms for generating graphs of varying properties. This includes Erdős-Rényi (ER) graphs (Erdős & Rényi, 1959), scale-free networks (SFN) (Barabási & Albert, 1999), graphs following the Barabási-Albert (BA) model (Albert & Barabási, 2002) and stochastic block model (SBM) (Holland et al., 1983), as well as star and complete graphs. Each of these graph generation algorithms exhibits different properties and correspond to graphs that appear in different applications. For instance, Erdős-Rényi graphs are typically sparse with low average degree, while Barabási-Albert graphs are dense and exhibit power-law degree distributions. We leverage the NetworkX library for generating our random graphs. Additionally, we extracted anonymized EgoNets from WikiData (Vrandečić & Krötzsch, 2014) by replacing the entity and relation names with generic names. We refer to this structure as *Anonymized Wikidata Extract (AWE)* in our experiments. We generate graphs with the number of nodes selected uniformly at random from the [5-30] interval. More details on the random graph generators used (with visualizations) are available in Appendix A.

Step 2: Assigning entity and relation names. Once we have an initial graph structure, we assign names to the nodes and relations to the edges. For each graph, we first decide a number of relation types to be assigned to the edges, and assign each of these relation types to one of one-to-one, one-to-many, many-to-one and many-to-many. Then, for each edge in the graph, we randomly assign between 1 to p ($=3$ in our experiments) relations types without violating the relation type arity.

Step 3: Generate temporal facts. Then, for each edge (u, v) labeled with a relation r , we assign a valid time interval $[t1, t2]$ that respects the relation types, and turn the tuple $(u, v, r, t1, t2)$ into a textual temporal fact using a template.

Step 4: Question generation. Having generated the random graphs, we then create questions about those graphs. We consider eight types of questions that are frequently used in day-to-day life and are common in various benchmarks. **EventAtTimeT**: asking which entity had some relation R with entity E at some T; **EventAtWhatTime**: asking at what time a relation R between two entities E1 and E2 started/ended; **NumEventsInTimeInterval**: asking how many entities had relation R with entity E between T1 to T2; **BeforeAfter**: asking which entity had relation R with E1 right before/after E2; **EventAtTimeOfAnotherEvent**: asking when E1 had relation R1 with E2, which entity had relation R2 with E3; **FirstLast**: asking which entity was the first to have relation R with E; **RelationDuration**: Asking the k-th time relation R happened between E1 and E2, how long did it last; and **timeline**: Asking to sort the entities that had relation R with E chronologically.

To create any of the above questions, we keep sampling graphs and fact(s) from the graph until a proper question of the desired type can be created for that graph and for that fact. For example, to create a *BeforeAfter* question, we keep sampling a graph G and fact $F = (S, R, O, T1, T2) \in G$ until we have a case where there is a unique entity E that was the R of O right before $[T1, T2]$.

Following the above two steps, we generated 10 questions per graph generation and per question type. We sorted the facts in five different ways as will be discussed later. This gives as a benchmark with a total of $7 * 8 * 5 * 10 = 2800$ questions, where 7 is the number of graph generation algorithms, 8 is the number of question types, 5 is the number of sorting algorithms, and 10 is the number of samples we generated. Example questions of each category type are shown in Table 2.

3.2 ToT-ARITHMETIC: A TEMPORAL ARITHMETIC DATASET

Our second task, ToT-Arithmetic, shifts from synthetic data to a real-world focus. This task moves beyond understanding the logic and semantics of time and delves into the practical application

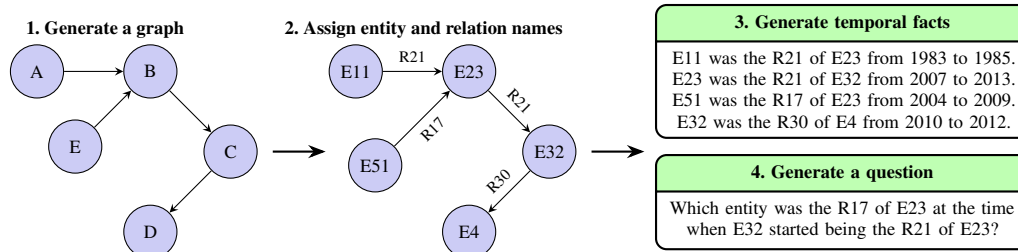


Figure 2: Steps for creating the ToT-Semantic dataset.

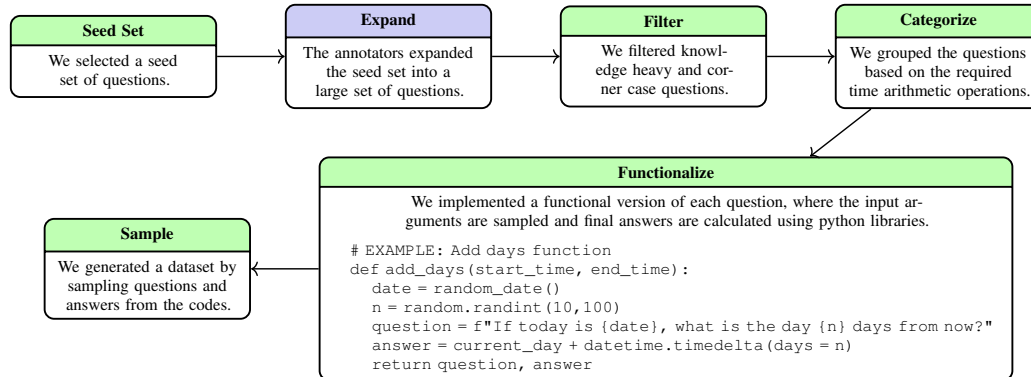


Figure 3: Steps for creating the ToT-Arithmetic dataset. The green and blue colors represent the operations done by the authors and the annotators respectively.

of mathematical operations within a temporal context. Through it, we are able to measure an LLM’s proficiency in temporal arithmetic and its practical utility in handling time-related computations.

To create a large time-arithmetic dataset that covers a wide variety of problems, we took the steps summarized in Figure 3. We explain each step in more detail below.

- **Seed Set:** By examining the existing benchmarks and the kind of temporal arithmetic questions that arise in them and through searching the web, we gathered a small set of initial questions.
- **Expand:** We presented our seed set to 15 annotators who were tasked to propose either new time arithmetic questions that were not in our seed set, or to provide questions corresponding to other scenarios or question templates where one requires to do similar time arithmetic operations to one of the questions in our seed set. We gathered a large list of questions through this process.

Table 2: Example for each question type in the ToT-Semantic dataset.

Question Type	Example
EventAtTimeT	Find the entity that evidently was the R17 of E69 in year 1932.
EventAtWhatTime	At what time did E69 start being the R90 of E22?
NumEventsInTimeInterval	Find the number of unique entities that were the R82 of E27 between 1952 to 1957. Relations that ended in 1952 or started in 1957 must be counted.
BeforeAfter	Immediately before E59, which entity was the R20 of E6?
EventAtTimeOfAnotherEvent	E94 was the R82 of which entity at the time when E83 started being the R20 of E59?
FirstLast	Which entity was the first that was the R35 of E91?
RelationDuration	When E24 was the R53 of E11 for the 2nd time, for how many years did the relation last? The duration can be computed by subtracting the start time from the end time.
Timeline	Which entities were the R17 of E69?

- **Filter:** We manually went through all the questions and filtered the ones that were focusing on corner cases, or that required extensive knowledge (*e.g.*, requiring to memorize the entire calendar).
- **Categorize:** We then grouped the remaining problems into seven categories, shown with examples in Table 3. Categories are formed based on the time arithmetic operations required, as follows: **AddSubtract:** adding or subtracting a number (corresponding to days, weeks, minutes, hours, etc.) from a date or time; **Compare:** comparing a number of dates/times provided in different formats chronologically; **Duration:** computing the difference between two dates/times; **Schedule:** finding mutual free spots within multiple blocked times; **Timezone:** involving dealing with different timezones; **Trick:** involving questions with slight twists; and **MultiOp:** involving questions where multiple of the above operations are needed.
- **Functionalizing:** Following (Srivastava et al., 2024), we implemented a functional version of each problem to enable sampling different values for each question and solving based on those values. A functional version of one of our simple problems is provided in Figure 3.
- **Sampling:** We then sampled questions and answers from our functionalized problems. We made the number of samples proportional to the number of different problems that fell under each category. Specifically, we sampled 350 for AddSubtract, 350 for Compare, 200 for Duration, 250 for Schedule, 100 for Timezone, 250 for Trick, and 350 for MultiOp. This resulted in a dataset with 1850 questions in total.

3.3 QUALITY CHECK

For both tasks, we did multiple rounds of quality checks where we verified 1) whether the generated labels are correct, and 2) whether the question is clear and the provided instructions are sufficient to know in what format the output should be produced. This procedure was done until no more issues could be found in the dataset.

Table 3: Examples for each question type in the `ToT-Arithmetic` dataset.

Category	Example
AddSubtract	Your driver’s license expires on 18 May, 2017. You receive a renewal notice saying it can be renewed 117 days in advance. What’s the earliest date you can renew your license?
Compare	E42 was discovered in 14 April, 52 BC and E11 was discovered in 05 October, 530 BC. Which was discovered earlier?
Duration	Stella and William were born on 1999-Dec-16 and 2000-Oct-03 respectively. When William was 400 days old, how old was Stella in days?
Schedule	Lucas is available from 11 to noon and also from 3:30 to 5. Asher is available from 11 to 12:30 and also from 4 to 5. They want to have a 30 minute meeting. The meeting has to start on the hour or half hour. How many possibilities are there for the meeting time?
Timezone	Flight departs location A at 11:08 (24hr) UTC(+0000). It reaches location B at 07:23:20 PM IST(+0530). What is the total time duration taken to fly?
Trick	If the date for the day before tomorrow in yyyy-mm-dd format is 2016-01-20, what is the date 27 days from now in the same format?
MultiOp	Alex solves 2 puzzles in 4 hours, 50 minutes, and 22 seconds. What is the time taken by them to solve 6 puzzles, at the same pace.

4 EXPERIMENTS AND RESULTS

In this study, we evaluate the performance of five frontier large language models (LLMs) on our benchmark. The models evaluated include Claude-3-Sonnet (Anthropic, 2024), Mistral Large (2407) (Team, 2024), GPT-4 (Achiam et al., 2023), Gemini 1.5 Pro (Reid et al., 2024), and GPT-4o OpenAI (2024). For GPT-4, we employed GPT-4 Turbo for the `ToT-Semantic` task, as it supports a larger context size, and standard GPT-4 for the `ToT-Arithmetic` task due to its superior performance. The same variant of GPT-4o was used for both tasks.

In our experiments, we aim to answer the following questions:

Table 4: LLM accuracy on temporal reasoning tasks by graph structure.

Graph	Claude-3-Sonnet	Mistral Large	GPT-4	Gemini 1.5 Pro	GPT-4o	Average
BA	48.50	63.00	63.25	62.75	72.00	61.90
Complete	34.00	32.75	40.25	52.50	51.00	42.10
ER	42.25	42.25	48.75	60.50	62.25	51.20
SBM	42.00	48.50	50.75	57.75	61.75	52.15
SFN	58.75	77.75	75.25	75.75	86.00	74.70
Star	59.50	77.50	80.25	74.25	81.75	74.65
AWE	68.75	88.50	92.00	87.50	94.00	86.15
Average Rank	4.75	3.50	2.75	1.43	1.12	

- **RQ1:** What is the effect of the temporal structure on the LLM performance?
- **RQ2:** What kind of temporal questions are easier/harder for LLMs to answer?
- **RQ3:** How important is the order of the facts in the model prompt and what is the best way of ordering the facts?
- **RQ4:** How well do the frontier models perform on two aspects of temporal reasoning: semantics and arithmetic?

4.1 INVESTIGATING THE IMPACT OF TEMPORAL STRUCTURE ON LLM TEMPORAL REASONING

In different applications where temporal reasoning arises, the structure of the facts can be different. Some tasks may provide all the information about an entity (corresponding to a star graph) and then ask questions about it, whereas in some applications such as social networks the structure of the facts may follow a power-law distribution. It is natural to question whether the inherent temporal structure of a problem might influence an LLM’s ability to reason over its data. Drawing inspiration from recent work analyzing graph neural networks (Palowitch et al., 2022; Tsitsulin et al., 2022; Yasir et al., 2023; Fatemi et al., 2024), this section aims to quantify how different temporal dependencies affect an LLM’s temporal reasoning capabilities using graph generators to create many different kinds of temporal structure.

The graph structure of the temporal relationships significantly affects LLM performance, as demonstrated in Table 4. Notably, GPT-4 accuracy more than doubled between complete graphs (40.25%) and AWE graphs (92.00%). Also, Mistral Large accuracy varied drastically across graph types, from 32.75% for complete graphs to 88.50% for AWE graphs. This highlights a critical gap in temporal reasoning research, which has largely overlooked the diverse graph structures and reasoning tasks found in real-world applications, instead focusing primarily on specific knowledge graphs (like YAGO11k). This may explain the superior performance of LLMs on AWE graphs in our experiments, with GPT-4o nearly solving the task with 94.00% accuracy.

4.1.1 INFLUENCE OF GRAPH SIZE ON LLM PERFORMANCE

A key question is how different models behave as a function of the size of a graph, measured in terms of the number of edges (facts) and nodes (entities). As illustrated in Figure 4, increasing either the number of edges or nodes in the `ToT-Semantic` dataset mostly leads to a decrease in LLM performance. We observe, however, that different models are affected differently. For example, for the smaller graphs with < 250 edges, GPT-4o outperforms the other models, whereas when the size increases to > 1000 edges, Gemini 1.5 Pro outperforms the other models. Moreover, we observe that while the performance of GPT-4o and Gemini 1.5 Pro does not degrade much after a certain point of increasing the number of edges (specifically, for the last three buckets), other models’ performances keep decreasing (with the exception of GPT-4 at the last bucket).

The above results raise the question of whether the graph structure’s impact observed in Section 4.1 is merely a consequence of varying graph sizes. To address this, we present the average number of

Table 5: Average number of nodes and edges by graph structure.

Graph	#nodes	#edges
BA	17.41	144.07
Complete	17.25	619.86
ER	16.18	316.4
SBM	17.51	368.15
SFN	17.52	53.46
Star	16.16	34.12
AWE	18.99	25.41
Average	17.29	223.07

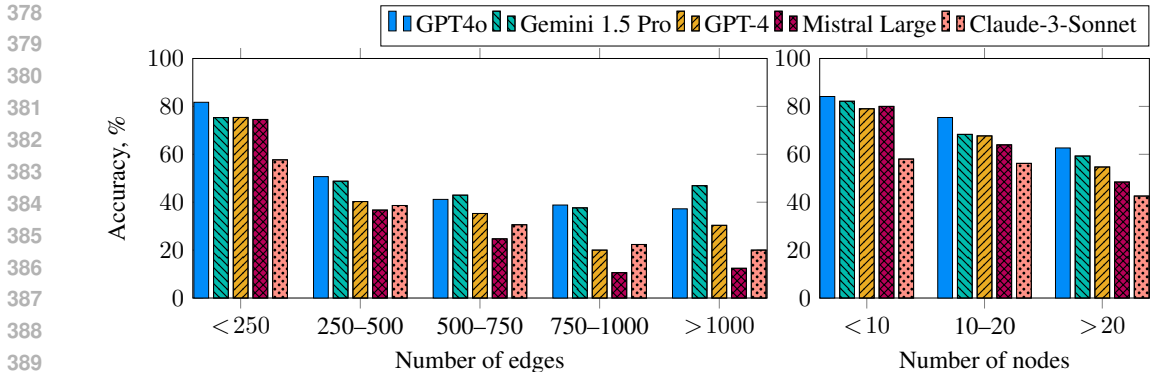


Figure 4: Accuracy of models for different number of edges and nodes.

nodes and edges for each graph structure in Table 5. While the average number of nodes does not appear to consistently influence LLM performance across structures, the number of edges does show some correlation. However, there are exceptions. For instance, SBM graphs have far more edges on average than ER graphs, yet the average performance of models across ER graphs is lower than SBM graphs. Also, SFN graphs have far more edges on average than Star graphs, yet GPT-4o performs better on SFN graphs than Star graphs. This indicates that both the number of edges and the specific structure of the graph play a significant role in determining LLM performance. As for number of nodes, AWE graphs have more nodes on average compared to the other graph structures, yet the average performance of models across AWE is the highest across all (see Table 4).

4.2 EFFECTS OF TEMPORAL QUESTION TYPE ON LLM TEMPORAL REASONING

Table 6: LLM accuracy on temporal reasoning by question category.

Temporal Question Type	Claude-3-Sonnet	Mistral Large	GPT-4	Gemini 1.5 Pro	GPT-4o	Average
EventAtTimeT	47.14	64.86	65.43	72.29	71.43	64.23
EventAtWhatTime	90.29	90.00	89.43	93.14	96.86	91.94
NumEventsInTimeInterval	29.71	57.14	61.43	59.14	63.71	54.23
BeforeAfter	53.14	56.57	55.43	52.86	64.00	56.40
EventAtTimeOfAnotherEvent	50.00	57.43	67.14	71.43	75.71	64.34
FirstLast	68.57	57.71	67.71	68.57	73.71	67.25
RelationDuration	41.43	76.57	80.00	84.57	88.86	74.29
Timeline	24.00	31.14	28.29	36.29	38.57	31.66
Average Rank	4.31	3.75	3.37	2.44	1.12	

In this experiment, we systematically investigate the impact of different temporal tasks on the reasoning ability of LLMs. We quantify this impact by evaluating model performance across a variety of tasks, as summarized in Table 6.

Task type and reasoning requirements. A key question in our investigation is whether the type of temporal task and the associated reasoning requirements influence LLM performance. The `ToT-Semantic` dataset includes questions of varying difficulty levels, which can be categorized based on the reasoning type: **Single-fact solutions:** Questions `EventAtTimeT` and `EventAtWhatTime` require retrieving one single fact and answering the question directly based on the retrieved fact. **Multi-fact solutions:** The remaining questions require retrieving multiple facts and performing operations (*e.g.*, counting, sorting) to extract relevant information and formulate an answer.

LLMs consistently demonstrate superior performance on tasks requiring the retrieval of a single fact compared to those necessitating the integration of multiple facts. This performance gap can be attributed to the increased cognitive demands associated with multi-fact tasks. While single-fact questions primarily rely on the identification and extraction of relevant information, multi-fact questions demand a deeper comprehension and synthesis of retrieved information.

Performance variations within question types. Even within zero-order reasoning tasks, LLMs demonstrate varying levels of proficiency. For example, EventAtTimeT and EventAtWhatTime are structurally similar, yet LLMs tend to excel at the latter. We hypothesize that this performance difference may be attributed to the fact that EventAtTimeT requires a simple time arithmetic operation to recognize that a timestamp T falls within a time interval $[T1, T2]$, whereas EventAtWhatTime does not require any time arithmetic operation.

Analysis on Timeline questions. Timeline questions are the most difficult category of questions for the models according to Table 6. An analysis of these questions reveals that they pose the greatest challenge across all tasks. To answer these questions, typically structured as “Sort the entities that were the R17 of E69 chronologically?”, the model needs to extract multiple entities (in the ToT-Semantic dataset, every timeline question has more than one entity in its label), and then do temporal arithmetic to sort them temporally. To further analyze the models on these questions, we calculated the average precision and recall for each model in Table 7, where precision shows what percentage of the entities extracted by the model are correct entities (i.e. must be included in the timeline) and recall shows what percentage of the correct entities have been extracted by the model. We report the results once averaged over all graph structures and once only for complete graphs (the most challenging graph structure). Gemini 1.5 Pro demonstrates superior precision and recall, aligning with its relatively high accuracy observed in Table 6. The only model outperforming Gemini 1.5 Pro on timeline questions is GPT-4o. The fact that the precision and recall of GPT-4o is lower than that of Gemini but its overall performance on timeline questions is higher shows that Gemini is better at retrieving the correct entities but worse at arithmetic operation (as also confirmed later in Section 4.4). Moreover, GPT-4, despite having higher accuracy than Claude-3-Sonnet on timeline questions, exhibits the lowest precision and recall. This suggests that GPT-4 frequently outputs fewer entities than are present in the true answers (50% of the times), leading to missed correct entities (lower recall) and a higher proportion of false positives among its predictions (lower precision).

Since complete graphs pose the greatest difficulty among all graph structures (Table 4), we provide a separate analysis of average precision and recall for these graphs in the final two columns of Table 7. Notably, all models except Gemini 1.5 Pro experienced declines in both precision and recall on complete graphs, whereas Gemini was primarily impacted in terms of recall.

Table 7: Precision and recall on timeline questions.

Graph structure	All		Complete	
	Precision	Recall	Precision	Recall
Claude-3-Sonnet	0.73	0.75	0.56	0.54
Mistral Large	0.62	0.65	0.30	0.33
GPT-4	0.60	0.56	0.36	0.23
Gemini 1.5 Pro	0.81	0.83	0.82	0.65
GPT-4o	0.78	0.74	0.69	0.51

Table 8: LLM accuracy on temporal reasoning tasks as a function of the order of the facts.

Order of facts	Claude-3-Sonnet	Mistral Large	GPT-4	Gemini 1.5 Pro	GPT-4o	Average
Shuffle	45.71	55.71	60.71	63.04	68.93	58.82
RelationAndStartTime	54.29	63.93	65.36	68.57	72.14	64.86
StartTimeAndRelation	47.68	59.11	60.54	64.64	65.89	59.57
StartTimeAndTarget	49.11	60.89	61.61	65.18	70.00	61.36
TargetAndStartTime	73.57	67.50	62.60	75.00	81.07	71.95

4.3 IMPACT OF TEMPORAL FACT ORDER ON LLM PERFORMANCE

A noteworthy question arises regarding the potential influence of fact order on LLM performance in temporal reasoning tasks. To investigate this, we conducted experiments on ToT-Semantic dataset. We sorted the facts using different methods: **Shuffle**: randomizing the order of facts; **RelationAndStartTime**: prioritizing facts based on their relation name, then by start time; **StartTimeAndRelation**: prioritizing facts based on start time, then by relation name; **StartTimeAndTarget**: prioritizing facts based on start time, then by the target entity; **TargetAndStartTime**: Prioritizing facts based on the target entity, then by start time.

Ideally, LLMs should exhibit robustness to the order in which facts are presented, given the independent nature of each fact. However, as shown in Table 8, our observations reveal a significant impact of fact order on LLM performance. Notably, performance is consistently lowest when facts are presented in a shuffled order and consistently highest when facts are sorted based on the target entity and start time (TargetAndStartTime). We also observe that some sorting strategies such as

486 StartTimeAndRelation are only slightly better than the shuffled order, thus revealing that not any
 487 kind of ordering is ideal for LLMs. This finding offers valuable practical insights into how facts
 488 should be structured when temporal reasoning is a key component of the LLM task. By organizing
 489 facts in a manner that aligns with the temporal flow of the narrative or task, we can potentially
 490 enhance LLM performance and ensure more accurate and reliable results. While previous work
 491 has shown that ordering premises in the correct order of chain-of-thought solution improves LLM’s
 492 logical reasoning (Chen et al., 2024; Saparov & He, 2022), our results extend that to general-purpose
 493 temporal orderings (independent of the chain-of-thought).

494 Table 9: LLM accuracy on the ToT-Arithmetic dataset by question type.
 495

496 Category	Claude-3-Sonnet	Mistral Large	GPT-4	Gemini 1.5 Pro	GPT-4o	Average
497 AddSubtract	58.57	61.14	76.28	71.14	76.29	68.68
498 Compare	39.14	62.29	63.14	55.43	66.57	57.30
499 Duration	15.00	17.50	16.00	13.50	15.00	15.40
500 Schedule	29.60	44.40	43.60	40.00	53.20	42.16
501 Timezone	74.00	87.00	88.00	90.00	92.00	86.20
502 Trick	40.40	44.80	45.60	41.20	53.20	45.04
503 MultiOp	26.57	54.86	46.86	62.57	46.86	47.54
504 Average Rank	4.71	2.71	2.43	3.28	1.57	

505 4.4 TEMPORAL SEMANTICS VS TEMPORAL ARITHMETIC

506 This study examined the performance of temporal arithmetic capabilities in LLMs using the
 507 ToT-Arithmetic dataset. Results, as shown in Table 9, indicate that the models consistently
 508 excelled in Timezone questions, while struggling the most with Duration questions. This superior
 509 performance in Timezone questions could be attributed to the abundance of information about various
 510 timezones available online, compared to other question types. Scheduling and Trick questions also
 511 proved challenging for LLMs, likely due to their creative nature and requirement for deeper reasoning.
 512 In contrast, AddSubtract results were relatively strong, potentially reflecting LLMs’ optimization for
 513 mathematical reasoning and their ability to apply that knowledge to temporal reasoning tasks.

514 **Analysis on Duration questions.** Analysis of Duration questions in the ToT-Arithmetic dataset
 515 revealed them to be the most challenging for the evaluated models. Notably, the most common
 516 error among incorrect answers was a deviation of precisely one day from the ground truth label.
 517 Specifically, when GPT-4 or Gemini 1.5 Pro erred on Duration questions, approximately 21% and
 518 25% of its responses were within one day of the ground truth, respectively. This suggests that LLMs
 519 can approximate the correct calculation but often stumble in the final steps, highlighting a gap in their
 520 ability to execute complex arithmetic with precision.

521 **Common failure: direction.** One frequent error in ToT-Arithmetic occurs when determining
 522 the number of months between two dates. For example, from February 11th, 2002, to October 11th,
 523 2002, the correct duration is eight months, but the model sometimes incorrectly calculates it as four
 524 months. This issue is particularly noticeable in questions that involve going back in time, such as:
 525 “Sam’s birthdate is October 11th, 1996. Today is February 25th, 2002. Calculate Sam’s age in days.”

526 **Common failure: leap year calculation.** Another frequent error in ToT-Arithmetic arises
 527 when determining the number of days between two dates that span multiple years. Incorrectly
 528 accounting for leap years, which have an extra day (February 29th), often leads to inaccurate results.
 529

530 5 CONCLUSION

531 In conclusion, we introduce Test of Time (ToT), a novel benchmark designed to assess LLMs’
 532 temporal reasoning abilities in a more comprehensive and controlled manner than existing work.
 533 Our two-pronged approach, encompassing both semantic and arithmetic tasks, enables a nuanced
 534 evaluation of temporal reasoning. Through extensive experiments with ToT, we have gained valuable
 535 insights into the strengths and weaknesses of current LLMs in these critical aspects of temporal
 536 reasoning. By open-sourcing our datasets and evaluation framework, we hope to stimulate further
 537 research and development in this field, ultimately contributing to the advancement of LLM capabilities
 538 in complex reasoning tasks.
 539

REPRODUCIBILITY STATEMENT

To ensure the reproducibility of our work, we provide the following resources and information:

Benchmark creation: A detailed description of the construction methodology for our temporal reasoning benchmark is available in Section 3. This includes the process of creation of both ToT-Semantic and ToT-Arithmetic.

LLM access: The LLMs evaluated in this study are publicly accessible via API calls. We specify the names of the LLMs and the corresponding versions used for our experiments in Section 4.

Evaluation procedure: Appendix D outlines the evaluation procedure used for our study along with some examples to better clarify the procedure.

Furthermore, we will make the code and the generated benchmark dataset publicly available upon publication to facilitate the reproduction of our results and encourage further research in this direction.

REFERENCES

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Janice Ahn, Rishu Verma, Renze Lou, Di Liu, Rui Zhang, and Wenpeng Yin. Large language models for mathematical reasoning: Progresses and challenges. *arXiv preprint arXiv:2402.00157*, 2024.
- Réka Albert and Albert-László Barabási. Statistical mechanics of complex networks. *Reviews of modern physics*, 74(1):47, 2002.
- Anthropic. Introducing the next generation of claude. <https://www.anthropic.com/news/claude-3-family>, 2024. Available at: <https://www.anthropic.com/news/claude-3-family>.
- Albert-László Barabási and Réka Albert. Emergence of scaling in random networks. *science*, 286(5439):509–512, 1999.
- Himanshu Beniwal, Mayank Singh, et al. Remember this event that year? assessing temporal information and reasoning in large language models. *arXiv preprint arXiv:2402.11997*, 2024.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *NeurIPS*, 2020.
- Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. Sparks of artificial general intelligence: Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*, 2023.
- Ines Chami, Sami Abu-El-Haija, Bryan Perozzi, Christopher Ré, and Kevin Murphy. Machine learning on graphs: A model and comprehensive taxonomy. *Journal of Machine Learning Research*, 23(89):1–64, 2022.
- Wenhu Chen, Xinyi Wang, and William Yang Wang. A dataset for answering time-sensitive questions. *arXiv preprint arXiv:2108.06314*, 2021.
- Xinyun Chen, Ryan A Chi, Xuezhi Wang, and Denny Zhou. Premise order matters in reasoning with large language models. *arXiv preprint arXiv:2402.08939*, 2024.
- Ziyang Chen, Xiang Zhao, Jinzhi Liao, Xinyi Li, and Evangelos Kanoulas. Temporal knowledge graph question answering via subgraph reasoning. *Knowledge-Based Systems*, 251:109134, 2022.
- Zheng Chu, Jingchang Chen, Qianglong Chen, Weijiang Yu, Haotian Wang, Ming Liu, and Bing Qin. Timebench: A comprehensive evaluation of temporal reasoning abilities in large language models. *arXiv preprint arXiv:2311.17667*, 2023.

- 594 Antonia Creswell, Murray Shanahan, and Irina Higgins. Selection-inference: Exploiting large
595 language models for interpretable logical reasoning. *arXiv preprint arXiv:2205.09712*, 2022.
596
- 597 Shib Sankar Dasgupta, Swayambhu Nath Ray, and Partha Talukdar. Hyte: Hyperplane-based
598 temporally aware knowledge graph embedding. In *Proceedings of the 2018 conference on empirical*
599 *methods in natural language processing*, pp. 2001–2011, 2018.
- 600 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of
601 deep bidirectional transformers for language understanding. In Jill Burstein, Christy Doran, and
602 Tamar Solorio (eds.), *Proceedings of the 2019 Conference of the North American Chapter of the*
603 *Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and*
604 *Short Papers)*, pp. 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational
605 Linguistics. doi: 10.18653/v1/N19-1423. URL <https://aclanthology.org/N19-1423>.
606
- 607 Bhuwan Dhingra, Jeremy R Cole, Julian Martin Eisenschlos, Daniel Gillick, Jacob Eisenstein, and
608 William W Cohen. Time-aware language models as temporal knowledge bases. *Transactions of*
609 *the Association for Computational Linguistics*, 10:257–273, 2022.
- 610 Paul Erdős and Alfred Rényi. On random graphs. *Publicationes Mathematicae Debrecen*, 6:290–297,
611 1959.
- 612 Bahare Fatemi, Layla El Asri, and Seyed Mehran Kazemi. Slaps: Self-supervision improves structure
613 learning for graph neural networks. *Advances in Neural Information Processing Systems*, 34:
614 22667–22681, 2021.
- 615 Bahare Fatemi, Sami Abu-El-Haija, Anton Tsitsulin, Mehran Kazemi, Dustin Zelle, Neslihan Bulut,
616 Jonathan Halcrow, and Bryan Perozzi. Ugs! : A unified framework for benchmarking graph
617 structure learning. *arXiv preprint arXiv:2308.10737*, 2023.
- 618 Bahare Fatemi, Jonathan Halcrow, and Bryan Perozzi. Talk like a graph: Encoding graphs for large
619 language models. In *ICLR*, 2024.
- 620 Vivek Gupta, Pranshu Kandoi, Mahek Bhavesh Vora, Shuo Zhang, Yujie He, Ridho Reinanda, and
621 Vivek Srikumar. Temptabqa: Temporal question answering for semi-structured tables. *arXiv*
622 *preprint arXiv:2311.08002*, 2023.
- 623 Wes Gurnee and Max Tegmark. Language models represent space and time. *arXiv preprint*
624 *arXiv:2310.02207*, 2023.
- 625 Jonathan Halcrow, Alexandru Mosoi, Sam Ruth, and Bryan Perozzi. Grale: Designing networks for
626 graph learning. In *Proceedings of the 26th ACM SIGKDD international conference on knowledge*
627 *discovery & data mining*, pp. 2523–2532, 2020.
- 628 Carl Hewitt. Planner: A language for proving theorems in robots. In *Proceedings of the 1st*
629 *International Joint Conference on Artificial Intelligence, IJCAI’69*, pp. 295–301, San Francisco,
630 CA, USA, 1969. Morgan Kaufmann Publishers Inc.
- 631 Paul W Holland, Kathryn Blackmond Laskey, and Samuel Leinhardt. Stochastic blockmodels: First
632 steps. *Social networks*, 5(2):109–137, 1983.
- 633 Xuming Hu, Junzhe Chen, Xiaochuan Li, Yufei Guo, Lijie Wen, Philip S Yu, and Zhijiang Guo. Do
634 large language models know about facts? *arXiv preprint arXiv:2310.05177*, 2023.
- 635 Jie Huang and Kevin Chen-Chuan Chang. Towards reasoning in large language models: A survey,
636 2023.
- 637 Zhen Jia, Abdalghani Abujabal, Rishiraj Saha Roy, Jannik Strötgen, and Gerhard Weikum. Tem-
638 pquestions: A benchmark for temporal question answering. In *Companion Proceedings of the The*
639 *Web Conference 2018*, pp. 1057–1062, 2018.
- 640 Zhen Jia, Soumajit Pramanik, Rishiraj Saha Roy, and Gerhard Weikum. Complex temporal question
641 answering on knowledge graphs. In *Proceedings of the 30th ACM international conference on*
642 *information & knowledge management*, pp. 792–802, 2021.

- 648 Zhen Jia, Philipp Christmann, and Gerhard Weikum. Tiq: A benchmark for temporal question an-
649 swering with implicit time constraints. In *Companion Proceedings of the ACM on Web Conference*
650 *2024*, pp. 1394–1399, 2024.
- 651
- 652 Mehran Kazemi, Hamidreza Alvari, Ankit Anand, Jialin Wu, Xi Chen, and Radu Soricut. Geomverse:
653 A systematic evaluation of large models for geometric reasoning. *arXiv preprint arXiv:2312.12241*,
654 2023a.
- 655 Mehran Kazemi, Najoung Kim, Deepti Bhatia, Xin Xu, and Deepak Ramachandran. LAMBADA:
656 Backward chaining for automated reasoning in natural language. In Anna Rogers, Jordan Boyd-
657 Graber, and Naoaki Okazaki (eds.), *Proceedings of the 61st Annual Meeting of the Association*
658 *for Computational Linguistics (Volume 1: Long Papers)*, pp. 6547–6568, Toronto, Canada, July
659 2023b. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.361. URL
660 <https://aclanthology.org/2023.acl-long.361>.
- 661
- 662 Mehran Kazemi, Quan Yuan, Deepti Bhatia, Najoung Kim, Xin Xu, Vaiva Imbrasaitė, and Deepak
663 Ramachandran. Boardgameqa: A dataset for natural language reasoning with contradictory
664 information. *Advances in Neural Information Processing Systems*, 36, 2023c.
- 665 Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay Ra-
666 masesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, et al. Solving quantitative
667 reasoning problems with language models. *Advances in Neural Information Processing Systems*,
668 35:3843–3857, 2022.
- 669 Adam Liska, Tomas Kocisky, Elena Gribovskaya, Tayfun Terzi, Eren Sezener, Devang Agrawal,
670 D’Autume Cyprien De Masson, Tim Scholtes, Manzil Zaheer, Susannah Young, et al. Streamingqa:
671 A benchmark for adaptation to new knowledge over time in question answering models. In
672 *International Conference on Machine Learning*, pp. 13604–13622. PMLR, 2022.
- 673
- 674 Jason Xinyu Liu, Ziyi Yang, Ifrah Idrees, Sam Liang, Benjamin Schornstein, Stefanie Tellex, and
675 Ankit Shah. Grounding complex natural language commands for temporal tasks in unseen
676 environments. In *Conference on Robot Learning*, pp. 1084–1110. PMLR, 2023.
- 677 John McCarthy. Programs with common sense. In *Proceedings of the Teddington Conference on the*
678 *Mechanization of Thought Processes*, pp. 75–91, London, 1959. Her Majesty’s Stationary Office.
679 URL <http://www-formal.stanford.edu/jmc/mcc59.html>.
- 680
- 681 Sumit Neelam, Udit Sharma, Hima Karanam, Shajith Iqbal, Pavan Kapanipathi, Ibrahim Abdelaziz,
682 Nandana Mihindukulasooriya, Young-Suk Lee, Santosh Srivastava, Cezar Pendus, et al. Sygma:
683 System for generalizable modular question answering overknowledge bases. *arXiv preprint*
684 *arXiv:2109.13430*, 2021.
- 685 Qiang Ning, Zhili Feng, Hao Wu, and Dan Roth. Joint reasoning for temporal and causal relations.
686 *arXiv preprint arXiv:1906.04941*, 2019.
- 687
- 688 Kai Nylund, Suchin Gururangan, and Noah A Smith. Time is encoded in the weights of finetuned
689 language models. *arXiv preprint arXiv:2312.13401*, 2023.
- 690
- 691 OpenAI. Hello gpt-4o. <https://openai.com/index/hello-gpt-4o/>, 2024. Accessed:
692 2024-10-01.
- 693 John Palowitch, Anton Tsitsulin, Brandon Mayer, and Bryan Perozzi. Graphworld: Fake graphs
694 bring real insights for gnns. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge*
695 *Discovery and Data Mining*, pp. 3691–3701, 2022.
- 696
- 697 Bryan Perozzi, Rami Al-Rfou, and Steven Skiena. Deepwalk: Online learning of social representa-
698 tions. In *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery*
699 *and data mining*, pp. 701–710, 2014.
- 700
- 701 Bryan Perozzi, Bahare Fatemi, Dustin Zelle, Anton Tsitsulin, Mehran Kazemi, Rami Al-Rfou, and
Jonathan Halcrow. Let your graph do the talking: Encoding structured data for llms. *arXiv preprint*
arXiv:2402.05862, 2024.

- 702 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi
703 Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text
704 transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551, 2020.
- 705
- 706 Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy Lillicrap, Jean-baptiste
707 Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, et al. Gemini
708 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint*
709 *arXiv:2403.05530*, 2024.
- 710 Benedek Rozemberczki, Peter Englert, Amol Kapoor, Martin Blais, and Bryan Perozzi. Pathfinder
711 discovery networks for neural message passing. In *Proceedings of the Web Conference 2021*, pp.
712 2547–2558, 2021.
- 713 Clayton Sanford, Bahare Fatemi, Ethan Hall, Anton Tsitsulin, Mehran Kazemi, Jonathan Halcrow,
714 Bryan Perozzi, and Vahab Mirrokni. Understanding transformer reasoning capabilities via graph
715 algorithms. *arXiv preprint arXiv:2405.18512*, 2024.
- 716
- 717 Abulhair Saparov and He He. Language models are greedy reasoners: A systematic formal analysis
718 of chain-of-thought. *arXiv preprint arXiv:2210.01240*, 2022.
- 719
- 720 Abulhair Saparov, Richard Yuanzhe Pang, Vishakh Padmakumar, Nitish Joshi, Mehran Kazemi,
721 Najoung Kim, and He He. Testing the general deductive reasoning capacity of large language
722 models using ood examples. *Advances in Neural Information Processing Systems*, 36, 2023.
- 723 Saurabh Srivastava, Anto PV, Shashank Menon, Ajay Sukumar, Alan Philipose, Stevin Prince, Sooraj
724 Thomas, et al. Functional benchmarks for robust evaluation of reasoning performance, and the
725 reasoning gap. *arXiv preprint arXiv:2402.19450*, 2024.
- 726
- 727 Zhaochen Su, Jun Zhang, Tong Zhu, Xiaoye Qu, Juntao Li, Min Zhang, and Yu Cheng. Timo:
728 Towards better temporal reasoning for language models. *arXiv preprint arXiv:2406.14192*, 2024.
- 729 Oyvind Tafjord, Bhavana Dalvi, and Peter Clark. ProofWriter: Generating implications, proofs,
730 and abductive statements over natural language. In *Findings of the Association for Com-
731 putational Linguistics: ACL-IJCNLP 2021*, pp. 3621–3634, Online, August 2021. Associa-
732 tion for Computational Linguistics. doi: 10.18653/v1/2021.findings-acl.317. URL <https://aclanthology.org/2021.findings-acl.317>.
- 733
- 734 Qingyu Tan, Hwee Tou Ng, and Lidong Bing. Towards benchmarking and improving the temporal
735 reasoning capability of large language models. *arXiv preprint arXiv:2306.08952*, 2023.
- 736
- 737 Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu
738 Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly capable
739 multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- 740 Mistral AI Team. Large enough. <https://mistral.ai/news/mistral-large-2407/>,
741 2024. Accessed: 2024-10-01.
- 742
- 743 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay
744 Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation
745 and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- 746 Anton Tsitsulin, Benedek Rozemberczki, John Palowitch, and Bryan Perozzi. Synthetic graph
747 generation to benchmark graph learning. *arXiv preprint arXiv:2204.01376*, 2022.
- 748
- 749 Siddharth Vashishtha, Adam Poliak, Yash Kumar Lal, Benjamin Van Durme, and Aaron Steven
750 White. Temporal reasoning in natural language inference. In *Findings of the Association for
751 Computational Linguistics: EMNLP 2020*, pp. 4070–4078, 2020.
- 752 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
753 Kaiser, and Illia Polosukhin. Attention is all you need. *NeurIPS*, 30, 2017.
- 754
- 755 Denny Vrandečić and Markus Krötzsch. Wikidata: a free collaborative knowledgebase. *Communica-
tions of the ACM*, 57(10):78–85, 2014.

- 756 Yiqi Wang, Wentao Chen, Xiaotian Han, Xudong Lin, Haiteng Zhao, Yongfei Liu, Bohan Zhai, Jianbo
757 Yuan, Quanzeng You, and Hongxia Yang. Exploring the reasoning abilities of multimodal large
758 language models (mlms): A comprehensive survey on emerging trends in multimodal reasoning.
759 *arXiv preprint arXiv:2401.06805*, 2024.
- 760 Yuqing Wang and Yun Zhao. Tram: Benchmarking temporal reasoning for large language models.
761 *arXiv preprint arXiv:2310.00835*, 2023.
- 762 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny
763 Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in
764 neural information processing systems*, 35:24824–24837, 2022.
- 765 Siheng Xiong, Ali Payani, Ramana Kompella, and Faramarz Fekri. Large language models can learn
766 temporal reasoning. *arXiv preprint arXiv:2401.06853*, 2024.
- 767 Sen Yang, Xin Li, Lidong Bing, and Wai Lam. Once upon a time in graph: Relative-time pretraining
768 for complex temporal reasoning. In *Proceedings of the 2023 Conference on Empirical Methods in
769 Natural Language Processing*, pp. 11879–11895, 2023.
- 770 Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan.
771 Tree of thoughts: Deliberate problem solving with large language models. *Advances in Neural
772 Information Processing Systems*, 36, 2023.
- 773 Mustafa Yasir, John Palowitch, Anton Tsitsulin, Long Tran-Thanh, and Bryan Perozzi. Examining
774 the effects of degree distribution and homophily in graph learning models, 2023.
- 775 Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine
776 really finish your sentence?, 2019.
- 777 Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min,
778 Beichen Zhang, Junjie Zhang, Zican Dong, et al. A survey of large language models. *arXiv
779 preprint arXiv:2303.18223*, 2023.
- 780 Ben Zhou, Kyle Richardson, Qiang Ning, Tushar Khot, Ashish Sabharwal, and Dan Roth. Temporal
781 reasoning on implicit events from distant supervision. *arXiv preprint arXiv:2010.12753*, 2020.

782 A DESCRIPTION OF GRAPH GENERATORS.

783 Here we detail each graph generator used to create the examples in TOT . We note that every collection
784 of temporal facts, where each fact is a relationship between two entities, can be expressed as a
785 temporal graph with nodes as entities. TOT specifically targets LLM reasoning ability over such
786 collections. We do not claim that graph generators are the only way to construct such a benchmark.
787 However, because all temporal fact collections contain an underlying graph, we propose a generation
788 framework based on graph models to produce benchmark examples. We argue that a framework that
789 exposes generation of the static graph backbone is more controllable and allows for a benchmark that
790 is more comprehensive w.r.t. the variety and complexity of temporal relationships between generated
791 entities.

800 First, we cover the six *random* graph generators used to create the synthetic examples. All random
801 graph generators are probabilistic models which take hyperparameters that control the expected
802 macro-properties of each graph (Palowitch et al., 2022):

- 803 • Erdős-Rényi (ER) (Erdős & Rényi, 1959): This model takes an edge probability parameter
804 p and generates each edge with probability p , i.i.d. over all possible edges.
- 805 • Scale-Free Networks (SFN) (Barabási & Albert, 1999): a graph is grown by a sequence
806 of steps, each step either (1) adding a new node and connecting it to an existing node, or
807 (2) adding an edge between two existing nodes. Input parameters control the probability
808 of these events. This process generates a “scale-free” power law of node degrees, in sharp
809 contrast to the ER model.

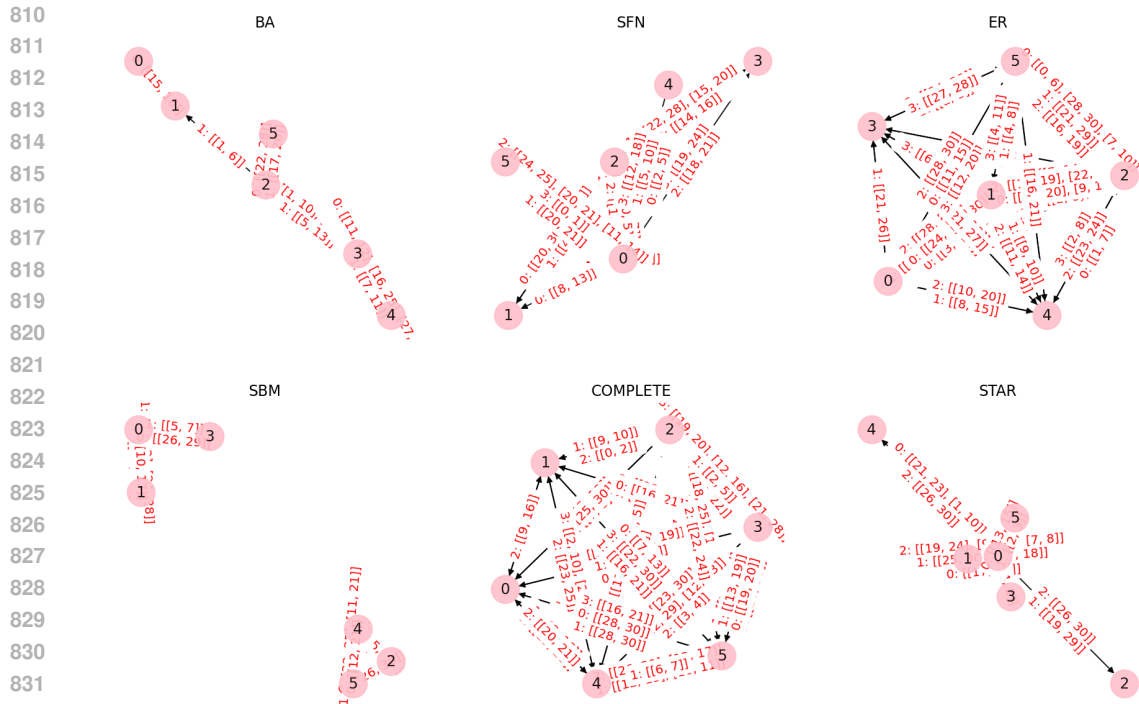


Figure 5: A visualization of a representative graph from each graph generator: Erdős-Rényi (ER), Scale-Free Networks (SFN), Barabási-Albert (BA), Stochastic Block Model (SBM), star-graph, and complete-graph.

- Barabási-Albert (BA) model (Albert & Barabási, 2002): a graph is grown by a sequence of steps, each step adding a new node to the graph, and connecting the node to m existing nodes with probability proportional to their current degree. Similar to SFN, this process also generates a “scale-free” graph with a particular distribution known as the Barabási-Albert model.
- Stochastic Block Model (SBM) (Holland et al., 1983): This graph model can be thought of as clustered ER. It divides n nodes into k clusters, and then connects two nodes with probability p if they are in the same cluster, else with probability q if they are in different clusters. k , p , and q are all hyperparameters.
- A star-graph generator creates a “star” graph on n nodes: node 0 is the center of the star, and all other nodes connect to it (and only it).
- A complete-graph generate creates a “complete” graph on n nodes, in which all nodes connect to each other node.

An example from each of the above graph generators is shown in Figure 5. In the figure, edges are annotated with temporal relationships in the format `relation_id: [[interval_1, ..., interval_k]]`. Note that each edge can have multiple relationships, and each relationship can have multiple intervals. The visualization shows the diversity of temporal knowledge graphs that our framework is able to generate. We note that while our study was limited to parametric graph generators in this work, the field of graph machine learning (Chami et al., 2022) offers many options for both modeling (Perozzi et al., 2014) and learning (Halcrow et al., 2020; Fatemi et al., 2021; Rozemberczki et al., 2021; Fatemi et al., 2023) link structure.

Second, we describe our Anonymized Wikidata Extract (AWE) strategy for creating anonymized questions from real-world data. We first identify the 78 most common relations in WikiData that specify time-bound entity relationships. Each relation encodes a temporal edge between two entities. To match the schema of our synthetic graphs, we convert each time specification on each edge to an interval. Then, for each entity in the graph, we extract the *ego-graph* of the entity by (1) collecting the entity and all its neighbors and (2) collecting all edges (along with temporal information) between

nodes collected in (1). This process produces a temporal graph with a schema identical to those produced from random graph generators. Before generating questions from the graphs, we anonymize them by (a) mapping each entity name to a unique identifier such as E_{679} ; and then (b) mapping each relation name to a unique identifier such as R_3 . We then generate questions from the graph as described in 3.1.

B DETAILS OF QUESTION GENERATION.

Given a graph with temporal facts, generating logically-consistent questions from our list of diverse question types (see Table 2) is non-trivial. To generate the total question set, we loop through generated graphs, choose a question type uniformly-at-random, and then attempt to fill the question type template with facts from the graph. The exact algorithmic procedure is given below. Note that the `SAMPLEFACTS` routine will vary significantly depending on the question type. For some questions, it is sufficient to generate a single fact and check if the question can be generated. For other question, multiple facts must be sampled (sometimes sequentially, in a BFS fashion) and checked for cohesion with the template. We do all of this in a brute-force manner.

Algorithm 1 Generate all questions from a certain question type template.

```

1: Procedure GENERATEQUESTIONS( $\mathcal{G}, n, \text{template}, m$ )
2:  $Q \leftarrow \phi$ 
3: for  $i \in [n]$  do
4:    $G \leftarrow \text{SAMPLEGRAPH}(\mathcal{G})$ 
5:    $q \leftarrow \text{GENERATEQUESTION}(G, \text{template}, m)$ 
6:   if  $q = \phi$  then
7:     continue
8:   end if
9:    $Q \leftarrow Q \cup \{q\}$ 
10: end for
11: return  $Q$ 

```

Algorithm 2 Generate a single question from a graph with maximum trials m .

```

1: Procedure GENERATEQUESTION( $G, \text{template}, m$ )
2:  $q = \phi$ 
3: for  $j \in [m]$  do
4:    $F \leftarrow \text{SAMPLEFACTS}(G, \text{template})$ 
5:    $q \leftarrow \text{template}(F)$ 
6:   if  $q \neq \phi$  then
7:     break
8:   end if
9: end for
10: return  $q$ 

```

C LARGE-SCALE TOT-SEMANTIC EXPERIMENTS

To facilitate a more comprehensive analysis and enable deeper insights, we expanded our synthetic dataset significantly. This enlarged dataset now encompasses approximately 50,000 examples, a substantial increase from the previous set of around 3,000 examples. We anticipate that this expanded resource will prove valuable for future research endeavors that necessitate a larger and more diverse synthetic dataset. Due to the computational demands of evaluating all LLMs on this large dataset, results are reported solely for Gemini 1.5 Pro.

Impact of Graph Structure on LLM Accuracy. Our initial experiment with this expanded dataset involved replicating the graph structure analysis. As illustrated in Table 10, graph structure continues to exert a significant influence on the final accuracy of the LLM, even within this larger dataset.

Table 10: LLM temporal reasoning by graph structure on the larger set of ToT-Semantic.

Graph Structure	Accuracy (%)
BA	70.96
Complete	51.07
ER	61.85
SBM	60.32
SFN	79.13
Star	73.77
AWE	88.72
Average	69.40

Table 11: Impact of graph structure and question type on a larger set of ToT-Semantic.

Temporal task	BA	Complete	ER	SBM	SFN	Star	AWE	Average Rank
EventAtTimeT	74.46	54.22	65.54	68.07	80.84	76.75	91.93	3.57
EventAtWhatTime	98.19	81.69	90.72	90.48	98.31	98.43	97.95	1.00
BeforeAfter	53.49	34.46	48.07	45.66	68.55	58.80	73.98	7.00
EventAtTimeOfAnotherEvent	76.99	52.89	62.53	65.18	84.82	85.78	90.48	3.79
FirstLast	70.84	49.04	61.69	55.66	87.23	68.80	92.53	4.43
NumEventsInTimeInterval	57.71	40.84	54.22	49.64	64.22	70.84	83.73	6.14
RelationDuration	88.55	80.60	83.49	82.77	87.47	88.80	90.48	2.36
Timeline	47.47	14.82	28.55	25.06	61.57	41.93	88.67	7.71

Impact of graph structure and temporal task on LLM performance. Our second experiment examined the accuracy of the model across various question types and graph generators. The expanded dataset provided sufficient examples per category, enabling more robust results. The results are reported in Table 11. Consistent with our earlier findings, single-fact questions generally outperformed multi-fact questions. Notably, the highest accuracy was observed for EventAtWhatTime in single-fact questions and RelationDuration in multi-fact questions. This alignment with the results from the smaller dataset reinforces their significance and suggests that the smaller dataset serves as a reliable proxy for the larger one.

Impact of Graph Structure and order of facts on LLM Performance. In this experiment, we evaluated LLM performance across various combinations of graph structure and fact order. The results, presented in Table Table 12, reveal that the target_and_start_time ordering consistently yields the best performance across the expanded dataset, regardless of graph structure. Conversely, the shuffle ordering consistently underperforms across most graph structures.

D EVALUATION PROCESS

We adopted a structured approach to ensure consistent evaluation. The LLM prompts incorporate specific guidelines for output formatting, requiring a JSON structure with fields like ‘explanation’ and ‘answer’. This standardized output facilitated automated evaluation through parsing the JSON, extracting the answer field(s), and comparing to the golden label. Here are examples of instructions in the prompt (please see below for the full prompt):

Table 12: Impact of graph structure and sorting type on a larger set of ToT-Semantic.

Order of facts	BA	Complete	ER	SBM	SFN	Star	AWE	Rank Average
relation_and_start_time	73.42	52.03	64.98	61.45	81.93	74.32	90.36	2.00
shuffle	66.72	44.65	54.74	54.14	74.17	72.74	85.02	4.71
start_time_and_relation	67.55	46.31	57.76	55.72	77.86	72.14	88.48	4.00
start_time_and_target	68.60	46.61	58.96	55.95	78.31	70.78	88.63	3.29
target_and_start_time	78.54	65.74	72.82	74.32	83.36	78.84	91.11	1.00

Example from ToT-Semantic: Answer the following question based on the temporal facts assuming the facts are unidirectional. Output only a valid JSON string with two fields: “explanation” and “answer”. Do not output anything else. The explanation field contains your reasoning. The answer field contains the value corresponding to your final answer.

Example from ToT-Arithmetic: Return your answer as a JSON in the following format: JSON = “explanation”: <your step by step solution>, “answer”: “day_of_week”.

This prompting method ensured clear instructions for the model’s output format. Our experiments showed consistent adherence to these instructions, demonstrating the effectiveness of our prompt design and leading to a robust and straightforward evaluation process.

Example from ToT-Semantic

Prompt: Here is a set of temporal facts:

E92 was the R17 of E69 from 1980 to 1988.

E69 was the R17 of E24 from 1972 to 1981.

E59 was the R17 of E69 from 1999 to 2000.

E69 was the R90 of E42 from 1997 to 2001.

E11 was the R17 of E69 from 1977 to 1980.

E92 was the R17 of E69 from 1995 to 1999.

E69 was the R90 of E24 from 1977 to 1978.

E69 was the R90 of E42 from 1980 to 1982.

E69 was the R17 of E42 from 1986 to 1994.

E92 was the R90 of E69 from 1993 to 1996.

E69 was the R90 of E42 from 1974 to 1975.

E59 was the R90 of E69 from 1999 to 2000.

E69 was the R90 of E42 from 1983 to 1987.

Answer the following question based on the temporal facts assuming the facts are unidirectional. Output only a valid JSON string with two fields: “explanation” and “answer”. Do not output anything else. The explanation field contains your reasoning. The answer field contains the value corresponding to your final answer.

Question: Find the entity that was the R90 of E69 immediately after E92 was the R90 of E69.

Model response (Gemini 1.5 pro): {“explanation”: “E92 was the R90 of E69 from 1993 to 1996. The entity that was the R90 of E69 immediately after E92 was E59 from 1999 to 2000.”, “answer”: “E59” }

Label: E59

Example from ToT-Arithmetic

Prompt: If right now, it is 22:00 (24hr) in Location A (EST -0500), what time is it in Location B (PST -0800)? Answer with time and day in json. Eg: JSON = “explanation”: <your step by step solution>, “day”: “same_day/+x/-x”, “time”: “HH:MM:SS”, where day=same_day if the answer is the same day, or +x/-x, where x is the number of days before/after the start day, and time is in 24-hour format. Assume standard time without daylight saving for all timezones. If seconds aren’t specified, assume it to be 00.

Model response (GPT-4o): {“explanation”: “EST is 3 hours ahead of PST. To convert 22:00 EST to PST, subtract 3 hours from 22:00. This results in 19:00 PST.”, “day”: “same_day”, “time”: “19:00:00” }

Label: {“day”: “same_day”, “time”: “19:00:00” }

1026 E QUALITY OF INSTRUCTIONS

1027

1028

1029

1030

1031

1032

1033

1034

1035

1036

1037

F INSTRUCTIONS TO PARTICIPANTS

1038

1039

1040

1041

1042

1043

1044

1045

1046

1047

1048

1049

1050

1051

1052

1053

1054

1055

1056

1057

1058

1059

1060

1061

1062

1063

1064

1065

1066

1067

1068

1069

1070

1071

G LIMITATION AND FUTURE WORK

1072

1073

1074

1075

1076

1077

1078

1079

The current work has several limitations that provide avenues for future research:

Single-Sentence Time Anchoring . This benchmark focuses on scenarios where the start and end times of a fact are both mentioned within a single sentence. However, in real-world scenarios, temporal information can be spread across multiple sentences or even documents. It is worth noting that this setup is easily convertible to the more general case where temporal information can be spread across multiple sentences. While we chose to focus on the single-sentence setup for this initial work, future research could readily adapt the benchmark to the multi-sentence scenario and explore the challenges and opportunities it presents.

Time Arithmetic Benchmark Compilation

Thank you for participating in our eval hour to help us expand our dataset to cover all the categories of time arithmetic that we can think of.

Terminology:

- **Time arithmetic:** Calculations with time values, often involving years, months, days, hours, minutes, seconds.
- **Category:** A high-level category of time arithmetic operations, such as addition/-subtraction, time conversion, etc.
- **Examples:** Real-life sentences that fall into a category. For instance, "Today is 27 July 2020 and I was told that my furniture will be delivered to me in exactly 60 days from now. On what date will the furniture be delivered?" is an example of addition.

Goal: Our goal is to cover as many real-life categories and subcategories related to time arithmetic as possible. We also want each subcategory to have multiple different real-life examples.

Levels of Importance of Contributions:

1. Discovering/adding a new category.
2. Adding new real-life examples within a subcategory (please contribute more in less densely populated areas).

Corner cases are useful, but please don't focus all your time on them. Discovering broader categories would be the most useful!

Please try to add new examples which are as different from existing ones as possible.

Thanks!

1080 **Exclusive Focus on Explicit Temporal Facts (By Design).** This benchmark intentionally focuses
1081 solely on explicit temporal facts (those with clear time anchors), excluding static facts (those without
1082 time anchors). This deliberate choice was made to ensure the benchmark specifically targets and
1083 assesses models' capabilities in temporal reasoning. However, future work could expand the scope to
1084 include static facts, offering a more comprehensive evaluation of both temporal and general factual
1085 reasoning.

1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133