

Too Long, Didn't Model: Decomposing LLM Long-Context Understanding With Novels

Sil Hamilton¹, Rebecca M. M. Hicke², Matthew Wilkens¹, David Mimno¹

¹Department of Information Science

²Department of Computer Science

Cornell University

{srh255,rmh327,wilkens,mimno}@cornell.edu

Abstract

Although the context length of large language models (LLMs) has increased to millions of tokens, evaluating their effectiveness beyond needle-in-a-haystack approaches has proven difficult. We argue that novels provide a case study of subtle, complicated structure and long-range semantic dependencies often over 128k tokens in length. Inspired by work on computational novel analysis, we release the Too Long, Didn't Model (TLDM) benchmark, which tests a model's ability to report plot summary, story-world configuration, and elapsed narrative time. We find that none of seven tested frontier LLMs retain stable understanding beyond 64k tokens. Our results suggest language model developers must look beyond "lost in the middle" benchmarks when evaluating model performance in complex long-context scenarios. To aid in further development we release the TLDM benchmark together with reference code and data.

1 Introduction

Large language model (LLM) context lengths have expanded to millions of tokens, theoretically enabling the analysis of long and complicated documents. However, recent research suggests LLMs do not properly integrate information across long contexts (Hsieh et al., 2024; Karpinska et al., 2024) and have difficulty keeping track of order within contexts (Merrill et al., 2024; Li et al., 2025). This failure mode has been hard to evaluate.

Current long-context benchmarks like Needle In a Haystack (Kamradt, 2023) and Passkey Retrieval (Mohtashami and Jaggi, 2023) evaluate a model's minimal ability to access "lost in the middle" data, but fail to test long-context *understanding*. Retrieving one relevant document from a sea of irrelevant documents does not replicate how we expect most users use long contexts: to integrate multiple documents to arrive at complex conclusions.

We therefore present the Too Long, Didn't Model (TLDM) benchmark: a pipeline for test-

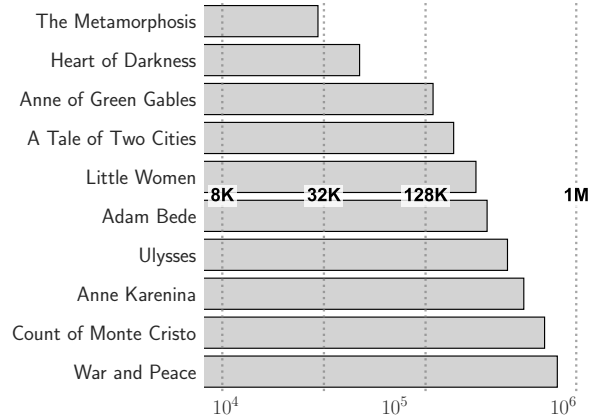


Figure 1: Token lengths of ten popular novels as tokenized by Gemma 3 contrasted with common maximum LLM context lengths (log) as indicated by dotted lines.

ing long-context LLMs on a set of forty English-language novels ranging in length from <32k to >128k tokens using a suite of narrative understanding tasks including summarization, storyworld description, and narrative time estimation. Our benchmark circumvents limitations imposed by the need for human annotations by selecting narrative tasks suitable for associative ground construction. For example, a novel summary can be decomposed to a series of concatenated chapter summaries. This property allows us to assess model stability in long-context regimes by measuring the difference between novel-level and chapter-level predictions.¹

2 Related Work

State sequencing. It is unknown whether self-supervised machine learning models are able to track state over long inputs. While theory suggests they can (Liu et al., 2023; Merrill and Sabharwal, 2025), certain models fail to do so in practice. Merrill et al. (2024) show state-space models (e.g. S4 and Mamba) struggle with state-tracking tasks. Transformers fare better: researchers have found

¹We make our data and code available [here](#).

empirical evidence for state tracking in tasks like entity tracking (Li et al., 2021; Kim and Schuster, 2023), permutation composition (Li et al., 2025), Othello (Li et al., 2023), and chess (Karvonen, 2024). These results indicate transformers can theoretically process long-context input.

Long-context benchmarking. It is popular to test frontier transformers on long-context inputs by subjecting them to benchmarks such as Needle in a Haystack (Kamradt, 2023) and passkey retrieval (Mohtashami and Jaggi, 2023). These benchmarks test long-context processing by having the model retrieve relevant documents randomly shuffled into sets of irrelevant documents. Early long-context models (e.g. GPT-4) performed poorly on these tasks, which encouraged language model developers to forefront long-context testing (Liu et al., 2024; Li et al., 2024; Zhang et al., 2024).

Literature as benchmark. Competent “lost in the middle” performance does not mean models can integrate information over long contexts. Researchers have therefore proposed benchmarks for assessing models on more complicated tasks such as question answering (Wang et al., 2024; Yuan et al., 2024), “multi-hop reasoning” (Roberts et al., 2024), instruction following (Bai et al., 2024), or all of the above (Chen et al., 2025; Hsieh et al., 2024). Benchmarks have likewise begun to turn to literature as a potential source of natural long-context data (Sun et al., 2021; Kim et al., 2024; Ahuja et al., 2025). One such benchmark, NOCHA (Karpinska et al., 2024) finds open-weight models achieve only near-random accuracy when querying texts averaging 127k tokens in length — but NOCHA (and other literary benchmarks) fail to provide insight into *when* in the context window models begin to fail to process information. Our benchmark proposes a unique set of state-oriented narrative understanding tasks to test how models represent state *up to a given point in a novel*.

3 Methods

Data. We align Project Gutenberg’s catalog (Hart, 1971) with the MultiHATHI corpus (Hamilton and Piper, 2023) to identify 6,219 English-only public domain works of fiction averaging 241 pages in length. We then randomly sample books containing clearly delineated chapters that appear to follow a single narrative. We keep the first ten suitable texts that fall into each of four length bins: <32k,

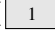
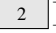
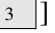



T1: No novel.

[]

T2: Novel, unaltered.

[     ]

T3: Novel, one chapter per user message.

[] [] [] [] [] []

T4: Novel, truncated to window of interest.

[   ]

T5: Novel, chapters shuffled.

[     ]

Figure 2: The five treatments considered in our study as applied to a six-chapter novel. Brackets indicate user message beginning and end while rectangles indicate chapters. Crosshatches indicate absence of input.

32k–64k, 64k–128k, and >128k tokens.²³ Our final dataset contains forty novels containing an average of 27 chapters each.

Tasks. We deploy three narrative understanding tasks that require processing large amounts of text:

1. **Summarization:** Summarize the narrative with one sentence per chapter.
2. **Storyworld description:** Return the last known physical location of every character in the narrative.
3. **Narrative time:** Estimate the narrative time passed in hours, days, months, or years.

Each task requires models to extract and report different types of narrative information; the first requires identifying salient plot points, the second requires entity tracking, and the third requires a model of story time independent of narration.⁴

Windows of interest. Each task-specific prompt instructs the model to perform the task for the first 25%, 50%, 75%, or 100% of chapters. This allows us to evaluate whether models can limit results to a subsection of text. Subsequent results will refer to these subspans as *windows of interest*.

Text treatments. We probe for the circumstances in which LLMs fail to process long contexts by permuting all texts with five treatments, presented in

²³The length of each text in number of tokens is calculated by applying the Gemma 2 tokenizer to the text file versions of each text with paratextual information removed.

³We list all sampled novels in Appendix A.

⁴Task-specific prompts are made available in the appendix.

Figure 2.⁵ Our first treatment (T1) forgoes the actual text for the novel title and author.⁶ Our second treatment (T2) passes in the input text unaltered. The third treatment (T3) wraps each chapter in a unique user message, with the intuition being that explicitly delineating chapters could aid models in parsing long inputs. Our fourth treatment (T4) truncates the input text to the window of interest. Our final treatment (T5) randomly shuffles the chapters to test whether models are able to reconstruct narratives from anachronous input text. We further consider all possible combinations of T3-5 for a total of nine treatments over each input text.

Evaluation. TLDM assesses the stability of model predictions made over long contexts relative to short-context responses. We do not have human-labeled ground truth, but instead compare individual model performance on short contexts to performance on long contexts. Contemporary LLMs are often pretrained with a context length of 4096 tokens before being generalized to longer contexts in post-training (Abdin et al., 2024; Yang et al., 2025; Su et al., 2023). This token range is approximately the length of the average English novel chapter.⁷ We therefore generate chapter-level outputs independently and concatenate them to create model-specific novel-level predictions.⁸ We then prompt the model with each text treated as described in section 3. We finally compute the difference between these full-text predictions and the concatenated short-context output using a similarity heuristic normalized to the range $[0, 1]$.⁹

4 Results

We test seven recent frontier models (Table 1) on the TLDM benchmark to evaluate the current state of the art. All seven models were released in 2025 and support from 128k to 10 million input tokens. We access GPT-4.1 (OpenAI, 2025) and DeepSeek V3 (DeepSeek et al., 2025) via Microsoft Azure, Mistral Small 3.1 (Mistral, 2025) via the Mistral

Dev.	Model	Context	Release	OW
Meta	Llama 4 Scout	10,000,000	4/2025	✓
OpenAI	GPT-4.1	1,000,000	4/2025	×
DeepSeek	DeepSeek V3	1,000,000	2/2025	✓
Google	Gemini 2.0 Flash	1,000,000	2/2025	×
Google	Gemma 3 27b	128,000	4/2025	✓
Alibaba	Qwen 3 32b	128,000	4/2025	✓
Mistral	Mistral Small 3.1	128,000	3/2025	✓

Table 1: Comparison of recent large language models sorted by context window size and release date. OW indicates open weights.

API, and Gemini 2.0 Flash (Team Gemini, 2025) & Gemma 3 27b (Team Gemma et al., 2025) via the Google AI Studio API. We then host Qwen 3 32b (Yang et al., 2025)¹⁰ and Llama 4 Scout (Meta, 2025) on two Nvidia H200 on AWS.¹¹ Values for each length bin are averaged over 10 novels.

Full-novel performance. We first report results where the model is asked to analyze the entire input text (whole unaltered novels, treatment 2 in Table 2). We find that all models exhibit comparable performance when processing volumes with <64k tokens but that performance begins to degrade as book lengths exceed 64k tokens. Performance degrades at different rates, with summary and storyworld scores dropping faster than time estimate scores. Open-weight models (particularly Gemma 3 27b and Qwen 3 32b) exhibit the steepest decline in performance. Of the models equipped to process over 128k tokens, we find GPT-4.1 is most consistent across all context lengths. Llama 4 Scout and Gemini 2.0 Flash are the next most resilient, achieving reliable performance in summary and time estimation over all lengths. However, no model performed well in estimating storyworlds, suggesting models grow increasingly inconsistent in their descriptions when processing individual chapters versus whole chapters. Finally, we find that performance scales with parameter count.

Treatment impact. Examining the effect of each treatment on average model performance reveals several trends. First, increasing novel length decreases model performance across all treatments (excluding title/author only, T1). Even when we only ask models to analyze a subset of the provided text (the “window of interest”), the same pattern

⁵We provide example summary responses in Appendix C.

⁶Note we pass in text author and title for all inputs independent of T1. This condition tests for text memorization.

⁷The mean chapter length in our corpus is 2,845 words or 3,696 Gemma 2 tokens.

⁸We concatenate chapter-level summaries; take the last recorded location of a character across all chapters (recurringly passing in characters from previous chapters to stabilize predictions); and sum per-chapter predictions in seconds.

⁹Semantic similarity for summaries, Jaccard similarity plus semantic similarity for storyworld descriptions, and absolute relative error for time.

¹⁰We disable chain-of-thought token sampling for Qwen 3 32b to maintain even footing with the other models.

¹¹We consume a total of \$600 in compute credits across all services.

Model	Summary					Storyworld					Time				
	<32	32-64	64-128	>128	B	<32	32-64	64-128	>128	B	<32	32-64	64-128	>128	B
GPT-4.1	0.80	0.81	0.81	0.82	0.27	0.17	0.27	0.16	0.09	0.00	0.58	0.54	0.54	0.38	0.35
Llama 4 Scout	0.76	0.77	0.74	0.77	0.29	0.10	0.10	0.07	0.02	0.00	0.55	0.60	0.57	0.58	0.24
Gemini 2.0 Flash	0.73	0.75	0.69	0.72	0.27	0.12	0.14	0.09	0.05	0.00	0.55	0.67	0.54	0.63	0.29
DeepSeek V3	0.75	0.80	0.69	0.30	0.24	0.14	0.12	0.08	0.00	0.00	0.50	0.58	0.61	0.45	0.28
Mistral Small 3.1	0.80	0.73	0.62	-	0.24	0.24	0.19	0.11	-	0.01	0.58	0.47	0.59	-	0.26
Gemma 3 27b	0.72	0.66	0.16	-	0.25	0.18	0.25	0.01	-	0.00	0.62	0.54	0.50	-	0.22
Qwen 3 32b	0.78	0.73	0.62	-	0.26	0.25	0.06	0.04	-	0.00	0.56	0.64	0.56	-	0.32

Table 2: LLM performance comparing full novels (T2) to per-chapter results across different volume lengths, with similarity scores for summaries, storyworlds, and times. Values are normalized to the closed interval $[0, 1]$. All reported lengths are in thousands of Gemma 2 tokens. Performance is consistent across models below the 64k–128k length bracket. Compare with random baselines (B) averaged over 1,440 random pairs on a per-model basis.

holds. In fact, differences in model performance between short and long texts increase as the window of interest increases. That is, there is a greater difference in average model performance for <32k and >128k novels when we request analysis of all chapters than when we ask for 25% of chapters. This pattern makes intuitive sense, as the models are forced to consider more text comparatively when examining longer windows of longer texts.

We see limited evidence of memorization from pre-training. Passing only a volume’s title and author to the models (T1) decreases average model performance by roughly a third on the summary task, and to near zero for storyworld; narrative time estimates are equivalent. Similarly, shuffling the chapters of a text (T5) reduces the performance on the summaries and storyworlds for shorter text windows but does not consistently affect time estimates. Truncating the texts to the window of interest (T4) and all treatment combinations with T4 improve performance for truncated windows of interest; this effect is strongest for the shortest windows and longest texts. In contrast, passing individual chapters as user messages (T3) has little consistent impact on performance. Finally, we note that the average model performance when reporting storyworlds falls as the window of interest increases from 25% of the novel to 100%.

5 Discussion and Conclusion

We present the Too Long, Didn’t Model (TLDM) benchmark for long-context understanding. We release initial benchmark scores for seven frontier LLMs released in early 2025. The benchmark includes three narrative understanding tasks: summarization, storyworld reporting, and narrative time estimation. All require models to infer information

over the full text of a novel. It evaluates models on novels of varying lengths (<32k to >128k), assesses the models’ ability to focus on a particular subset of texts, and determines the impact of various experimental treatments (e.g., shuffling chapters). Assessing models via this benchmark therefore provides a comprehensive understanding of their ability to perform complex long-context analysis.

No tested model is perfect at long contexts. Our preliminary model evaluations show that, despite having context windows of up to 10M tokens, models’ performance declines considerably with longer texts (those above 64k tokens), especially for non-summary tasks. The true context of these models is still limited for complex understanding tasks.

Model scale benefits long-context understanding. We find long-context abilities improve linearly with model size. This indicates smaller, open-weight models that can be run on laptops continue to perform worse than do larger commercial models.

Text linearity aids long-context models. We find long-context models are impacted by document order — especially so when focussing on limited narrative windows. This suggests some models develop particularly inelastic mechanisms for tracking linear narratives.

Next steps. We encourage future researchers to investigate whether mechanistic interpretability could yield answers to questions raised herein. It would be valuable to determine how models are currently performing long-context narrative analysis and tracking narrative state, and whether their representational mechanisms are similar to those of humans. Doing so would help researchers develop strategies for better predicting which long-context tasks LLMs are most appropriate for.

Limitations

There are several key limitations to this work. The first is the lack of true ground truth values. The expense and time needed to produce validated human ground truth for full novel-level annotations means that the TLDM benchmark compares novels only to their own short context performance; thus, while we are able to evaluate how model performance extends comparatively to long contexts, we lack maximally robust assessments of model vs. human capabilities. Second, all texts included in the benchmark are in English, meaning we do not evaluate models’ multilingual performance. Finally, compute restrictions limit the number and variety of models we are able to evaluate in this paper.

Acknowledgments

This work was supported by NEH grant HAA-290374-23, AI for Humanists, granted to Matthew Wilkens and David Mimno.

References

- Marah Abdin, Jyoti Aneja, Harkirat Behl, Sébastien Bubeck, Ronen Eldan, Suriya Gunasekar, Michael Harrison, Russell J. Hewett, Mojan Javaheripi, Piero Kauffmann, James R. Lee, Yin Tat Lee, Yuanzhi Li, Weishung Liu, Caio C. T. Mendes, Anh Nguyen, Eric Price, Gustavo de Rosa, Olli Saarikivi, and 8 others. 2024. [Phi-4 Technical Report](#). *Preprint*, arXiv:2412.08905.
- Kabir Ahuja, Melanie Sclar, and Yulia Tsvetkov. 2025. [Finding Flawed Fictions: Evaluating Complex Reasoning in Language Models via Plot Hole Detection](#). *Preprint*, arXiv:2504.11900.
- Yushi Bai, Xin Lv, Jiajie Zhang, Yuze He, Ji Qi, Lei Hou, Jie Tang, Yuxiao Dong, and Juanzi Li. 2024. [LongAlign: A Recipe for Long Context Alignment of Large Language Models](#). *Preprint*, arXiv:2401.18058.
- Pei Chen, Hongye Jin, Cheng-Che Lee, Rulin Shao, Jingfeng Yang, Mingyu Zhao, Zhaoyu Zhang, Qin Lu, Kaiwen Men, Ning Xie, Huasheng Li, Bing Yin, Han Li, and Lingyun Wang. 2025. LongLeader: A Comprehensive Leaderboard for Large Language Models in Long-context Scenarios. In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 8734–8750, Albuquerque, New Mexico. Association for Computational Linguistics.
- DeepSeek, Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Daya Guo, Dejian Yang, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, and 33 others. 2025. [DeepSeek-V3 Technical Report](#). *Preprint*, arXiv:2412.19437.
- Sil Hamilton and Andrew Piper. 2023. [Multihathi: A complete collection of multilingual prose fiction in the hathitrust digital library](#). *Journal of Open Humanities Data*.
- Michael S. Hart. 1971. [Project Gutenberg](#).
- Cheng-Ping Hsieh, Simeng Sun, Samuel Kriman, Shantanu Acharya, Dima Rekesh, Fei Jia, and Boris Ginsburg. 2024. Ruler: What’s the real context size of your long-context language models? In *First Conference on Language Modeling*.
- Greg Kamradt. 2023. [Needle In A Haystack - Pressure Testing LLMs](#).
- Marzena Karpinska, Katherine Thai, Kyle Lo, Tanya Goyal, and Mohit Iyyer. 2024. [One thousand and one pairs: A “novel” challenge for long-context language models](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 17048–17085, Miami, Florida, USA. Association for Computational Linguistics.
- Adam Karvonen. 2024. Emergent world models and latent variable estimation in chess-playing language models. In *First Conference on Language Modeling*.
- Najoung Kim and Sebastian Schuster. 2023. [Entity tracking in language models](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3835–3855, Toronto, Canada. Association for Computational Linguistics.
- Yekyung Kim, Yapei Chang, Marzena Karpinska, Aparna Garimella, Varun Manjunatha, Kyle Lo, Tanya Goyal, and Mohit Iyyer. 2024. [FABLES: Evaluating faithfulness and content selection in book-length summarization](#).
- Belinda Z. Li, Zifan Carl Guo, and Jacob Andreas. 2025. [\(How\) Do Language Models Track State?](#) *Preprint*, arXiv:2503.02854.
- Belinda Z. Li, Maxwell Nye, and Jacob Andreas. 2021. [Implicit representations of meaning in neural language models](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1813–1827, Online. Association for Computational Linguistics.
- Kenneth Li, Aspen K Hopkins, David Bau, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. 2023. [Emergent world representations: Exploring a sequence model trained on a synthetic task](#). In *The Eleventh International Conference on Learning Representations*.

- Tianle Li, Ge Zhang, Quy Duc Do, Xiang Yue, and Wenhui Chen. 2024. [Long-context LLMs Struggle with Long In-context Learning](#).
- Bingbin Liu, Jordan T. Ash, Surbhi Goel, Akshay Krishnamurthy, and Cyril Zhang. 2023. [Transformers Learn Shortcuts to Automata](#). *Preprint*, arXiv:2210.10749.
- Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2024. [Lost in the Middle: How Language Models Use Long Contexts](#). *Transactions of the Association for Computational Linguistics*, 12:157–173.
- William Merrill, Jackson Petty, and Ashish Sabharwal. 2024. The Illusion of State in State-Space Models. In *Proceedings of the 41st International Conference on Machine Learning*, volume 235. PMLR.
- William Merrill and Ashish Sabharwal. 2025. [A Little Depth Goes a Long Way: The Expressive Power of Log-Depth Transformers](#). *Preprint*, arXiv:2503.03961.
- Meta. 2025. The Llama 4 Herd: The beginning of a new era of natively multimodal AI innovation. <https://ai.meta.com/blog/llama-4-multimodal-intelligence/>.
- Mistral. 2025. Mistral Small 3.1. <https://mistral.ai/news/mistral-small-3-1>.
- Amirkeivan Mohtashami and Martin Jaggi. 2023. [Random-access infinite context length for transformers](#). In *Thirty-seventh Conference on Neural Information Processing Systems*.
- OpenAI. 2025. Introducing GPT-4.1 in the API. <https://openai.com/index/gpt-4-1/>.
- Jonathan Roberts, Kai Han, and Samuel Albanie. 2024. Needle threading: Can llms follow threads through near-million-scale haystacks? *arXiv preprint arXiv:2411.05000*.
- Jianlin Su, Yu Lu, Shengfeng Pan, Ahmed Murtadha, Bo Wen, and Yunfeng Liu. 2023. [RoFormer: Enhanced Transformer with Rotary Position Embedding](#). *Preprint*, arXiv:2104.09864.
- Simeng Sun, Kalpesh Krishna, Andrew Mattarella-Micke, and Mohit Iyyer. 2021. [Do Long-Range Language Models Actually Use Long-Range Context?](#) *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 807–822.
- Team Gemini. 2025. Gemini 2.0: Flash, Flash-Lite and Pro. <https://developers.googleblog.com/en/gemini-2-family-expands/>.
- Team Gemma, Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane Rivi re, Louis Rouillard, Thomas Mesnard, Geoffrey Cideron, Jean-bastien Grill, Sabela Ramos, Edouard Yvinec, Michelle Casbon, Etienne Pot, Ivo Penchev, and 14 others. 2025. [Gemma 3 Technical Report](#). *Preprint*, arXiv:2503.19786.
- Minzheng Wang, Longze Chen, Cheng Fu, Shengyi Liao, Xinghua Zhang, Bingli Wu, Haiyang Yu, Nan Xu, Lei Zhang, Run Luo, and 1 others. 2024. Leave no document behind: Benchmarking long-context llms with extended multi-doc qa. *arXiv preprint arXiv:2406.17419*.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, and 41 others. 2025. [Qwen3 Technical Report](#). *Preprint*, arXiv:2505.09388.
- Tao Yuan, Xuefei Ning, Dong Zhou, Zhijie Yang, Shiyao Li, Minghui Zhuang, Zheyue Tan, Zhuayu Yao, Dahua Lin, Boxun Li, Guohao Dai, Shengen Yan, and Yu Wang. 2024. [LV-Eval: A Balanced Long-Context Benchmark with 5 Length Levels Up to 256K](#). *Preprint*, arXiv:2402.05136.
- Xinrong Zhang, Yingfa Chen, Shengding Hu, Zihang Xu, Junhao Chen, Moo Khai Hao, Xu Han, Zhen Leng Thai, Shuo Wang, Zhiyuan Liu, and Maosong Sun. 2024. [\\$infty\\$Bench: Extending Long Context Evaluation Beyond 100K Tokens](#).

A Corpus Contents

The books included in the TLDM benchmark are:

<32k

- *Beasley’s Christmas Party* by Booth Tarkington
- *The Battle Of The Strong (A Romance of Two Kingdoms): Volume 2* by Gilbert Parker
- *The Caxtons: Part 12* by Edward Bulwer-Lytton
- *Godolphin: Volume 5* by Edward Bulwer-Lytton
- *The Romance of a Christmas Card* by Kate Douglas Wiggin
- *The Story of a China Cat* by Laura Lee Hope
- *Better Dead* by J. M. Barrie
- *The Tale of Jasper Jay Tuck-Me-In Tales* by Arthur Scott Bailey
- *Our Little Hawaiian Cousin* by Mary Hazelton Wade
- *Christmas at Thompson Hall* by Anthony Trollope

32k–64k

- *Alexander’s Bridge and The Barrel Organ* by Willa Cather and Alfred Noyes
- *Tom Swift and His Undersea Search or The Treasure on the Floor of the Atlantic* by

Victor Appleton

- *Kilmeny of the Orchard* by Lucy Maud Montgomery
- *What Will He Do With It: Book 10* by Edward Bulwer-Lytton
- *The Tragedy of the Korosko* by Arthur Conan Doyle
- *Dorothy Dainty's Gay Times* by Amy Brooks
- *Thistle and Rose: A Story for Girls* by Amy Walton
- *Ruth Fielding Homeward Bound: A Red Cross Worker's Ocean Perils* by Alice B. Emerson
- *Isla Heron* by Laura E. Richards
- *Frank Reade, Jr., Fighting the Terror of the Coast* by Anonymous

64k–128k

- *Lost in the Fog* by James De Mille
- *Dora Deane; Or, The East India Uncle* by Mary Jane Holmes
- *Going Some* by Rex Beach
- *A Pirate of Parts* by Richard Neville
- *The Backwoodsmen* by Charles G. D. Roberts
- *The Watchers: A Novel* by A. E. W. Mason
- *Discourses of Keidansky* by Bernard G. Richards
- *In Queer Street* by Fergus Hume
- *Dick Merriwell's Assurance; Or, In His Brother's Footsteps* by Burt L. Standish
- *The Earl's Promise, A Novel: Volume 2* by Mrs. J. H. Riddell

128k+

- *The Moon Pool* by Abraham Merritt
- *Under Two Flags* by Ouida
- *Born in Exile* by George Gissing
- *Esther Waters* by George Moore
- *Desert Conquest; or, Precious Waters* by A. M. Chisholm
- *The Dust Flower* by Basil King
- *Wager of Battle: A Tale of Saxon Slavery in Sherwood Forest* by Henry William Herbert
- *Betty Alden: The first-born daughter of the Pilgrims* by Jane G. Austin
- *Perch of the Devil* by Gertrude Atherton
- *The Brooklyn Murders* by G. D. H. Cole

B Example Prompts (T1)

What follows are three example prompts for the novel-level prediction of the 1898 novel "The Battle of the Strong: A Romance of Two Kingdoms"

by Gilbert Parker. Each prompt is run with a new user session.

Summary.

Source: "The Battle of the Strong: A Romance of Two Kingdoms" by Gilbert Parker.

Situation: You were given a narrative. You will now be given a task about the narrative. Complete the task. Keep your response brief and to the point.

Task: Summarize the narrative with one sentence per chapter. Describe what happens. Do not reference the narrative itself.

Limit your response to the narrative from chapter 1 up until, and including, chapter 12.

Storyworld description.

Source: "The Battle of the Strong: A Romance of Two Kingdoms" by Gilbert Parker.

Situation: You were given a narrative. You will now be given a task about the narrative. Complete the task. Keep your response brief and to the point.

Task: List each character in the narrative and their physical location in the story.

Here are a list of possible characters in the narrative: {characters}. The list might be blank.

If the character in the above list is present in the narrative, note their last location in the narrative. If a character is in the narrative but is not in the above list, still note their last location in the narrative. Be consistent with the name. Only list characters present in the narrative. Only list individuals. Do not list groups of characters.

Limit your response to the narrative from chapter 1 up until, and including, chapter 12.

JSON Schema follows.

Narrative time estimation.

Source: "The Battle of the Strong: A Romance of Two Kingdoms" by Gilbert Parker.

Situation: You were given a narrative. You will now be given a task about the narrative. Complete the task. Keep your response brief and to the point.

Task: Predict how much time in minutes, hours, or days elapsed in this narrative. Specify an exact number and unit. Only respond with a number and a unit, e.g. 30 minutes, 6 hours, or 1 day.

Limit your response to the narrative from chapter 1 up until, and including, chapter 12.

C Example GPT-4.1 responses.

Here we provide example summaries for *Christmas at Thompson Hall* by Anthony Trollope. These summaries were generated with GPT-4.1. We specifically provide the (concatenated) ground, and the summaries produced under treatments 2 (the whole unaltered novel), 3 (each chapter passed in as a separate user message), and 5 (chapters randomly shuffled) — with a full window of interest.

Ground. 1. During a harsh winter journey to attend a family Christmas gathering, an Englishwoman in a Paris hotel braves confusing corridors and an awkward encounter with staff to retrieve mustard as a remedy for her ailing husband's throat.

2. Mrs. Brown, after preparing a mustard plaster for her supposedly ill husband in a Paris hotel, accidentally applies it to the throat of a sleeping stranger in the wrong room, panics at her mistake, considers—but ultimately avoids—confessing to the night porter, and returns distraught and chilled to her own room.

3. Mrs. Brown, desperate to escape a French hotel after mistakenly putting a mustard plaster on a stranger instead of her ailing husband Mr. Brown, anxiously concocts excuses to persuade her reluctant husband to leave early while fearing discovery by the night porter and the injured stranger.

4. Mrs. Brown, after mistakenly entering a stranger's hotel room at night and inadvertently placing a mustard plaster on his throat instead of her husband's, faces public embarrassment and suspicion but ultimately explains the misunderstanding, allowing her and Mr. Brown to narrowly avoid

greater trouble and continue their journey, despite lingering awkwardness with the aggrieved stranger and staff.

5. A series of misunderstandings and mistaken identity occur when Mrs. Brown accidentally applies a mustard plaster to a stranger, Mr. Jones, at a hotel, leading to awkwardness and eventual reconciliation when it is revealed that Mr. Jones is engaged to Mrs. Brown's sister and all parties meet at Thompson Hall.

Treatment 2 (whole input). 1. During a bitterly cold Christmas journey to England, Mrs. Brown struggles with her husband's alleged illness and obtains mustard from their Paris hotel to treat his throat.

2. In the dark, Mrs. Brown mistakenly applies the mustard poultice to a stranger, Mr. Jones, instead of her husband, and flees in panic rather than wake or help him.

3. Distraught and desperate to leave the hotel to avoid a confrontation, Mrs. Brown invents excuses to hurry her husband's departure despite guilt and fear of discovery.

4. Mrs. Brown is confronted about the incident when the porter discovers her handkerchief in Mr. Jones's room, leading to an embarrassing explanation and apology to Mr. Jones, which allows them to continue their journey.

5. Upon arrival at Thompson Hall, Mrs. Brown and her husband awkwardly discover that the stranger from Paris, Mr. Jones, is her sister's fiancé, leading to the revelation of her mistake and eventual forgiveness during the family's Christmas celebration.

Treatment 3 (separate user messages). 1. Anxious to reach her family for Christmas despite her husband's complaints and illness, Mrs. Brown brings him to a cold Paris hotel, where he demands a mustard plaster for his throat.

2. In the middle of the night, Mrs. Brown becomes lost in the hotel while seeking mustard, then mistakenly applies the plaster to a stranger's throat instead of her husband's.

3. Discovering her error and unable to face the consequences, Mrs. Brown desperately tries to avoid detection and persuade her husband that they must leave quickly for England.

4. The Browns are confronted by hotel staff and the porter with evidence of the mishap, forcing Mrs. Brown to confess, after which her husband

explains the accident to the injured stranger, Mr. Jones, allowing them to depart.

5. Unexpectedly reunited at Thompson Hall, it is revealed that Mr. Jones is the fiancé of Mrs. Brown's sister, Jane, and after an awkward Christmas gathering, Mrs. Brown's mistake becomes a family anecdote, ultimately forgiven by all.

Treatment 5 (shuffled chapters). Certainly! Here is a one-sentence summary for each of chapters I to V:

Chapter I: On a cold Christmas in Paris, Mrs. Brown ventures through a labyrinthine hotel late at night to obtain mustard for her ailing husband's sore throat.

Chapter II: After a series of mishaps, Mrs. Brown mistakenly enters the wrong room and, believing the sleeping man is her husband, applies a mustard plaster to a stranger's throat.

Chapter III: Realizing her error and fearing exposure, Mrs. Brown slips away from the stranger, lies to her husband about her absence, and anxiously resolves they must escape early the next morning.

Chapter IV: As the Browns try to leave, Mrs. Brown is confronted by hotel staff who have evidence of her midnight blunder, leading to an awkward explanation and a necessary apology to the afflicted Mr. Jones before their delayed departure.

Chapter V: Traveling together by coincidence, the Browns and Mr. Jones arrive at Thompson Hall, where, amidst family gatherings and mounting suspicion, the truth of Mrs. Brown's comical mistake is ultimately revealed and reluctantly forgiven.

D Treatment Results

We provide all treated experiment results in [Figure 3](#).

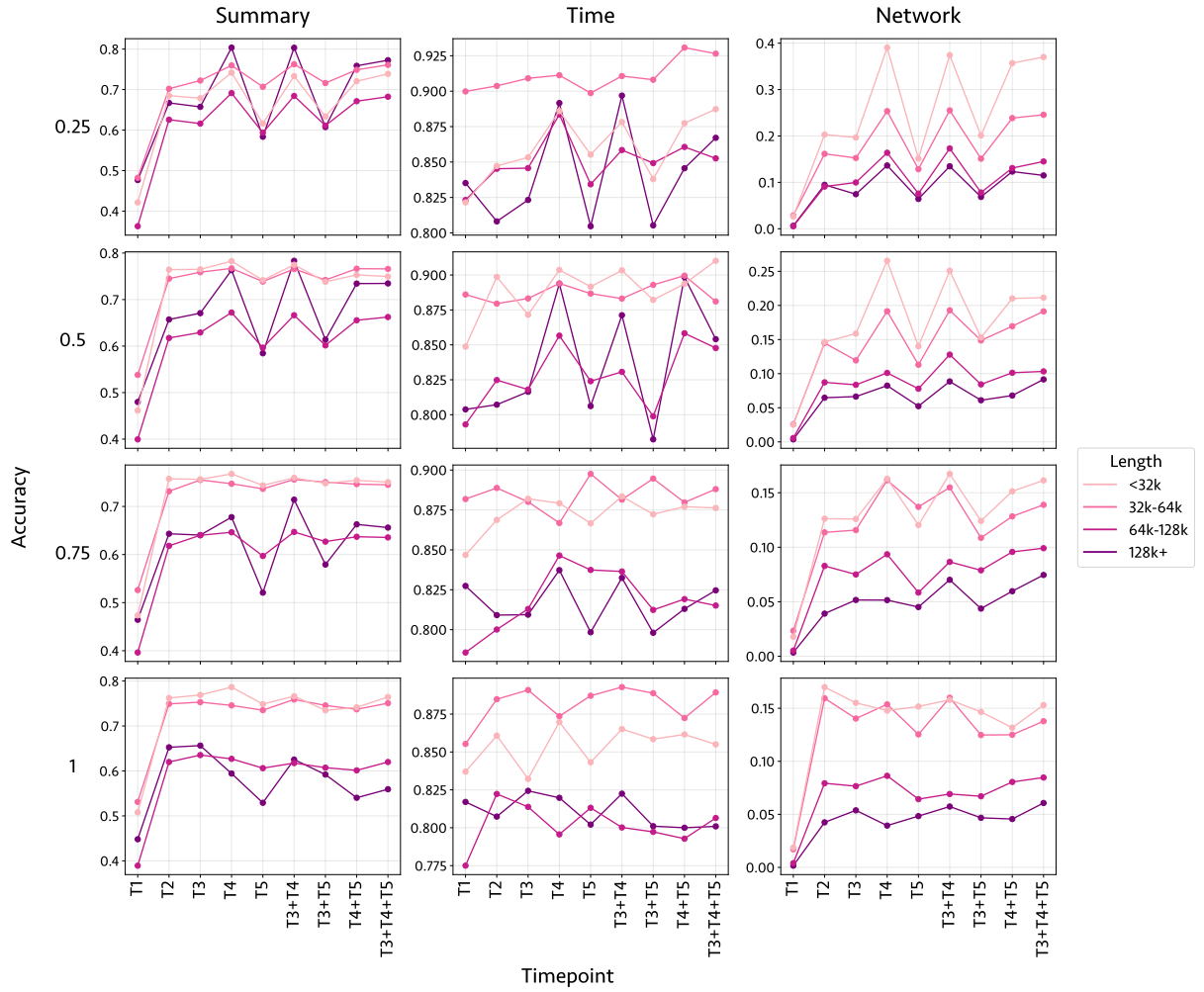


Figure 3: Accuracy scores per task averaged over all models. Values for each window of interest are presented. Note performance consistently degrades as input length increases.