
Benchmarking Temporal Reasoning: Can Large Language Models Navigate Time When Stories Refuse to Follow a Straight Line?

Feifei Sun¹ Ziyi Tong¹ Houjing Wei¹ Cheng Peng¹ Teeradaj Racharak²

Minh Le Nguyen¹

¹Japan Advanced Institute of Science and Technology (JAIST), Japan,

²Advanced Institute of So-Go-Chi (Convergence Knowledge) Informatics, Tohoku University, Japan,

Abstract

Temporal reasoning remains a challenging task for Large Language Models (LLMs), particularly when confronted with nonlinear narratives and mixed time systems, where events are presented out of chronological order. While human cognition effortlessly reconstructs temporal sequences in such narratives, LLMs often exhibit inconsistent reasoning and fail to infer the correct event order. In this paper, we present a comprehensive study on sentence-level event ordering to evaluate emerging frontier LLMs in temporal reasoning tasks. We contribute (i) a novel dataset derived from historical records, blending absolute and relative time expressions across varied granularities; (ii) a benchmark covering emerging frontier LLMs including GPT family, DeepSeek series, Qwen models, and open-source models; and (iii) an absolute-relative time conversion table to support future research on mixed time systems. Our experiments¹ reveal substantial limitations across current models, with a consistent performance decline when relative time disrupts chronological signals. We further provide a detailed benchmark analysis across multiple dimensions, including model types, sentence length, temporal granularity, and format violations. Our findings offer key insights and valuable resources to advance temporal reasoning research in LLMs.

1 Introduction

Temporal reasoning is a fundamental component of natural language understanding, underpinning applications such as question answering, narrative comprehension, and timeline construction. Despite rapid progress in Large Language Models (LLMs), reasoning over temporal sequences—especially within nonlinear narratives—remains a persistent challenge. Unlike humans, who can effortlessly reconstruct event orders from fragmented or non-chronological inputs, LLMs often struggle when faced with mixed time systems involving both absolute and relative time expressions.

Nonlinear narratives, characterized by disrupted temporal flow and interleaved time references, are common in historical texts, biographies, and storytelling. These contexts require models not only to interpret explicit time expressions but also to infer implicit event dependencies across varying temporal granularities (e.g., year, month, day). While existing benchmarks have explored temporal reasoning through question answering or multi-task datasets Jia et al. [2018], Qin et al. [2021], Chu et al. [2023], Wang and Zhao [2023], Tan et al. [2023], they often underrepresent event ordering as a standalone capability. As LLMs continue to advance, dedicated benchmarks for this fundamental yet

¹<https://github.com/fantastic-Feifei/MTS-benchmark>

fragile skill—particularly under naturalistic and temporally ambiguous conditions—are increasingly needed.

In this work, we address this gap by formulating sentence-level event ordering as a core temporal reasoning task under nonlinear narrative settings. We construct a benchmark derived from historical records sourced from Wikidata, where each sentence is temporally anchored and spans a range of granularities. To simulate realistic narrative complexity, we include both absolute and relative time expressions, capturing scenarios where temporal cues are implicit, vague, or mixed.

We evaluate a suite of leading frontier LLMs, including models from the GPT, DeepSeek, Qwen, and LLaMA families, along with Mistral-7B, focusing on their ability to recover event order, recognize temporal dependencies, and reason effectively under disrupted chronological signals.

To support future research, we also release a curated table of over 6,000 absolute-to-relative time expression that links structured time expressions (e.g., “1945”) with natural references (e.g., “the end of World War II”), offering a reusable resource for investigating mixed-time systems.

Our work makes the following contributions: we propose sentence-level event ordering as a benchmark task for evaluating temporal reasoning in nonlinear narratives; we construct a novel dataset based on historical texts, enriched with both absolute and relative time annotations across varied temporal granularities; we present a comprehensive benchmark study involving both leading frontier models (e.g., GPT-4, Deepseek, QWQ) and strong open-source baselines (e.g., LLaMA 3.3, Mistral, LLaMA 2-13B), systematically evaluating their ability to reason over mixed time systems; and we release an absolute-relative time conversion table to support further research in temporal inference.

Guided by these contributions, we investigate the following research questions:

- *How do different model architectures perform in temporal reasoning tasks?*
- *How do temporal granularity and event sequence length influence reasoning accuracy?*
- *Is there an interaction between time type and reasoning complexity?*
- *To what extent do relative time expressions affect model performance?*

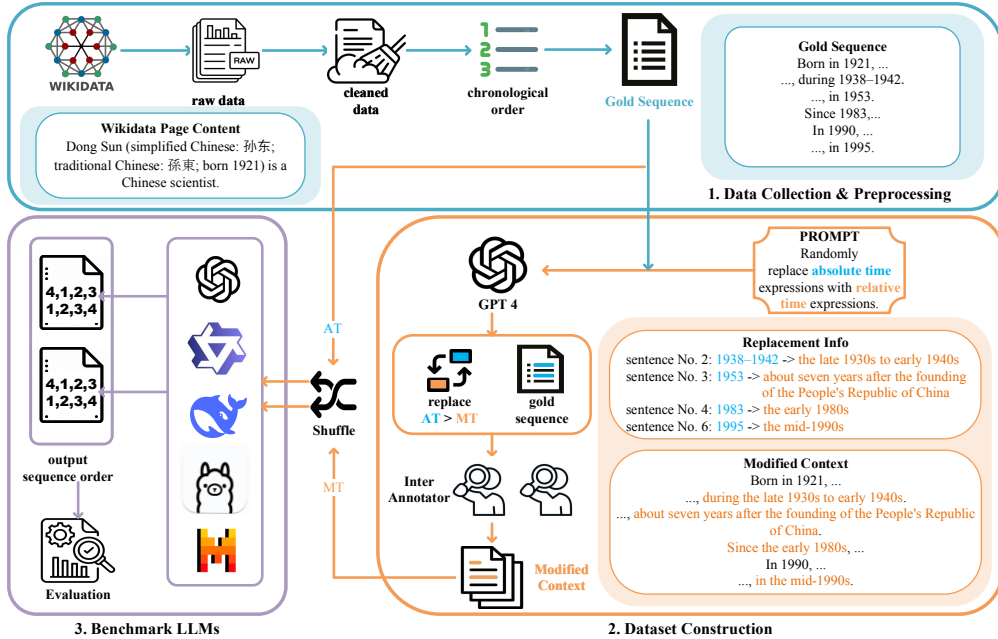


Figure 1: Overview of our benchmark pipeline. (1) we collect and clean biographical passages from wikidata, extracting temporally anchored sentences to form gold-standard event sequences. (2) a subset of absolute time expressions is rewritten into natural relative expressions (via GPT-4), producing modified contexts with replacement mappings, which are then validated by annotators. (3) multiple llms are evaluated on sentence-level event ordering under absolute-time (AT) and mixed-time (MT) settings, where models output a comma-separated index list (e.g., 2,1,4,3).

2 Related Work

Temporal Question Answering Temporal reasoning (TR) has long been recognized as a core challenge in natural language processing, essential for tasks involving event sequencing, duration inference, and causal understanding. Early QA-style benchmarks, such as TempQuestions Jia et al. [2018] and TimeDial Qin et al. [2021], focus on reasoning under explicit, implicit, and ordinal temporal constraints. Other datasets, like that of Chen et al. [2021], explore temporal drift through Wikipedia–Wikidata alignment, revealing the sensitivity of language models to subtle time-based context changes. TempReason Tan et al. [2023] expands the temporal QA paradigm to a multi-level framework, encompassing time-time, time-event, and event-event reasoning. This line of work demonstrates the increasing complexity of temporal understanding required by modern QA systems.

However, while these QA datasets reflect diverse forms of temporal reasoning, they often embed event ordering as a latent step within broader reasoning chains, making it difficult to isolate and evaluate this capability directly. In contrast, our work treats event ordering as a standalone task, enabling focused assessment of model performance under temporally ambiguous and nonlinear narrative conditions.

Comprehensive Temporal Benchmarks Recent benchmarks such as TimeBench Chu et al. [2023] and TRAM Wang and Zhao [2023] evaluate a broad spectrum of temporal reasoning skills by combining multiple tasks—such as duration estimation, temporal arithmetic, frequency detection, and causal inference—into large-scale evaluation suites. TempReason Tan et al. [2023] adopts a more structured design with three reasoning levels, but remains grounded in the question answering paradigm.

In contrast, we focus on sentence-level event ordering—an underexplored yet challenging sub-task—under hybrid time conditions that mix absolute and relative expressions. This design enables a finer-grained evaluation of LLMs’ ability to recover global temporal structure from fragmented, nonlinear narratives.

While existing work has addressed absolute or relative temporal reasoning in isolation, the distinct challenges of mixed time—such as implicit anchoring, granularity mismatch, and nonlinearity—remain underexplored. We outline these issues and their implications for benchmark construction in Section 3.2.

Instruction Sensitivity and Model Coverage Recent work has shown that instruction tuning alone may not ensure reliable execution of structured or temporally grounded tasks Lou et al. [2024], especially in scenarios requiring compositional reasoning or strict output format adherence Chia et al. [2023], Wang et al. [2022], Xu et al. [2023]. Although instruction-tuned models demonstrate strong performance in QA and classification, they often struggle in tasks demanding sequence-level reasoning or alignment with latent structural constraints Peng et al. [2023], Min et al. [2023].

Our benchmark contributes to this line of research by providing a comparative analysis of instruction-following behaviors across model families—including underexplored but high-performing models such as DeepSeek and Qwen—under temporally sensitive, zero-/one-shot prompting settings. While many prior studies focus on GPT-family models or open-domain QA tasks Kimura et al. [2021], Chen et al. [2021], Saxena et al. [2021], Dhingra et al. [2022], Tan et al. [2023], Gupta et al. [2023], Jia et al. [2024], Xiong et al. [2024], Fatemi et al. [2024], Deroy and Maity [2024], Su et al. [2024], Yuan et al. [2024], Zhang et al. [2024], Deng et al. [2024], Ruiz et al. [2025], recent open-source models like DeepSeek and Qwen—despite their strong reasoning capabilities—remain underexplored in temporal settings. Our benchmark fills this gap by providing targeted evaluations of instruction-following behavior across both frontier and open models under mixed-time conditions.

3 Benchmark Setup

3.1 Task Overview

We formulate temporal reasoning in nonlinear narratives as a sentence-level event ordering task. Given a short passage composed of n unordered sentences $P = \{s_1, s_2, \dots, s_n\}$, where each s_i describes an event associated with a time expression t_i , the model is tasked with inferring the correct

chronological order of the events. The time expressions can be absolute (e.g., “in 1923”) or relative (e.g., “three years later”), or a combination of both.

The expected output is a permutation π over the indices $\{1, \dots, n\}$ such that the reordered sequence $\{s_{\pi(1)}, s_{\pi(2)}, \dots, s_{\pi(n)}\}$ respects the underlying temporal timeline implied by the input. This task requires interpreting time expressions, resolving references, and aligning events across possibly fragmented or non-chronological inputs.

To operationalize this task, we construct a benchmark dataset that provides both absolute-time (AT) and mixed-time (MT) input settings, along with gold chronological sequences used for evaluation. An overview of the full data construction and evaluation pipeline is shown in Figure 1.

3.2 Challenges of Mixed Temporal Reasoning

Temporal reasoning in mixed time systems introduces challenges beyond standard timeline inference. First, relative expressions (e.g., “the following year”) require anchoring to implicit reference points, which are often unstated. Second, absolute and relative expressions may co-occur, requiring joint interpretation and temporal alignment. Third, varying temporal granularity—some events given as years, others as full dates—creates ambiguity in sequencing. Finally, nonlinear narratives frequently present events out of order, demanding global integration of dispersed time cues.

Table 1: Time expression types used in our benchmark. The latter two categories are treated as *relative* for MT setting.

Expression Type	Example
Absolute Time	“in 1945”, “in March 2007”, “on July 20, 1969”
Relative Time	“three years later”, “shortly after the war”
Event-Anchored Time	“the end of World War II”, “during the Great Depression”

3.3 Experimental Factors

To systematically investigate how different aspects of temporal structure affect model performance, we design benchmark settings along the following dimensions:

Mixed time expressions: introducing temporal ambiguity by randomly replacing a subset of absolute time expressions with relative references using an LLM-based rewriting strategy. We allow minor imprecision or implicit temporal references—such as GPT-4 occasionally grounding expressions like “this year” as 2023 irrespective of narrative context—as long as they do not alter the overall event order. This design choice reflects the inherent ambiguity in mixed-time narratives and evaluates whether models can still recover global chronological structure under such conditions.

Temporal granularity: comparing passages with coarse-grained (year-only) versus fine-grained (month or day included) time annotations.

Event sequence length: varying the number of events from 4 to 40 to examine how model performance scales with narrative length, and whether reasoning abilities degrade as the temporal chain becomes longer.

These experimental factors enable a fine-grained analysis of model sensitivity to temporal complexity under diverse and naturalistic conditions.

3.4 Dataset Construction

We construct our dataset from Wikidata Vrandečić and Krötzsch [2014] by extracting 15,000 historical and contemporary figures born after 1900, focusing on occupations such as scientists, historians, and politicians to ensure temporal and professional diversity. For each entity, we retrieve the English Wikipedia page and extract time-anchored event sentences using regex-based patterns. Sentences are filtered for grammaticality, relevance, and valid absolute dates, then chronologically sorted to form gold-standard event sequences. We retain passages containing 4 to 40 events to balance sequence complexity and data coverage. The detailed dataset construction process is provided in Appendix C.

Table 2: Comparison of dataset statistics before and after conversion.

Statistic	Value	Statistic	Value
Total passages	4,824	Total passages	4,824
Avg. events per passage	7.99	Avg. relative per passage	4.53
Temporal granularity — year	70.72%	Avg. absolute per passage	3.46
Temporal granularity — month	23.46%	Relative time ratio (relative / all)	56.73%
Temporal granularity — day	5.82%		

(a) Absolute-time dataset.

(b) Mixed-time dataset.

To simulate mixed-time narratives, we randomly convert a subset of absolute expressions into relative or descriptive forms using GPT-4o. A controlled prompt ensures the rewrites are semantically faithful and logically consistent with surrounding context. To assess the quality of these rewrites, two NLP expert annotators—also co-authors of this work—independently evaluate 200 randomly sampled passages on three dimensions: (i) *Info Accuracy*, (ii) *Context Logic*, and (iii) *Naturalness*. Agreement scores are high for accuracy (79.5%) and contextual coherence (71.5%), while naturalness exhibits moderate variance (quadratic weighted Cohen’s $\kappa = 0.19$). These results confirm that most rewritten expressions are reliable for constructing mixed-time inputs.

The final dataset comprises 4,824 passages with an average of 8 events each. In the mixed-time setting, 56.7% of expressions are rewritten as relative forms. Distributions by event count and temporal granularity are shown in Table 2 and Table 3. We also release a time expression conversion table (e.g., “1945” \rightarrow “the end of World War II”) to support future work on temporal paraphrasing and normalization (see Appendix D).

3.5 Benchmark Settings and Models

We evaluate LLMs on a sentence-level temporal ordering task. Given a passage with shuffled event sentences, the model must predict the correct chronological order as a permutation of sentence indices. We define two task variants:

Absolute-Time Task (AT): Passages contain only absolute time expressions (e.g., “in 1945”).

Mixed-Time Task (MT): Some absolute expressions are rewritten as relative references (e.g., “the end of World War II”) using a GPT-based strategy. See Table 1 for the formal definition of time expression types.

All models are evaluated using a one-shot instruction-style prompt with a single illustrative example (Appendix A). We include both closed-source and open-source models:

Closed-source Frontier Models: Including GPT-4, GPT-3.5, Deepseek-v3 Liu et al. [2024], Deepseek-r1 [Guo et al., 2025], Qwen2.5-7B [Qianwen et al., 2024], and QwQ-32B [Team, 2025], are accessed via official APIs, where hardware details are unavailable.

Open-source Models: Including LLaMA3.3-70B [Grattafiori et al., 2024], LLaMA2-13B [Touvron et al., 2023], and Mistral-7B Jiang et al. [2023], are run locally with the Ollama framework on a VM with 32 CPUs, 128 GB RAM, and one-quarter of an NVIDIA H100 GPU.

3.6 Evaluation Metrics

We report the following evaluation metrics:

Exact Match (EM): The percentage of outputs that exactly match the gold-standard permutation, reflecting the model’s ability to recover the *global temporal structure* of the passage. We additionally report the error rate, defined as $1 - \text{EM}$, which captures the proportion of incorrect predictions.

Kendall’s τ : Rank correlation between predicted and gold orders. This captures the *local temporal consistency* between event pairs.

Pairwise Accuracy: Fraction of correctly ordered sentence pairs.

We further apply:

Table 3: Distribution of passages by event count intervals. Most passages contain fewer than 10 events, aligning with the practical reasoning capacity of current LLMs, while longer passages are retained to test extended event ordering.

Event count range	# Passages	Percentage
4–9	3,817	79.13%
10–14	593	12.29%
15–19	184	3.81%
20–29	133	2.76%
30–39	56	1.16%
≥40	41	0.85%

Table 4: Performance comparison across models under absolute-time (AT) and mixed-time (MT) conditions. Percentage change is relative to AT: $\frac{MT-AT}{AT} \times 100\%$. Drops are highlighted in red, gains in green.

Model	EM (AT)	EM (MT)	Kendall’s τ (AT)	Kendall’s τ (MT)
QwQ-32B	0.54	0.33 (↓39%)	0.73	0.53 (↓27%)
Deepseek-r1	0.52	0.32 (↓38%)	0.70	0.53 (↓24%)
Deepseek-v3	0.33	0.21 (↓36%)	0.51	0.38 (↓25%)
GPT-4	0.31	0.15 (↓52%)	0.50	0.34 (↓32%)
LLaMA3.3-70B	0.21	0.13 (↓38%)	0.40	0.30 (↓25%)
GPT-3.5 turbo	0.12	0.07 (↓42%)	0.21	0.17 (↓19%)
Qwen2.5-7B	0.07	0.05 (↓29%)	0.20	0.14 (↓30%)
LLaMA2-13B	0.01	0.01 (↓0%)	0.00	0.05 (↑–)
Mistral-7B	0.00	0.01 (↑–)	0.05	0.06 (↑20%)

McNemar’s Test: For EM significance across AT and MT conditions.

Wilcoxon Signed-Rank Test: For Kendall’s τ significance across AT and MT.

Malformed outputs are excluded. We also analyze EM and Kendall’s τ scores by passage length and model family in Section 4.

Appendix F provides dataset visualizations, including event count distributions (Figure 8) and the temporal granularity of time expressions (Figure 7).

4 Results and Analysis

We analyze model performance on sentence-level event ordering under AT and MT conditions, covering nine models across proprietary and open-source families. Evaluation uses EM, Kendall’s τ , and significance testing to assess sensitivity to temporal ambiguity. We further analyze model performance from three key perspectives—temporal granularity, event sequence length, and the presence of relative time expressions—to systematically address our four research questions.

4.1 Overall Model Performance

To address our first research question concerning the performance of different model architectures in temporal reasoning tasks, we begin by comparing overall accuracy across all evaluated models.

Deeper Reasoning Comes at the Cost of Instruction Following

While both the Qwen and Deepseek model families achieve superior performance in temporal reasoning tasks compared to other models, we observe a notable divergence in output format adherence within each family. As shown in Figure 2, the stronger reasoning variants—QwQ-32B and Deepseek-r1—exhibit significantly higher rates of invalid format outputs than their smaller counterparts. This

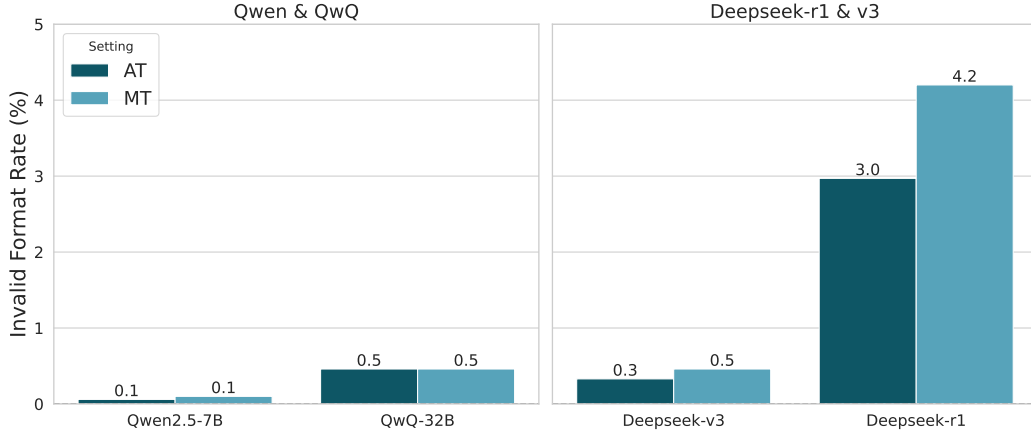


Figure 2: Invalid output rate (%) for Qwen and DeepSeek models under AT and MT. DeepSeek-r1 shows notably higher error rates, especially in MT, indicating reduced stability when processing relative time inputs.

pattern is consistent with findings from instruction-following literature Lou et al. [2024], which highlight that larger models, despite superior reasoning abilities, are more likely to deviate from strict output constraints—particularly in settings without strongly grounded demonstrations. An illustrative example of such a violation is provided in Appendix G.1.

These results reveal a trade-off between deep reasoning and strict instruction adherence. As models develop more complex inference capabilities, they may favor semantic interpretation over rigid output formatting, particularly under ambiguous prompts. This tension between interpretive depth and structural control is further evidenced by increased format violations, detailed in Appendix G.1.

4.2 Temporal Granularity Analysis

To address part of our second research question regarding the effect of temporal granularity on reasoning performance, we analyze how the granularity of time expressions influences model accuracy under both AT and MT conditions. Passages are grouped into two levels: those with only year-level expressions (*coarse-grained*) and those that include month or day annotations (*fine-grained*).

Fine-grained timestamps lead to more stable reasoning

Models perform more robustly on fine-grained passages, where temporal cues are more precise. These timestamps help disambiguate events that occur in the same year but at different times, enabling better alignment and control over sentence reordering.

To quantify this effect, we compare error rates between AT and MT across both granularity levels. As shown in Figure 10, the performance gap between AT and MT is consistently larger under coarse-grained inputs. For example, GPT-4 and QwQ-32B both show over 25% error rate increase when relative time replaces coarse absolute timestamps.

Stronger Models Within Families Are More Affected by Coarse-Grained Time Inputs

All models show performance degradation when temporal inputs are coarsened from day/month to year-level granularity. Notably, the strongest models—QwQ-32B and DeepSeek-r1—exhibit the largest MT–AT error increases under coarse-grained conditions (Figure 3), suggesting a reliance on fine-grained temporal cues. As specificity declines, these models may resort to overgeneralized reasoning, increasing deviation from the gold standard. This aligns with Yang et al. Yang et al. [2024], who show that temporally aware embeddings enhance reasoning but amplify sensitivity to time granularity. In contrast, weaker models appear less affected, likely due to simpler, more conservative reasoning. Detailed results are in Appendix G.2.

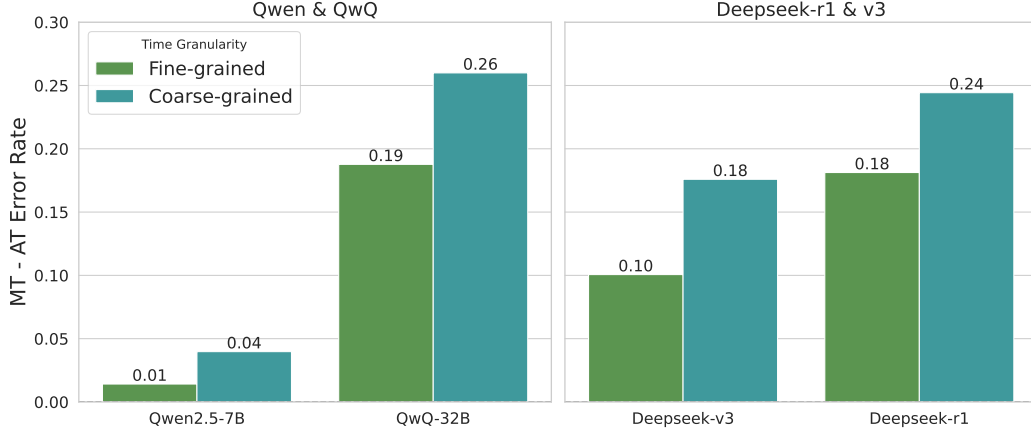
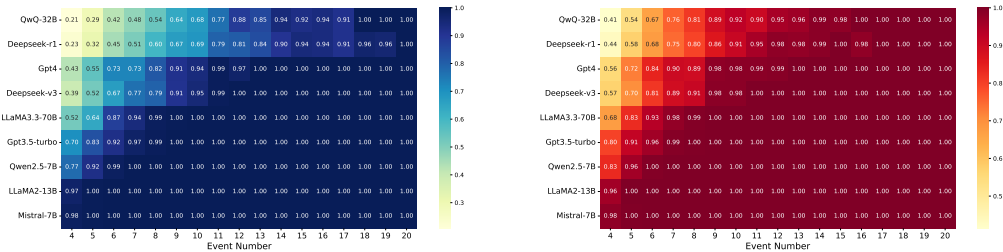


Figure 3: MT-AT error rate increase under different time granularities for Qwen and DeepSeek models. Both show greater degradation with coarse-grained inputs, with QwQ-32B and DeepSeek-r1 most affected, suggesting reduced robustness to underspecified temporal cues.

Relative time expressions are *less harmful* when granularity is high

The negative impact of switching to relative time is most severe under vague or underspecified temporal conditions. When time granularity is higher, relative expressions carry more specific temporal meaning—mitigating ambiguity and supporting more stable reasoning.

These findings highlight the interaction between surface-level time granularity and deeper temporal reasoning ability. Improving model robustness to coarse-grained relative time may require explicit training on relational semantics and underspecified narratives.



(a) AT Error Rates by Model and Event Number

(b) MT Error Rates by Model and Event Number

Figure 4: Comparison of AT and MT error rates across different models and event numbers. Error rate is defined as $1 - \text{Exact Match (EM)}$, representing the proportion of outputs that fail to exactly match the gold permutation.

4.3 Event Sequence Length Analysis

To address our third research question regarding the interaction between time type and reasoning complexity, we now examine how event sequence length influences temporal reasoning performance.

Longer sequences sharply degrade model performance

Figure 4 visualizes error rates for all evaluated models under both the AT and MT settings. We observe a clear trend: as the number of events rises, nearly all models experience a steady and often steep increase in error rate.

Most models begin with reasonably low error rates (e.g., 0.2–0.4) on short passages (4–6 events), particularly under the AT setting. However, accuracy degrades quickly, and by 12 events, even the best-performing models (e.g., GPT-4, Deepseek-R1, Qwen-32B) approach near-complete failure in the MT setting.

Relative time increases vulnerability to sequence length

The contrast between AT and MT is particularly striking: while AT error rates often increase more gradually, MT error rates rise faster and reach 1.0 earlier. This pattern reveals that relative time reasoning is disproportionately affected by sequence length—likely because models must track more implied temporal links without the support of explicit anchors.

Among all models, Qwen-32B and DeepSeek-R1 stand out for maintaining lower MT error rates in the 4–8 event range, while others such as Mistral and LLaMA variants fail almost immediately. The robustness of these models may stem from better generalization over temporal language, or implicit pretraining biases favoring temporal coherence.

By 15–20 events, nearly all models saturate at an error rate of 1.0 in both AT and MT conditions. These results indicate that existing models struggle to maintain coherence in long event chains, and relative-time reasoning becomes brittle under increased temporal complexity.

4.4 Absolute vs. Mixed Time Comparison

To directly address our last research question, we compare model performance between passages with AT and MT.

Mixed-Time Settings Introduce Substantial Difficulty

All models exhibit performance drops under the MT setting, though the magnitude varies. GPT-4 and GPT-3.5 Turbo experience steep EM reductions of over 50%, suggesting a strong reliance on explicit absolute-time cues. In contrast, frontier models like QwQ-32B and Deepseek-R1 show more graceful degradation, with EM drops around 20%, and Kendall’s τ remaining above 0.50.

Reliance on explicit timestamps amplifies degradation

Models like GPT-4 perform well under AT conditions but degrade sharply in MT, suggesting strong reliance on explicit date cues. In contrast, DeepSeek-R1 and QwQ maintain more stable performance, indicating better generalization to natural temporal variation. As shown in Table 4, color-coded drops in EM and Kendall’s τ highlight that smaller models (e.g., Mistral, LLaMA2-13B) not only perform poorly overall, but also show minimal AT–MT difference—suggesting weak temporal sensitivity. These findings underscore the need to evaluate LLMs under both controlled and realistic temporal settings to fully assess their reasoning capabilities.

5 Conclusion

In this work, we introduce a benchmark for evaluating LLMs’ temporal reasoning in narratives that mix absolute and relative time expressions. Unlike prior datasets restricted to a single time type or simplified settings, our benchmark reflects the hybrid temporal structures commonly found in real-world biographies.

Our experiments show that even strong models (e.g., QwQ-32B, DeepSeek-R1) struggle to maintain temporal coherence under mixed-time conditions, especially with coarse granularity or long-range dependencies. This suggests that current LLMs still rely heavily on surface cues and exhibit limited relational temporal reasoning.

To facilitate further research, we also release a large-scale aligned table of absolute-to-relative time conversions as a new resource for temporal normalization and contextual rewriting.

We hope this work motivates future advances in time-aware inference and constraint-driven model alignment, particularly for improving generalization in relative-time settings.

Acknowledgements

We would like to thank Naoya Inoue for providing valuable early-stage guidance and helping shape the initial direction of this research.

We thank all reviewers for their insightful and constructive feedback, which has helped us improve the clarity, precision, and completeness of this work.

Feifei Sun acknowledges the support of the China Scholarship Council (CSC) under scholarship number 202408050023.

Limitations

While our benchmark offers a robust platform for evaluating temporal reasoning in LLMs, several limitations remain.

First, the dataset is built from Wikipedia-style biographies, which—though rich in timestamped events—do not cover all narrative types. Domains such as scientific writing or fiction may exhibit different temporal patterns.

Second, we adopt a sentence-level event abstraction, omitting finer discourse phenomena like simultaneity or intra-sentential shifts. Time expressions are automatically extracted and occasionally noisy, which may affect alignment.

Third, relative expressions are generated via GPT-4o rewrites. While this improves lexical diversity, it introduces ambiguity—e.g., “2004” may become “the following year,” requiring prior context. Annotators observed occasional grounding errors (e.g., “this year” interpreted as 2023), but such cases are accepted if event order is preserved.

Fourth, our evaluation focuses on global metrics (EM, Kendall’s τ), which may overlook partial correctness in passages with underspecified temporal cues.

Fifth, we evaluate models under zero- and one-shot prompting only, without fine-tuning or architectural changes (e.g., temporal embeddings), which may further improve performance.

We also observe frequent instruction-following failures in open-source models. Despite format constraints, models like Mistral-7B often produce verbose outputs. One-shot prompting improves compliance, but we do not compare prompting strategies systematically due to budget constraints.

Finally, performance collapses on very long passages (e.g., >30 events), likely due to compounded reasoning and context-length challenges. These cases are excluded from analysis and underscore the need for better long-context temporal reasoning.

References

- Zhen Jia, Abdalghani Abujabal, Rishiraj Saha Roy, Jannik Strötgen, and Gerhard Weikum. Tempquestions: A benchmark for temporal question answering. In *Companion Proceedings of the The Web Conference 2018*, pages 1057–1062, 2018.
- Lianhui Qin, Aditya Gupta, Shyam Upadhyay, Luheng He, Yejin Choi, and Manaal Faruqui. Timedial: Temporal commonsense reasoning in dialog. *arXiv preprint arXiv:2106.04571*, 2021.
- 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.
- Yuqing Wang and Yun Zhao. Tram: Benchmarking temporal reasoning for large language models. *arXiv preprint arXiv:2310.00835*, 2023.
- Qingyu Tan, Hwee Tou Ng, and Lidong Bing. Towards benchmarking and improving the temporal reasoning capability of large language models. *arXiv preprint arXiv:2306.08952*, 2023.
- Wenhu Chen, Xinyi Wang, and William Yang Wang. A dataset for answering time-sensitive questions. *arXiv preprint arXiv:2108.06314*, 2021.

- Renze Lou, Kai Zhang, and Wenpeng Yin. Large language model instruction following: A survey of progresses and challenges. *Computational Linguistics*, 50(3):1053–1095, 2024.
- Yew Ken Chia, Pengfei Hong, Lidong Bing, and Soujanya Poria. Instructeval: Towards holistic evaluation of instruction-tuned large language models. *arXiv preprint arXiv:2306.04757*, 2023.
- Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A Smith, Daniel Khashabi, and Hannaneh Hajishirzi. Self-instruct: Aligning language models with self-generated instructions. *arXiv preprint arXiv:2212.10560*, 2022.
- Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, and Daxin Jiang. Wizardlm: Empowering large language models to follow complex instructions. *arXiv preprint arXiv:2304.12244*, 2023.
- Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. Instruction tuning with gpt-4. *arXiv preprint arXiv:2304.03277*, 2023.
- Sewon Min, Kalpesh Krishna, Xinxu Lyu, Mike Lewis, Wen-tau Yih, Pang Wei Koh, Mohit Iyyer, Luke Zettlemoyer, and Hannaneh Hajishirzi. Factscore: Fine-grained atomic evaluation of factual precision in long form text generation. *arXiv preprint arXiv:2305.14251*, 2023.
- Mayuko Kimura, Lis Kanashiro Pereira, and Ichiro Kobayashi. Towards a language model for temporal commonsense reasoning. In *Proceedings of the Student Research Workshop Associated with RANLP 2021*, pages 78–84, 2021.
- Apoorv Saxena, Soumen Chakrabarti, and Partha Talukdar. Question answering over temporal knowledge graphs. *arXiv preprint arXiv:2106.01515*, 2021.
- Bhuwan Dhingra, Jeremy R Cole, Julian Martin Eisenschlos, Daniel Gillick, Jacob Eisenstein, and William W Cohen. Time-aware language models as temporal knowledge bases. *Transactions of the Association for Computational Linguistics*, 10:257–273, 2022.
- Vivek Gupta, Pranshu Kandoi, Mahek Bhavesh Vora, Shuo Zhang, Yujie He, Ridho Reinanda, and Vivek Srikumar. Temptabqa: Temporal question answering for semi-structured tables. *arXiv preprint arXiv:2311.08002*, 2023.
- Zhen Jia, Philipp Christmann, and Gerhard Weikum. Faithful temporal question answering over heterogeneous sources. In *Proceedings of the ACM Web Conference 2024*, pages 2052–2063, 2024.
- Siheng Xiong, Ali Payani, Ramana Kompella, and Faramarz Fekri. Large language models can learn temporal reasoning. *arXiv preprint arXiv:2401.06853*, 2024.
- Bahare Fatemi, Mehran Kazemi, Anton Tsitsulin, Karishma Malkan, Jinyeong Yim, John Palowitch, Sungyong Seo, Jonathan Halcrow, and Bryan Perozzi. Test of time: A benchmark for evaluating llms on temporal reasoning. *arXiv preprint arXiv:2406.09170*, 2024.
- Aniket Deroy and Subhankar Maity. A short case study on understanding the capabilities of gpt for temporal reasoning tasks. *Authorea Preprints*, 2024.
- Zhaochen Su, Juntao Li, Jun Zhang, Tong Zhu, Xiaoye Qu, Pan Zhou, Yan Bowen, Yu Cheng, et al. Living in the moment: Can large language models grasp co-temporal reasoning? *arXiv preprint arXiv:2406.09072*, 2024.
- Chenhan Yuan, Qianqian Xie, Jimin Huang, and Sophia Ananiadou. Back to the future: Towards explainable temporal reasoning with large language models. In *Proceedings of the ACM Web Conference 2024*, pages 1963–1974, 2024.
- Xinliang Frederick Zhang, Nick Beauchamp, and Lu Wang. Narrative-of-thought: Improving temporal reasoning of large language models via recounted narratives. *arXiv preprint arXiv:2410.05558*, 2024.
- Irwin Deng, Kushagra Dixit, Vivek Gupta, and Dan Roth. Enhancing temporal understanding in llms for semi-structured tables. *arXiv preprint arXiv:2407.16030*, 2024.

- Alfredo Garrachón Ruiz, Tomás de la Rosa, and Daniel Borrajo. On the temporal question-answering capabilities of large language models over anonymized data. *arXiv preprint arXiv:2504.07646*, 2025.
- Denny Vrandečić and Markus Krötzsch. Wikidata: a free collaborative knowledgebase. *Communications of the ACM*, 57(10):78–85, 2014.
- Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*, 2024.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- Peng Qianwen, Gao Yanzipeng, Li Xiaoqing, Min Fanke, Li Mingrui, Wang Zhichun, and Liu Tianyun. . In Hongfei Lin, Hongye Tan, and Bin Li, editors, *Proceedings of the 23rd Chinese National Conference on Computational Linguistics (Volume 3: Evaluations)*, pages 294–301, Taiyuan, China, July 2024. Chinese Information Processing Society of China. URL <https://aclanthology.org/2024.ccl-3.33/>.
- Qwen Team. Qwq-32b: Embracing the power of reinforcement learning, March 2025. URL <https://qwenlm.github.io/blog/qwq-32b/>.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- Albert Q Jiang, A Sablayrolles, A Mensch, C Bamford, D Singh Chaplot, Ddl Casas, F Bressand, G Lengyel, G Lample, L Saulnier, et al. Mistral 7b. arxiv. *arXiv preprint arXiv:2310.06825*, 10, 2023.
- Wanqi Yang, Yanda Li, Meng Fang, and Ling Chen. Enhancing temporal sensitivity and reasoning for time-sensitive question answering. *arXiv preprint arXiv:2409.16909*, 2024.

A Prompt Templates

We provide below the two prompt templates used in our study: one for rewriting absolute time expressions into relative ones, and another for evaluating temporal reasoning via event ordering. Both prompts follow a standardized instruction style to ensure consistency across model families.

(1) Relative Time Conversion Prompt

Time Replacement Prompt

You are a time conversion assistant. Your task is to replace exactly `{num_to_keep}` absolute time expressions with relative time expressions.

- Absolute time refers to any date in year, month-year, or full-date format.
- Retain `{num_to_keep}` absolute times, convert the rest into natural relative references.
- Avoid repeating the same phrasing.
- Do not simply compute or state time differences.

Return:

- **Modified Context:** the rewritten passage.
- **Replacement Information:** lines showing original \rightarrow relative expressions.

(2) Event Ordering Prompt (Benchmark)

One-Shot Benchmark Prompt

The following is a set of shuffled sentences. Please infer the correct order and return the sentence order as a sequence of numbers.

- Instructions:** - Only return a comma-separated sequence of numbers.
- Do not include any explanations, additional text, or line breaks.
- The sequence should reflect the correct order of the given sentences.

Example:

Input:

1. The sun rises in the east.
2. It is early morning.
3. The birds are singing.

Correct output:

2,1,3

Now, process the following sentences:

`{context}`

Please output only the sequence of numbers.

B Example of Relative Time Conversion

Below are three representative examples showing how absolute time expressions are converted into relative expressions using our prompt-based generation pipeline.

We acknowledge that certain time replacements (e.g., replacing “2015” with “a few years after joining MIT”) may introduce implicit event dependencies, such as the need to infer the timing of the prior event (i.e., joining MIT). However, our task primarily evaluates whether models can recover the correct chronological order of events rather than verifying precise temporal anchoring of each individual expression.

Therefore, as long as the replacement does not alter the relative order of events in the passage, such substitutions are considered acceptable within our task framework. These more ambiguous or indirect expressions are intentionally included to simulate the diversity and complexity of naturally occurring narratives with mixed temporal expressions.

Example 1

Original Passage (Gold Sequence)

- (1) His father, Babalyk, born in 1860, was the only child in the family.
- (2) He studied at a Kazakh school, then in the Tatar language school, then in **1941–1943** he graduated from the gymnasium in the city of Tacheng.
- (3) In **1943–1947**, while studying at the university in Ürümqi, he was arrested for nationalist actions and imprisoned.
- (4) After the founding of the Communist State, he became governor of Ili Kazakh Autonomous Prefecture in June 1955, and held that office until 1958.

Converted Passage (Mixed Time Expressions)

- (1) His father, Babalyk, born in 1860, was the only child in the family.
- (2) He studied at a Kazakh school, then in the Tatar language school, then **during the early 1940s** he graduated from the gymnasium in the city of Tacheng.
- (3) **Around the mid-1940s**, while studying at the university in Ürümqi, he was arrested for nationalist actions and imprisoned.
- (4) After the founding of the Communist State, he became governor of Ili Kazakh Autonomous Prefecture in June 1955, and held that office until 1958.

Replacement Mapping

1. Sentence 2: **1941–1943** → *during the early 1940s*
2. Sentence 3: **1943–1947** → *around the mid-1940s*

Example 2

Original Passage (Gold Sequence)

- (1) Zaharia was a gold medalist at the International Collegiate Programming Contest, where his team University of Waterloo placed fourth in the world and first in North America in 2005.
- (2) While at University of California, Berkeley's AMPLab in **2009**, he created Apache Spark as a faster alternative to MapReduce.
- (3) In 2013 Zaharia was one of the co-founders of Databricks where he is chief technology officer.
- (4) He joined the faculty of MIT in **2015**, and then became an assistant professor of computer science at Stanford University in 2016.
- (5) In 2019 he was spearheading MLflow at Databricks, while still teaching.

Converted Passage (Mixed Time Expressions)

- (1) Zaharia was a gold medalist at the International Collegiate Programming Contest, where his team University of Waterloo placed fourth in the world and first in North America in 2005.
- (2) While at University of California, Berkeley's AMPLab **several years later**, he created Apache Spark as a faster alternative to MapReduce.
- (3) In 2013 Zaharia was one of the co-founders of Databricks where he is chief technology officer.
- (4) **A few years after joining MIT**, he became an assistant professor of computer science at Stanford University in 2016.
- (5) In 2019 he was spearheading MLflow at Databricks, while still teaching.

Replacement Mapping

1. Sentence 2: 2009 → *several years later*
2. Sentence 4: 2015 → *a few years after joining MIT*

Example 3

Original Passage (Gold Sequence)

- (1) In 1937, Schulze moved to Peenemünde Army Research Center; in 1939, he was appointed chief of the Propulsion Unit, a position he held until 1945.
- (2) Classified as wards of the state, the seven men landed at Fort Strong on September 29, 1945; all but von Braun, Schulze included, were then transferred to Aberdeen Proving Ground to translate and catalog 14 tons of V-2 documents taken from Germany.
- (3) By 1946, Schulze was among the Operation Paperclip scientists employed at Fort Bliss.
- (4) He moved to Alabama, where he was naturalized in Birmingham on November 11, 1954.

Converted Passage (Mixed Time Expressions)

- (1) In 1937, Schulze moved to Peenemünde Army Research Center; in 1939, he was appointed chief of the Propulsion Unit, a position he held until **the end of World War II**.
- (2) Classified as wards of the state, the seven men landed at Fort Strong **during the late 1940s**; all but von Braun, Schulze included, were then transferred to Aberdeen Proving Ground to translate and catalog 14 tons of V-2 documents taken from Germany.
- (3) By the year after World War II ended, Schulze was among the Operation Paperclip scientists employed at Fort Bliss.
- (4) He moved to Alabama, where he was naturalized in Birmingham on November 11, 1954.

Replacement Mapping

1. Sentence 1: 1945 → *the end of World War II*
2. Sentence 2: September 29, 1945 → *during the late 1940s*

C Dataset Construction

C.1 Source and Collection

We construct our dataset from Wikidata by targeting historical and contemporary figures born after 1900. Using a SPARQL-based extraction pipeline, we retrieve 15,000 entities that satisfy the following criteria: (1) the entity must be an instance of human (Q5), (2) it must have a known birth date, and (3) it must be linked to an English Wikipedia page. To focus on domains rich in temporally annotated life events, we further filter by occupation, including scientists, historians, politicians, military personnel, inventors, engineers, Nobel laureates, astronauts, heads of state, journalists, and Olympic champions.

Each extracted entity includes its Wikidata ID, name, Wikipedia title, and article URL. These records serve as entry points for downstream extraction of event-time pairs, which form the core of our temporal reasoning benchmark.

C.2 Event and Time Extraction

For each entity collected from Wikidata, we retrieve its corresponding English Wikipedia page and process the article text to extract individual events anchored by absolute time expressions. We use regular expressions to identify a broad range of temporal formats, including full dates (e.g., “January 1, 2000”), partial dates (e.g., “April 1985”), ISO date formats (e.g., “2001-05-17”), and standalone years (e.g., “1942”).

We clean the raw Wikipedia text by removing titles, citations, markup, and references. The text is then segmented into sentences using rule-based heuristics. To ensure sentence quality and relevance, we retain only those sentences that: (1) contain a recognized date, (2) begin with a capital letter and end with a period, and (3) include at least one verb indicative of an event (e.g., “was born”, “joined”, “studied”, “published”).

Each retained sentence is treated as an individual event. For every passage, we generate a gold-standard event sequence by sorting the extracted sentences in chronological order based on the detected time expressions. To ensure the task remains non-trivial and meaningful for ordering-based evaluation, we discard passages with fewer than four valid event sentences. Filtering rules were designed for both robustness and precision, prioritizing sentences with well-formed structure and explicit temporal anchors.

This process results in a collection of time-anchored narratives with sentence-level event units and corresponding gold sequences, which serve as ground truth for model evaluation in the event ordering task.

C.3 Relative Time Generation

To simulate real-world narrative styles where absolute and relative time expressions often co-occur, we introduce controlled temporal ambiguity by converting some absolute time expressions into relative ones using a generation-based approach. Specifically, for each passage, we randomly sample a number of absolute time expressions to retain, and instruct a large language model (GPT-4o) to rewrite the remaining time references into natural relative expressions.

We design a specialized prompt that enforces the following constraints: (1) preserve exactly k absolute time expressions in the passage; (2) rewrite the remaining absolute times into relative references (e.g., “In 1923” → “Three years after his graduation”); (3) avoid simple numeric offsets or repetitive phrasing; and (4) return both the modified context and a mapping of replacements for traceability. The value of k is randomly chosen for each passage to encourage diversity in absolute-relative composition.

D Time Expression Conversion Table

To support future research on temporal rewriting and normalization, we release a conversion table that records all absolute-to-relative time expression rewrites applied during the construction of the our dataset. Each entry in the table represents a single replacement performed by GPT-4o during the mixed-time generation process.

Table 5 presents representative examples. The rewrites range from grounded historical interpretations (e.g., “1945” → “the end of World War II”) to relative references that depend on the surrounding narrative timeline (e.g., “2004” → “the following year”).

We caution that not all rewritten expressions are context-independent. While some rewrites refer to widely understood historical periods (e.g., “1945” → “the end of World War II”), others depend on the internal narrative timeline. For example, “2004” is sometimes rewritten as “the following year”, which is contextually appropriate only if the previous sentence refers to “2003”. Such replacements, though semantically coherent in context, may not be suitable for standalone use.

Moreover, during our double-annotation process (see Section 3.4 and Section E for details), we adopt a practical criterion: a rewritten time expression is considered valid as long as it preserves the overall temporal order of the passage, even if the substitution is not lexically precise. This design choice reflects our focus on evaluating temporal reasoning rather than surface-level rewriting fidelity.

We therefore encourage users to consult the context when applying this conversion table in downstream tasks such as generation, normalization, or rule extraction. The conversion table is best viewed as a supporting resource rather than a standalone ground truth.

For reproducibility, we defined a set of rule-based rewriting principles covering temporal granularity preservation, approximate time estimation, and contextual anchoring. These rules guided all absolute-to-relative conversions to ensure consistency across the dataset. A detailed description and implementation are provided in the released GitHub repository.

Table 5: Examples of absolute-to-relative time expression conversions, including grounded historical, approximate, and contextual rewrites.

Original time	Rewritten expression
1970	early 1970s
1979	late 1970s
2000	the turn of the millennium
2004	the following year
1967	several decades ago
1976	a little over four decades ago
1993	approximately three decades back
2009	fourteen years ago
2012	eight years ago
1948	a little over 75 years ago
1949	about 74 years back
1951	early 1950s
1956	approximately mid-1950s
1989	the last decade of the 1980s
Oct. 2, 2003	early October 2003
Jan. 23, 2004	late January 2004
Jan. 23, 2004	soon after
Jan. 28, 2004	at the end of January 2004
Feb. 1, 2004	shortly after

E Annotation Protocol and Analysis

To evaluate the quality of GPT-generated relative time expressions, we conducted a double annotation study on 200 sampled passages. The two annotators were the co-authors of this paper, both with expertise in NLP. Each annotator was presented with the original passage (*Gold Sequence*), the GPT-modified passage (*Gpt Modified Context*), and a detailed list of substitutions (*Replacement Info*).

Although some relative expressions are not exact translations of the original absolute timestamps, we consider the replacement acceptable as long as the temporal sequence of events remains unaffected. This evaluation criterion was reflected in the annotation guidelines for the “Info Accuracy” dimension. This decision aligns with our task definition, where the primary goal is to evaluate models’ ability to reconstruct the correct temporal order, rather than the surface accuracy of individual time expressions.

For each passage, annotators were instructed to evaluate the following three dimensions:

1. **Info Accuracy (Y/N)**: Whether the relative expression generated by GPT-4 accurately reflects the semantics of the original absolute timestamp.
 - **Y**: The relative time correctly corresponds to the absolute time and aligns with the provided substitution info.
 - **N**: The expression is semantically incorrect, overly vague, or omits critical temporal details.
2. **Context Logic (Y/N)**: Whether the modified relative expression fits logically and temporally within the surrounding passage.
 - **Y**: The expression is coherent in context and does not break the narrative or event sequence.
 - **N**: The expression introduces chronological contradictions or disrupts temporal flow.
3. **Naturalness Score (1–5)**: Fluency and readability of the modified sentence, regardless of correctness.
 - 5: Fully natural and indistinguishable from human-written text.
 - 4: Mostly fluent with only minor disfluency.
 - 3: Somewhat awkward but understandable.
 - 2: Clearly unnatural with evident phrasing issues.

Table 6: Example annotation showing the original passage, rewritten passage, replacement mapping, and evaluation fields.

Field	Example annotation
Original passage	In March 2007 she was elected to the fellowship of the Royal Society of Edinburgh. In 2018 she was appointed Head of the School of Informatics at Edinburgh, taking over from Johanna Moore, until succeeded by Helen Hastie in 2023. In 2018, Hillston was elected the membership of the Academia Europaea. Hillston was elected a Fellow of the Royal Society in May 2022. Since January 1st 2023 Hillston has been Editor-in-Chief of <i>Proceedings of the Royal Society A</i> .
Rewritten passage	In March 2007 she was elected to the fellowship of the Royal Society of Edinburgh. In 2018 she was appointed Head of the School of Informatics at Edinburgh, taking over from Johanna Moore, until succeeded by Helen Hastie this year . In 2018, Hillston was elected the membership of the Academia Europaea. Hillston was elected a Fellow of the Royal Society in May of last year . Since January 1st 2023 Hillston has been Editor-in-Chief of <i>Proceedings of the Royal Society A</i> .
Replacement mapping	Sentence 2: 2023 → this year Sentence 4: May 2022 → May of last year
Accuracy (Y/N)	Y
Coherence (Y/N)	N
Naturalness score (1–5)	3
Error type	InfoLoss
Free-form comment	Current year is assumed to be 2023, which disrupts contextual coherence.

- 1: Machine-like and syntactically poor.

Annotators were also encouraged to optionally tag common issues using a predefined label set:

- **infoless**: Key temporal information is missing.
- **vague**: Time span is ambiguous (e.g., “many years later”).
- **inconsistent**: Logical contradiction in event ordering.
- **HardUnderstand**: Converted sentence is semantically unclear.
- **Other**: Additional problems not captured by the above categories.

To ensure consistency, annotators jointly reviewed 5–10 initial examples and were encouraged to leave free-form comments for both high-quality and problematic samples. The estimated annotation time per passage ranged from 1–3 minutes.

Here, “free-form” refers to an open comment field in the annotation interface, where annotators could optionally write their reflections on the quality of time expression rewriting, such as naturalness, contextual alignment, or specific GPT-related issues.

E.1 Annotation Results and Agreement

E.1.1 Issue-Type Distribution and Annotator Differences.

Figure 5 summarizes the absolute counts of four issue types (HardUnderstand, Inconsistent, InfoLoss, Vague) labeled independently by annotator a and annotator b. Overall, the distributions reveal both areas of alignment and notable systematic differences. First, HardUnderstand and InfoLoss exhibit relatively low and comparable frequencies across annotators, suggesting stronger agreement in identifying linguistically opaque or information-dropping responses. In contrast, Inconsistent and especially Vague show substantial divergence. Annotator a assigned noticeably more Inconsistent

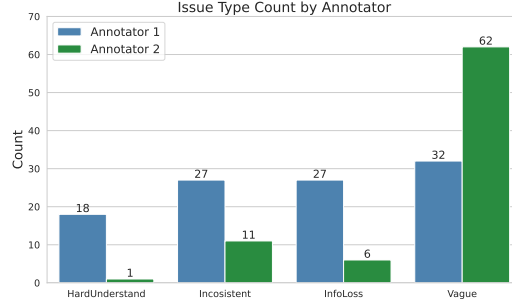


Figure 5: Absolute counts of issue types labeled by annotator a and annotator b.

labels (27 vs. 11), indicating a stricter sensitivity to logical contradiction or mismatch between context and model output. Conversely, annotator b labeled Vague far more frequently (62 vs. 32), reflecting a lower tolerance for underspecified or non-committal responses.

These discrepancies highlight differences in annotation granularity between the two annotators: annotator a tends to prioritize factual or logical coherence, whereas annotator b is more sensitive to lack of specificity. Importantly, such complementary biases help justify our double-annotation design, ensuring coverage of both logical correctness and pragmatic clarity. The subsequent agreement analysis accounts for these asymmetric tendencies, and disagreements were resolved through targeted discussion to achieve consistent final labels.

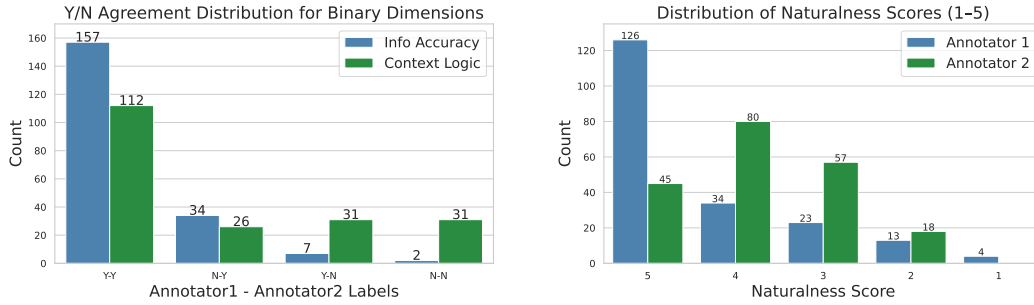


Figure 6: Annotation agreement patterns across evaluation dimensions.

E.1.2 Annotation Agreement Analysis Across Evaluation Dimensions.

Figure 6a illustrates the agreement distribution for the two binary dimensions—Information Accuracy and Context Logic—based on Y/N labels. For both dimensions, the majority of cases fall into the Y–Y category, with annotators jointly marking 157 instances as accurate for information and 112 instances as logically consistent. This concentration of positive agreement suggests that both annotators broadly shared similar criteria for identifying correct factual grounding and coherent event sequencing. However, the N–Y and Y–N counts reveal meaningful asymmetry: for Information Accuracy, annotator 1 produced more positive judgments (34 N–Y vs. 7 Y–N), while for Context Logic the asymmetry is reversed (26 N–Y vs. 31 Y–N). These patterns indicate subtle differences in annotator sensitivity—annotator 1 was stricter about contextual logic, whereas annotator 2 tended to be more conservative in judging factual accuracy.

Figure 6b presents the distribution of naturalness scores (1–5) assigned by the two annotators. Annotator 1 shows a strong skew toward higher scores, with the majority of responses rated as 5 (126 instances) and relatively few low scores. In contrast, annotator 2 exhibits a flatter distribution, assigning markedly fewer 5 scores (45) and substantially more 4 and 3 scores (80 and 57, respectively). This divergence suggests that annotator 2 applied a more fine-grained or stricter standard when evaluating fluency and stylistic naturalness, whereas annotator 1 more often regarded model outputs

as fully natural. Despite these differences in rating severity, both annotators show broadly aligned overall trends, with higher scores being far more frequent than lower ones.

Overall, the agreement patterns across both binary and scalar dimensions highlight complementary annotator tendencies: annotator 1 is comparatively generous in naturalness while stricter in logic; annotator 2 is more sensitive to stylistic nuance but more conservative in accuracy judgments. These differences underscore the value of double annotation and guided adjudication in obtaining stable final labels.

E.2 Error Type Distribution

To better understand annotator preferences and tendencies in error labeling, we compare the absolute count of common issue types annotated by each annotator. As shown in Figure 5, Annotator A overwhelmingly labeled vague expressions (62 instances), while Annotator B distributed their annotations more evenly across multiple categories.

Specifically, Annotator B marked 27 instances each of `InfoLoss` and `Inconsistent`, as well as 18 instances of `HardUnderstand`, compared to Annotator A’s respective counts of 6, 11, and 1. These differences suggest that Annotator A is particularly sensitive to ambiguity and imprecision in temporal phrasing, whereas Annotator B applies stricter standards in identifying information loss and logical inconsistency.

Despite these differences in emphasis, both annotators consistently identified problematic passages, reinforcing the value of error-type labels in guiding future improvements. The complementary nature of these annotation styles also offers useful insights into the diverse aspects of failure in GPT-based time rewriting.

Due to limited computational budget, we did not conduct adjudication to resolve annotation disagreements, which may leave some borderline cases open to interpretation.

Distribution of Temporal Granularity in Time Expressions

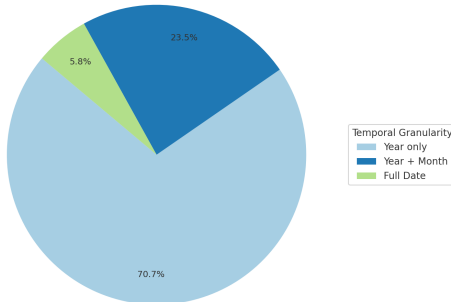


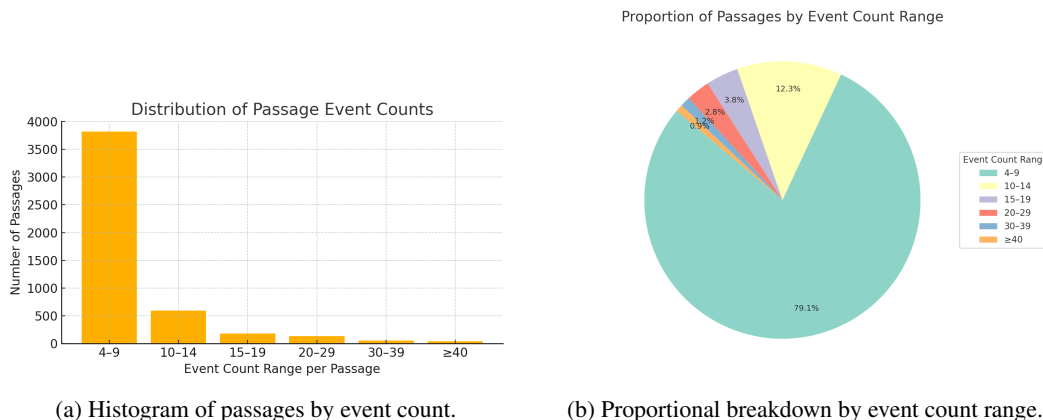
Figure 7: Distribution of temporal granularity among all absolute time expressions.

F Dataset Distribution Visualizations

To better illustrate the internal structure of our dataset, we present a set of visualizations that highlight the distribution of event counts and temporal granularity across all passages.

Figure 7 shows the distribution of temporal granularity for all absolute time expressions in the dataset. The majority (70.7%) of expressions specify only the year (e.g., “in 1987”), while 23.5% include both year and month (e.g., “July 1987”), and only 5.8% provide a full date (e.g., “July 15, 1987”). This skew toward coarse-grained time references reflects common patterns in Wikipedia-style biographical writing and suggests that many temporal relations must be inferred from sparse cues.

Figure 8 presents the distribution of event counts across passages. Figure 8a displays a histogram showing that most passages contain fewer than 10 events, with a peak in the 4–9 range. Figure 8b



(a) Histogram of passages by event count.

(b) Proportional breakdown by event count range.

Figure 8: Event count distribution across passages in the dataset.

provides a proportional breakdown, confirming that 79.1% of passages fall into the 4–9 event range. Passages with more than 20 events are relatively rare, accounting for less than 5% of the dataset.

This distribution suggests that the dataset is centered around passages with fewer than 10 events, maintaining a manageable level of complexity for most temporal reasoning tasks. At the same time, a small number of long-sequence passages (with 20 or more events) are included to support long-tail evaluation and stress-test models under extended temporal contexts.

G All model result

G.1 Format Violation Analysis

We first provide a representative example of a format violation. As shown in the box below, instead of returning a comma-separated list of sentence indices as instructed, the model outputs a verbose sequence of full event descriptions. This behavior constitutes a clear deviation from the expected format and illustrates a common failure mode among instruction-sensitive models. Such violations not only complicate automated evaluation but also indicate potential weaknesses in instruction comprehension, particularly when temporal reasoning is embedded in ambiguous inputs.

Example of Invalid Output Format

Gold Order: [5, 4, 1, 3, 2]
Expected Format: A comma-separated list of indices, e.g., 5, 2, 1, 4, 3
Model Output:

1970, She was established what was for nearly a decade the only protein crystallography laboratory in Israel., In 1970, Her parents ... Then, from 1979 to 1984 she was a group leader ... On Saturday, 18 October 2014, Professor Yonath ... She was visiting professor at the University of Chicago ...

Violation Type: Verbose explanation instead of index list

Figure 9 shows format violation rates under AT and MT settings. Violations include missing indices, extraneous text, or malformed outputs.

While most models perform well (violation <1%), **Mistral-7B** (51.3% AT, 34.0% MT) and **LLaMA2-13B** (14.7% AT, 12.7% MT) show significant instability. In contrast, models like Qwen2.5-7B, Deepseek-v3, and GPT-3.5-turbo maintain consistently low violation rates.

MT settings generally increase format errors, highlighting the destabilizing effect of relative expressions on instruction-following. Notably, Mistral-7B and LLaMA2-13B often generate verbose explanations instead of plain index lists.

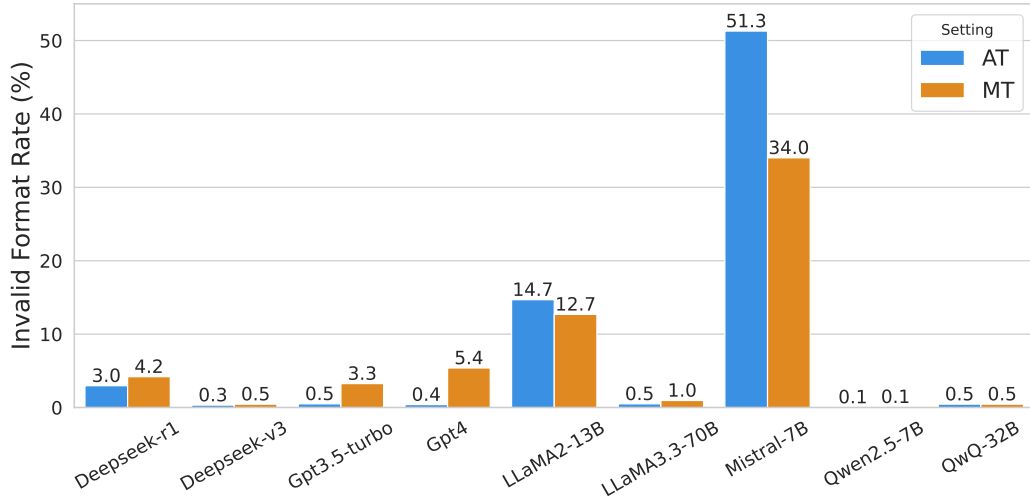


Figure 9: Prompt format violation rates across models in both AT and MT settings.

These findings suggest that instruction adherence is not solely determined by model size or reasoning ability, and remains fragile under ambiguous temporal input.

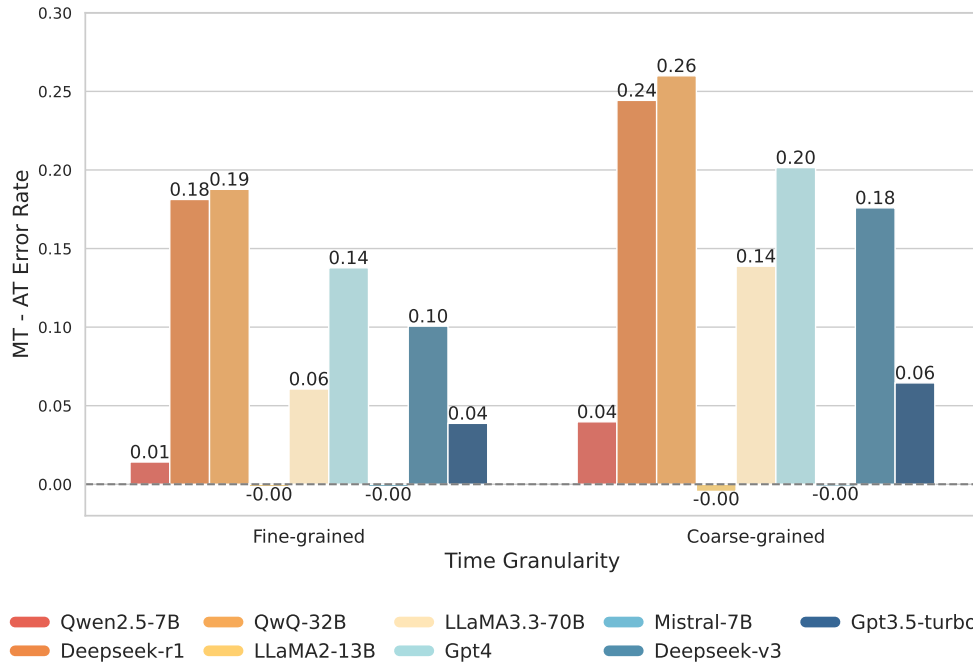


Figure 10: EM error rate increase (MT - AT) across time granularity levels. Coarse-grained (year-only) passages show stronger degradation under mixed-time input, especially for models such as Deepseek-r1 and QwQ-32b.

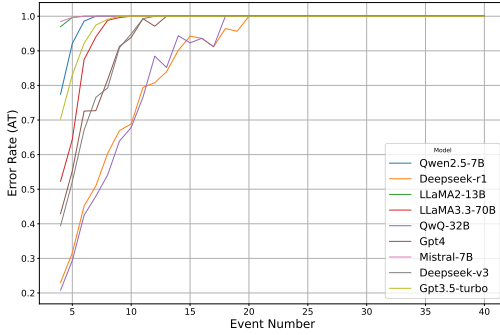
G.2 Granularity Analysis

Figure 10 reveals that when only year-level timestamps are present, models rely heavily on numerical comparison (e.g., 1995 vs. 2000) under AT. Once these cues are replaced with vague relative phrases like “a few years later,” performance degrades sharply. The absence of fine-grained resolution compounds the difficulty of interpreting relative time.

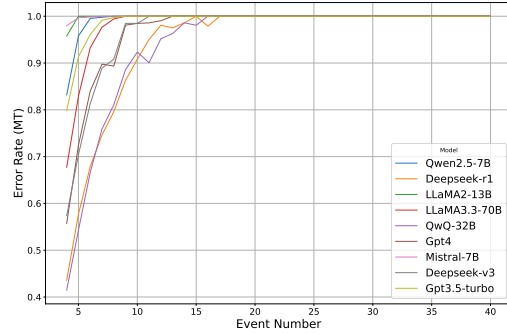
Interestingly, under fine-grained conditions, the performance gap between AT and MT narrows. While absolute timestamps are more complex (e.g., full dates), the corresponding relative phrases (e.g., “early that year,” “a few months earlier”) are often more informative. These naturalistic expressions provide additional linguistic cues that partially compensate for the loss of exact time, helping models maintain ordering accuracy.

H Full Error Rate Curves Across Event Numbers

Figure 11 provide a comprehensive view of model scalability when handling increasing event chains. While the main text focuses on results with up to 15 events (where most meaningful distinctions occur), we include these extended plots to show that beyond this point, most models saturate to an error rate of 1.0, suggesting a consistent upper bound on current models’ capacity for temporal reasoning in complex narratives.



(a) error rate (AT) vs. event number.



(b) error rate (MT) vs. event number.

Figure 11: Error rate trends under AT and MT; most models saturate at 1.0 beyond 15 events, indicating scalability limits