Distance between Relevant Information Pieces Causes Bias in Long-Context LLMs

¹Tsinghua University ²ModelBest Inc. ³Renmin University of China ⁴Apple Inc.

trc20@mails.tsinghua.edu.cn, goody1027@gmail.com

Abstract

Positional bias in large language models (LLMs) hinders their ability to effectively process long inputs. A prominent example is the "lost in the middle" phenomenon, where LLMs struggle to utilize relevant information situated in the middle of the input. While prior research primarily focuses on single pieces of relevant information, real-world applications often involve multiple relevant information pieces. To bridge this gap, we present LONGPIBENCH, a benchmark designed to assess positional bias involving multiple pieces of relevant information. Thorough experiments are conducted with five commercial and six open-source models. These experiments reveal that while most current models are robust against the "lost in the middle" issue, there exist significant biases related to the spacing of relevant information pieces. These findings highlight the importance of evaluating and reducing positional biases to advance LLM's capabilities¹.

1 Introduction

Large language models (LLMs) (Zhao et al., 2023; Minaee et al., 2024) have made significant progress in various natural language processing tasks (Hendrycks et al., 2021; Han et al., 2021). In particular, applications such as code repository analysis (Chen et al., 2021) and information extraction (Kočiský et al., 2018) often require processing long texts, with context lengths reaching up to 200,000 tokens (Li et al., 2024; Zhang et al., 2024). To address these demands, researchers have focused on enhancing LLMs' ability to handle extended inputs effectively (Chen et al., 2023; Han et al., 2024). As a result, multiple LLMs have been developed (Dubey et al., 2024; Team et al., 2024; OpenAI, 2024) which support context lengths of up to one million tokens.

Recent studies have shown that the position of relevant information significantly affects the performance of long-context LLMs (Liu et al., 2023; Lei et al., 2024; Hsieh et al., 2024). In "needle in a haystack" tasks, models struggle to utilize information located in the

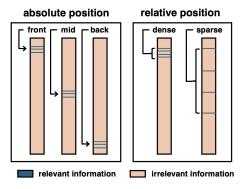


Figure 1: Illustration of absolute position and relative position. Absolute position refers to the location of relevant information within the entire context sequence, while relative position represents the distribution and distance between multiple relevant information pieces.

middle of the input, a phenomenon known as the "lost in the middle" effect (Liu et al., 2023). This evaluation method is commonly used to analyze positional bias (Hengle et al., 2024; Nelson et al., 2024). These analyses (Liu et al., 2023) focused on single relevant information pieces and their positions in the input sequence (front, middle, back), referred to as **absolute positions**.

However, real-world tasks like data analysis (Zhang et al., 2024) often involve multiple pieces of relevant information. This introduces a new characteristic: the distance between relevant information pieces, or how densely they are distributed, termed as **relative position**. Evidence from two types of extreme cases indicates that varying relative position may lead to significant bias, impairing LLMs' long-context performance (Lei et al., 2024). However, this kind of biases have not been systematically studied so far, which highlights the need for thorough investigation.

To bridge the gap, we introduce LONGPIBENCH, a benchmark designed to evaluate positional bias with multiple relevant pieces. It assesses positional bias in two categories: (1) **absolute positions**, referring to the location of relevant information within the entire context, and (2) **relative positions**, referring to the distribution and distance between multiple relevant information pieces. It includes diverse tasks of different complexity and spans four input lengths from 32K to 256K tokens. To the best of our knowledge, LONGPIBENCH is the most comprehensive benchmark for isolating and ana-

^{*} indicates equal contribution.

[†] Corresponding Author.

¹ Code and data are available.

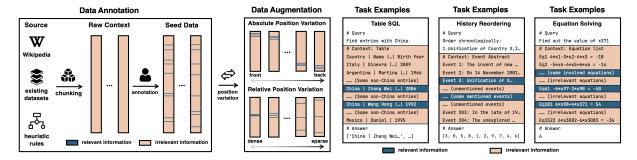


Figure 2: Construction and task examples of LONGPIBENCH. We manually annotated seed data and varied the positions of relevant information for data augmentation.

lyzing positional bias in long text models.

We evaluated eleven popular LLMs. Our experimental analysis yields several key findings: (1) Unlike earlier models significantly affected by the "lost in the middle" issue, most current models demonstrate enhanced robustness against this phenomenon. (2) However, they exhibit significant biases related to the spacing of relevant information. Specifically, as the distance between relevant pieces increases, model performance declines sharply before stabilizing. (3) Additionally, we discuss the impact of model size and query-aware contextualization on this issue.

These findings emphasize the importance of evaluating and mitigating positional biases to advance long-context LLM capabilities.

2 LONGPIBENCH

LONGPIBENCH is a dataset designed to evaluate positional bias with multiple relevant information pieces. As shown in Figure 2, we first manually annotated several seed examples and then augmented them by varying the positions of relevant information. More details can be found in Appendix A.

2.1 Core Statistics

LONGPIBENCH contains 3 different tasks, 4 different input length levels²: (32k, 64k, 128k, and 256k). To comprehensively analyze the impact of positional bias, we set 16 different absolute and relative location levels respectively. The benchmark is composed of 7,680 instances, each containing 10 pieces of relevant information. The whole dataset comprises to 922M tokens.

2.2 Seed Data Annotation

We manually labeled 20 seed data points for three tasks: *Table SQL*, *Timeline Reordering*, and *Equation Solving*, which cover a range of complexities in long-context tasks. Each instance contains 10 relevant pieces of information. This selection was based on an examination of long-context application scenarios, where the number of relevant elements typically falls around the order of magnitude of ten, although it varies across different tasks (Bai et al., 2023; Wang et al., 2024; Dong et al.,

2024). Detailed task definitions, examples, and other pertinent details are provided in Appendix A.

2.3 Data Augmentation

To analyze the positions of relevant information, we augmented the data by altering the absolute and relative positions of the relevant pieces while keeping all other features unchanged.

We broke down the context into elements based on natural information units: table entries for *Table SQL*, event entries for *Timeline Reordering*, and equation lines for *Equation Solving*. We labeled each element as relevant or irrelevant by selecting relevant elements, forming queries around them, and adding irrelevant ones. By introducing varying amounts of irrelevant information, he context lengths are controlled at four levels: 32K, 64K, 128K, and 256K. We then shuffled the element positions to introduce positional variations.

Absolute Position To analyze the impact of absolute position on LLM performance, we manipulated where relevant information appears in the context. Each context was divided into 16 equal segments from start to end. We placed all 10 relevant pieces within a single segment to keep their relative positions consistent. By moving this segment from the first to the last position, we varied the absolute position from the start to the end of the input. The average position of these relevant pieces served as the absolute position metric which is calculated as:

$$\text{Average Location} = \left(\frac{l-1}{N-1}\right) \times L,$$

where l is the current level, N is the total number of levels (16), and L is the length of the context.

This setup allowed us to assess how model performance changes as relevant information is placed further back in the context.

Relative Position To examine the effect of spacing between relevant information pieces on LLM performance, we created 16 levels of distribution density. Each level represents a different spacing configuration among the 10 relevant pieces. At the densest level, all relevant pieces are adjacent with no irrelevant information between them. At the sparsest level, they are evenly distributed throughout the context with equal intervals of

²measured with GPT2Tokenizer (Radford et al., 2019)

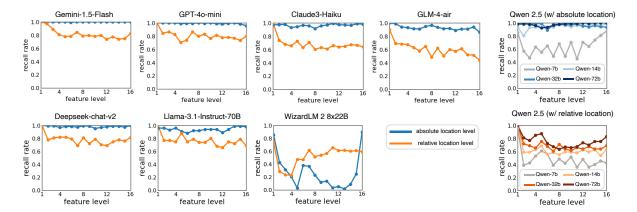


Figure 3: The impact of relevant information's absolute and relative position on eleven popular long context LLMs. A higher absolute position level indicates locations closer to the end of input, while a higher relative position level indicates a greater distance between relevant pieces of information.

irrelevant information. Intermediate levels gradually increase spacing from adjacent to evenly spaced. The distance between each relevant piece is calculated as:

$$\text{Distance} = \left(\frac{L}{n-1}\right) \times \left(\frac{l-1}{N-1}\right),$$

where n is the number of relevant pieces (10), l is the current level ranging from 1 to N, N is the total number of levels (16), and L is the length of the context.

To control for absolute position effects, we randomized the starting position of the first relevant piece in each example. This ensures that any observed performance differences are due to relative spacing rather than absolute positions within the context.

2.4 Quality Control

To ensure the integrity and reliability of the benchmark, we implemented comprehensive quality control measures throughout its construction.

Manual Correction Automatically synthesized data generated by rules or language models (LLMs) can be noisy and may not accurately reflect real-world scenarios. To address this, we manually annotated queries and answers based on the provided context, ensuring that the data is both accurate and contextually relevant.

Knowledge Masking In knowledge-intensive tasks like the Timeline Reordering task, there's a risk of data leakage if models rely on pre-existing knowledge rather than the provided context. To prevent this, we implemented knowledge masking by anonymizing or fictionalizing events using GPT-4 (OpenAI et al., 2024). This approach ensures that models cannot leverage memorized information and must rely solely on the context presented to them.

3 Experimental Setup

To evaluate the influence of context information positioning on long-text large language models (LLMs), we conducted experiments using popular long-context language models on a subset of LONGPIBENCH.

3.1 Models

We assessed a total of eleven LLMs, comprising six open-source and five commercial options. The selection of open-source models includes the 70B model from Llama-3.1-Instruct series (Dubey et al., 2024), the 7B, 14B, 32B, 72B models from Qwen-2.5 family (Qwen, 2024), the 8×22B model of WizardLM-2 (Xu et al., 2023). The commercial models we selected are GPT-4o-mini (OpenAI, 2024), Claude-3-Haiku (Anthropic, 2024), Gemini-1.5-Flash (Team et al., 2024), GLM-4-air (GLM et al., 2024), and Deepseek-Chatv2 (DeepSeek-AI et al., 2024). The selected models are good representatives of popular and top-performance long-context models.

3.2 Configuration

During the pre-test stage, we found that *Timeline Re-ordering* and *Equation Solving* proved to be excessively challenging, with even the top-performing LLMs achieving accuracy rates between 0.0 and 0.2 (more in Appendix B). Therefore, we reserve these two tasks for more advanced LLM analysis in the future. We currently focus our experiments on the *Table SQL* task. Additionally, since 32k tokens is the minimal context length supported by all tested LLMs, we standardized the context length to 32k³ tokens for all experiments. This uniform configuration ensures consistency across different models, allowing for a fair comparison of their performance under identical conditions. Detailed discussions on statistics, evaluation metrics, parameter settings, and prompt configurations are provided in Appendix C.

4 Results and Discussion

In this section, we analyze the impact of absolute and relative positional bias. And we further analyze these phenomena from two perspectives: the number of parameters and query-aware contextualization. Detailed Experimental Results are available in Appendix D.

³The minimal context size is 64k, but some tokenizers expand our 64k inputs to nearly 80k, exceeding the limit.

4.1 Impact of Absolute Position

As illustrated by the **blue lines** in Figure 3, we progressively shift the interval of relevant information from the beginning to the end and observed that while a few open-source models like Qwen 2.5 (7B) (Qwen, 2024) and WizardLM 2 ($8 \times 22B$) (Xu et al., 2023) still suffer from the severe "lost in the middle" phenomenon, commercial models and larger open-source models do not exhibit effects related to absolute position.

This outcome significantly surpasses previous evaluations (Liu et al., 2023), indicating that current long-context models have achieved greater robustness against variations in absolute position of relevant information.

4.2 Impact of Relative Position

As illustrated by the **orange lines** in Figure 3, we progressively increase the distance between relevant pieces of information and observe that all open-source and commercial models exhibit a significant bias toward different relative positions. This bias is characterized by an initial rapid decline in performance followed by a more gradual decrease. Even in straightforward retrieval tasks, relative position bias can lead to a 20–30% reduction in recall rates for competent commercial models.

These findings indicate that the relative positioning among multiple relevant pieces of information is a serious and unresolved issue, which may substantially undermine the effectiveness of long-text language models in practical applications.

4.3 Further Analysis

Effect of Parameter Size When selecting models for evaluation, we included four variants from the Qwen 2.5 Family (Qwen, 2024) with differing parameter sizes. These models exhibit no significant differences in architecture, training methods, or training data. By analyzing their performance under identical positional information features, we can isolate the impact of parameter size on the robustness to positional bias.

As illustrated in Figure 3, for absolute position bias, we found that simply increasing the model parameters from 7B to 14B—while keeping architecture, training methods, and data constant substantially mitigates the "lost in the middle" (Liu et al., 2023) issue. This suggests that robustness to absolute positions may be an "emergent ability" (Wei et al., 2022) and increasing the number of parameters can significantly enhances it.

In contrast, regarding biases related to relative positional information, augmenting the number of parameters only yielded minor quantitative improvements and did not alter the pronounced bias trend. This trend remains largely unchanged even in commercial models with approximately hundreds of billions of parameters. These findings indicate that merely increasing parameter size is insufficient to develop robustness to relative positions, and new techniques may be necessary.

Effect of Query-Aware Contextualization Liu et al. (2023) demonstrated that the placement of the query

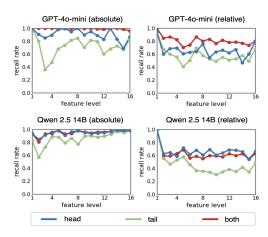


Figure 4: Impact of query placement (beginning, end, both) on the performance of GPT-4o-mini (OpenAI, 2024) and Owen-2.5-14B (Owen, 2024) models.

(beginning or end of the context) significantly affects the performance of decoder-only models due to unidirectional attention. When the query is placed after the context, the LLM cannot attend to the query token while processing the context tokens.

As shown in Figure 4, our experiments on GPT-4o-mini (OpenAI, 2024) and Qwen-2.5-14B (Qwen, 2024) corroborate this observation and confirm that it also holds for bias caused by relative position changes. Specifically, when the query is positioned at the end of the context, the model's performance is significantly worse compared to scenarios where the query is placed at the beginning or both at the beginning and the end. However, the difference between having the query solely at the beginning versus having it at both the beginning and the end varies depending on the model. This indicates that for decoder-only long-text models, positioning the query before the context is of paramount importance.

5 Conclusion

This study investigates a new category of positional bias involving multiple relevant pieces of information in long-context LLMs through three key contributions.

- (1) **Benchmark Development**: We introduce LONG-PIBENCH, the most comprehensive benchmark for evaluating positional bias in long-text LLMs, assessing both absolute and relative biases.
- (2) Comprehensive Evaluation: Using LONG-PIBENCH, we evaluated eleven popular LLMs, investigated the "lost in the middle" phenomenon, and identified novel yet significant biases related to the relative positioning of multiple relevant pieces of information.
- (3) Insightful Findings: Our experiments show that while modern LLMs have improved robustness against absolute positional biases, they are highly sensitive to the distance between relevant pieces of information. Performance declines sharply as the distance increases before stabilizing. We also explore how model size and query-aware contextualization impact these biases.

These findings emphasize the necessity of continuously mitigating positional biases in long-text models.

Limitation

Lack of In-depth Analysis Our systematic experiments demonstrate that two types of positional bias exist when multiple related pieces of information are present in the context. We also analyzed how these biases relate to the number of parameters and query contextualization. However, we are currently unable to explain the reasons behind these two positional biases.

Focus on Specific Models The evaluation was conducted on a set of eleven popular large language models (LLMs), including both open-source and commercial options. However, the findings are limited to these models. The study does not account for the performance of other emerging or less popular models, which might exhibit different results regarding positional biases.

Ethical Considerations

Human Annotation Our seed construction process involves manual annotation. This annotation was carried out by some of the authors, who are researchers with substantial knowledge in LLM evaluation. Consent was obtained from the individuals whose data we are using or curating. The data collection protocol was approved.

Data Security Some data in our Table SQL task may appear to pertain to personal information. However, this data is not actual personal information. Instead, it is generated by us through specific heuristics, eliminating the risk of personal information leakage.

Use of AI assistants We use GPT-40 (OpenAI, 2024) for expression modification and grammar sanity check during the composition process.

References

Anthropic. 2024. Introducing claude 3.5 sonnet. Accessed: 2024-09-15.

Yushi Bai, Xin Lv, Jiajie Zhang, Hongchang Lyu, Jiankai Tang, Zhidian Huang, Zhengxiao Du, Xiao Liu, Aohan Zeng, Lei Hou, Yuxiao Dong, Jie Tang, and Juanzi Li. 2023. Longbench: A bilingual, multitask benchmark for long context understanding. arXiv preprint arXiv:2308.14508.

Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. 2021. Evaluating large language models trained on code. *Preprint*, arXiv:2107.03374.

Yukang Chen, Shengju Qian, Haotian Tang, Xin Lai, Zhijian Liu, Song Han, and Jiaya Jia. 2023. Longlora: Efficient fine-tuning of long-context large language models. *arXiv preprint arXiv:2309.12307*.

DeepSeek-AI, Aixin Liu, Bei Feng, Bin Wang, Bingxuan Wang, Bo Liu, Chenggang Zhao, Chengqi Dengr, Chong Ruan, Damai Dai, et al. 2024. Deepseek-v2: A strong, economical, and efficient mixture-of-experts language model. *Preprint*, arXiv:2405.04434.

Zican Dong, Tianyi Tang, Junyi Li, Wayne Xin Zhao, and Ji-Rong Wen. 2024. BAMBOO: A comprehensive benchmark for evaluating long text modeling capacities of large language models. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 2086–2099, Torino, Italia. ELRA and ICCL.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *Preprint*, arXiv:2407.21783.

Team GLM, Aohan Zeng, Bin Xu, Bowen Wang, Chenhui Zhang, Da Yin, Diego Rojas, Guanyu Feng, Hanlin Zhao, Hanyu Lai, et al. 2024. Chatglm: A family of large language models from glm-130b to glm-4 all tools. *Preprint*, arXiv:2406.12793.

Chi Han, Qifan Wang, Hao Peng, Wenhan Xiong, Yu Chen, Heng Ji, and Sinong Wang. 2024. LM-infinite: Zero-shot extreme length generalization for large language models. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 3991–4008, Mexico City, Mexico. Association for Computational Linguistics.

Xu Han, Zhengyan Zhang, Ning Ding, Yuxian Gu, Xiao Liu, Yuqi Huo, Jiezhong Qiu, Yuan Yao, Ao Zhang, Liang Zhang, et al. 2021. Pre-trained models: Past, present and future. *Preprint*, arXiv:2106.07139.

Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. Measuring massive multitask language understanding. *Preprint*, arXiv:2009.03300.

Amey Hengle, Prasoon Bajpai, Soham Dan, and Tanmoy Chakraborty. 2024. Multilingual needle in a haystack: Investigating long-context behavior of multilingual large language models. *Preprint*, arXiv:2408.10151.

Cheng-Yu Hsieh, Yung-Sung Chuang, Chun-Liang Li, Zifeng Wang, Long T. Le, Abhishek Kumar, James Glass, Alexander Ratner, Chen-Yu Lee, Ranjay Krishna, et al. 2024. Found in the middle: Calibrating positional attention bias improves long context utilization. *Preprint*, arXiv:2406.16008.

Tomáš Kočiský, Jonathan Schwarz, Phil Blunsom, Chris Dyer, Karl Moritz Hermann, Gábor Melis, and Edward Grefenstette. 2018. The NarrativeQA reading comprehension challenge. *Transactions of the Association for Computational Linguistics*, 6:317–328.

Fangyu Lei, Qian Liu, Yiming Huang, Shizhu He, Jun Zhao, and Kang Liu. 2024. S3eval: A synthetic, scalable, systematic evaluation suite for large language models. *Preprint*, arXiv:2310.15147.

- Jiaqi Li, Mengmeng Wang, Zilong Zheng, and Muhan Zhang. 2024. LooGLE: Can long-context language models understand long contexts? In *Proceedings* of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 16304–16333, Bangkok, Thailand. Association for Computational Linguistics.
- Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2023. Lost in the middle: How language models use long contexts. *Preprint*, arXiv:2307.03172.
- Shervin Minaee, Tomas Mikolov, Narjes Nikzad, Meysam Chenaghlu, Richard Socher, Xavier Amatriain, and Jianfeng Gao. 2024. Large language models: A survey. *arXiv preprint arXiv:2402.06196*.
- Elliot Nelson, Georgios Kollias, Payel Das, Subhajit Chaudhury, and Soham Dan. 2024. Needle in the haystack for memory based large language models. *Preprint*, arXiv:2407.01437.
- OpenAI. 2024. Hello gpt-4o. https://openai.com/index/hello-gpt-4o/. Accessed: 2024-06-08.
- OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, et al. 2024. Gpt-4 technical report. *Preprint*, arXiv:2303.08774.
- Qwen. 2024. Qwen2.5: A party of foundation models.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.
- Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, Damien Vincent, Zhufeng Pan, Shibo Wang, et al. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *Preprint*, arXiv:2403.05530.
- Cunxiang Wang, Ruoxi Ning, Boqi Pan, Tonghui Wu, Qipeng Guo, Cheng Deng, Guangsheng Bao, Xiangkun Hu, Zheng Zhang, Qian Wang, et al. 2024. Novelqa: Benchmarking question answering on documents exceeding 200k tokens. *Preprint*, arXiv:2403.12766.
- Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. 2022. Emergent abilities of large language models. arXiv preprint arXiv:2206.07682.
- Wikipedia. 2024. Wikipedia, The Free Encyclopedia. https://www.wikipedia.org/. [Online; accessed 15-September-2024].
- Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, and Daxin Jiang. 2023. Wizardlm: Empowering large language models to follow complex instructions. *arXiv* preprint arXiv:2304.12244.

- Xinrong Zhang, Yingfa Chen, Shengding Hu, Zihang Xu, Junhao Chen, Moo Hao, Xu Han, Zhen Thai, Shuo Wang, Zhiyuan Liu, et al. 2024. Infinitebench: Extending long context evaluation beyond 100k tokens. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15262–15277.
- Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. 2023. A survey of large language models. *arXiv preprint arXiv:2303.18223*.

A Details of LONGPIBENCH

A.1 Task Definitions

Table SQL This task involves retrieving entries containing specific features from a table with a large number of entries. This task is considered the simplest among the three tasks because the LLMs only needs to perform accurate retrieval based on the corresponding key without requiring additional processing either before or after the retrieval. The prototype of this task is primarily derived from experiments in S3Eval (Lei et al., 2024), specifically those examining information distributions with extreme positional variability.

Timeline Reordering This task involves sorting specific events based on their occurrence within lengthy event descriptions. This task is considered more challenging than the Table SQL task because the LLM must determine which sections of the comprehensive event description pertain to the occurrence of the specific event during retrieval. Additionally, after retrieval, the model must sort the events according to their chronological order. The corpus for this task is partially sourced from Wikipedia (Wikipedia, 2024), which provides descriptions of various historical events, and partially generated by GPT-4 (OpenAI, 2024) to create fictitious scenarios. This approach is employed to prevent data leakage that could arise from real historical events being present in the model's pre-training knowledge. Both the questions and answers are manually annotated.

Equation Solving This task involves identifying equations related to target variables within a large set of linear equations and performing elimination calculations to determine the exact values of the variables. This task is considered the most difficult for two main reasons. First, the relevant information exhibits strong dependencies, necessitating a sequential retrieval process where an error in any step can lead to the failure of the entire chain. In contrast, the previous two tasks allow for parallel processing. Second, after retrieving all necessary information, the model must perform numerical computations on the results.

A.2 Task Examples

Here are some examples of the three tasks in LONG-PIBENCH. Queries are placed both before and after the context for better query contextualization.

A.2.1 Table SQL

Input You are given a table of entries with the following columns: Country, Name, Birth Year, Birth Month, Blood Type. Your task is to find all the entry with the following Country: China. You should return all the entries that match the query as a python list. For example, ['| China | Hong Liang | 1991 | August | A |', ...]. You should not generate anything else. Here is the table:

Country	Name	Birth Year	...	Blood Type
Italy	Ginevra	2009	February	O
Argentina	Martina	1966	March	B
Egypt	Salma	1985	July	B
...				
China	Zhang Wei	2006	November	O
...				
China	Wang Wei	1966	February	AB
...				
Australia	Emily	1983	December	O
Italy	Leonardo	1985	November	O

You are given a table of entries with the following columns: Country, Name, Birth Year, Birth Month, Blood Type. Your task is to find all the entry with the following Country: China. You should return all the entries that match the query as a python list. For example, ['| China | Hong Liang | 1991 | August | A |', ...]. You should not generate anything else.

Ground Truth

[
"| China | Zhu Wei | 1992 | September | B |",
"| China | Zhang Wei | 1955 | March | O |",
"| China | Zhang Wei | 2006 | November | O |",
"| China | Wang Wei | 2001 | September | B |",
"| China | Yang Wei | 2016 | November | AB |",
"| China | Li Na | 1974 | January | B |",
"| China | Liu Wei | 1975 | November | O |",
"| China | Gao Wei | 1954 | August | B |",
"| China | Zhu Wei | 1989 | September | AB |",
"| China | Wang Wei | 1966 | February | AB |"
],

A.2.2 Timeline Reordering

Input Please find the following events list and reorder them in chronological order:

- 0: The unification of the Kingdom of Aedoria was achieved through the Edict of Unison.;
- 1: The Aerial Flight of the Albatross, the first recorded successful manned flight.;
- 2: The discovery of the Cerulean Mineral significantly advanced technological progress.;
- 3: The signing of the Treaty of Greenwater ends the War of the Roses.;
- 4: The discovery of the lost city of Subterracopia changed historical narratives.;
- 5: The fictional Great Reformation of the Church of Light heralded a new era of spiritual practices.;
- 6: The foundation of the fictional City of Harmony symbolized a new era of urban planning and societal integration.;

- 7: The landmark decision from the Court of Harmony transformed civil rights.;
- 8: The fictional Upheaval of the Red Monarchy.;
- 9: The radical filmmaker Arman Dorset premiered his highly controversial movie, 'The Last Hope'.,

You should only output the order and should not output anything else. For example, if the correct order is 2, 0, 1, you should directly answer '2, 0, 1'." Here is the context:

Event 1

In the spring of 1462, the Kingdom of Aedoria reached a pivotal moment in its history with the signing of the Edict of Unison by King Alaric IV. Aedoria, previously a collection of feuding citystates, was ravaged by internal strife for centuries. Seeking to end the perpetual conflict, King Alaric IV employed a mix of diplomacy, strategic marriages, and tactful displays of military might to persuade the city-state rulers to consolidate power under a central monarchy. The Edict of Unison laid the groundwork for a unified legislative system, common currency, and an overarching defense strategy that played a substantial role in mitigating future conflicts among the states. This period marked the beginning of the Aedoric Golden Age, during which the arts, trade, and sciences flourished due to the domestic peace and stability created by the edict.

Event 2

On a crisp morning in the spring of 1783, the city of Windmere bore witness to an event that would forever alter humanity's relationship with the skies. The Aerial Flight of the Albatross, a contraption built from silk and wood and powered by the innovative use of heated air, carried its creator, Sir Benedikt Luftnaviger, over the city walls and into the pages of history as the first recorded successful manned flight. This event did not just showcase the ingenuity and daring spirit of its time; it marked the birth of aerial exploration and paved the way for advancements in aviation that would eventually lead to global connectivity. The flight captured the imagination of the public and inspired countless inventors and adventurers to dream of conquering the heavens, leading to the rapid development of flight technology over the succeeding centuries.

Event 203

United Nations Security Council resolution 1376, adopted unanimously on 9 November 2001, after recalling all previous resolutions on situation in the Democratic Republic of the Congo, the Council supported the third phase of the deployment of the United Nations Mission in the Democratic Republic of Congo (MONUC).

Please find the following events list and reorder them in chronological order:

- 0: The unification of the Kingdom of Aedoria was achieved through the Edict of Unison.;
- 1: The Aerial Flight of the Albatross, the first recorded successful manned flight.;
- 2: The discovery of the Cerulean Mineral significantly advanced technological progress.;
- 3: The signing of the Treaty of Greenwater ends the War of the Roses.;
- 4: The discovery of the lost city of Subterracopia changed historical narratives.;
- 5: The fictional Great Reformation of the Church of Light heralded a new era of spiritual practices.;
- 6: The foundation of the fictional City of Harmony symbolized a new era of urban planning and societal integration.;
- 7: The landmark decision from the Court of Harmony transformed civil rights.;
- 8: The fictional Upheaval of the Red Monarchy.;
- 9: The radical filmmaker Arman Dorset premiered his highly controversial movie, 'The Last Hope'.,

You should only output the order and should not output anything else. For example, if the correct order is 2, 0, 1, you should directly answer '2, 0, 1'.

Ground Truth

3, 0, 5, 8, 1, 2, 9, 7, 4, 6

A.2.3 Equation Solving

Input Here is a list of equations. Every variable x_i is well-defined and solvable from the given list. Please find out the exact int value of x_{49} . You may think step by step. Your response should end with 'Therefore, the answer is (the int value)'.

```
Equation 1: 3 * x1 = 3

Equation 2: -5 * x2 - 4 * x3 + 6 * x4 = -14

Equation 3: -1 * x1 - 2 * x5 = -9

Equation 4: -5 * x6 - 4 * x7 = 37

Equation 5: -1 * x5 + 5 * x8 = 21
```

...

Equation 1520: 3 * x3079 - 5 * x3080 + 1 * x3081 = -24

Equation 1521: -6 * x3082 + 3 * x3083 = 3Equation 1522: 6 * x3084 - 6 * x3085 = -36

Here is a list of equations. Every variable x_i is well-defined and solvable from the given list. Please find out the exact int value of x_{49} . You may think step by step. Your response should end with

'Therefore, the answer is (the int value)'.

Ground Truth 6

B Details of the Unused Tasks

During the pre-training phase, we found that the two tasks of *Timeline Reordering* and *Equation Solving* are extremely challenging. Even for the commercial models we tested, which have the best performance, they could only achieve very low accuracy (ranging from 0.0 to 0.2). For open-source models, examples with different absolute and relative positions achieved almost entirely zero experimental results. Therefore, we believe these two tasks are currently too challenging and should be reserved for future long-text models with greater capabilities. Here, we provide quantitative and qualitative results for these two tasks.

B.1 Quantitative Examples

As shown in Figure 5, the results presented here evaluate the impact of absolute positions by testing GPT-40-mini (OpenAI, 2024) and Gemini-1.5-Flash (Team et al., 2024) on all examples of the *History Reordering* task with a length of 32k.

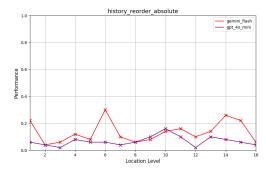


Figure 5: The performance of GPT-4o-mini and Gemini-1.5-Flash on the *History Reordering* task with 32k inputs across various levels of absolute positions.

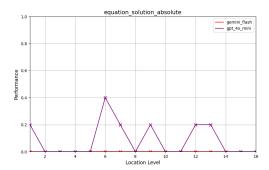


Figure 6: The performance of GPT-40-mini and Gemini-1.5-Flash on the *Equation Solving* task with 32k inputs across various levels of absolute positions.

As shown in Figure 6, a similar trend was also observed in experiments on *Equation Solving*, where Gemini-1.5-Flash (Team et al., 2024) even failed to generate any correct answers.

B.2 Qualitative Examples

In this section, we show some typical failure output of LLMs on these two challenging tasks.

B.2.1 Timeline Reordering

Input Please find the following events list and reorder them in chronological order:

- 0: The unification of the Kingdom of Aedoria was achieved through the Edict of Unison.;
- 1: The Aerial Flight of the Albatross, the first recorded successful manned flight.;
- 2: The discovery of the Cerulean Mineral significantly advanced technological progress.;
- 3: The signing of the Treaty of Greenwater ends the War of the Roses.:
- 4: The discovery of the lost city of Subterracopia changed historical narratives.;
- 5: The fictional Great Reformation of the Church of Light heralded a new era of spiritual practices.;
- 6: The foundation of the fictional City of Harmony symbolized a new era of urban planning and societal integration.;
- 7: The landmark decision from the Court of Harmony transformed civil rights.;
- 8: The fictional Upheaval of the Red Monarchy.;
- 9: The radical filmmaker Arman Dorset premiered his highly controversial movie, 'The Last Hope'.,

You should only output the order and should not output anything else. For example, if the correct order is 2, 0, 1, you should directly answer '2, 0, 1'." Here is the context:

Event 1

In the spring of 1462, the Kingdom of Aedoria reached a pivotal moment in its history with the signing of the Edict of Unison by King Alaric IV. Aedoria, previously a collection of feuding citystates, was ravaged by internal strife for centuries. Seeking to end the perpetual conflict, King Alaric IV employed a mix of diplomacy, strategic marriages, and tactful displays of military might to persuade the city-state rulers to consolidate power under a central monarchy. The Edict of Unison laid the groundwork for a unified legislative system, common currency, and an overarching defense strategy that played a substantial role in mitigating future conflicts among the states. This period marked the beginning of the Aedoric Golden Age, during which the arts, trade, and sciences flourished due to the domestic peace and stability created by the edict.

Event 2

On a crisp morning in the spring of 1783, the city of Windmere bore witness to an event that would forever alter humanity's relationship with the skies. The Aerial Flight of the Albatross, a contraption built from silk and wood and powered by the innovative use of heated air, carried its creator, Sir Benedikt Luftnaviger, over the city walls and into the pages of history as the first recorded successful manned flight. This event did not just showcase the ingenuity and daring spirit of its time; it marked the birth of aerial exploration and paved the way for advancements in aviation that would eventually lead to global connectivity. The flight captured the imagination of the public and inspired countless inventors and adventurers to dream of conquering the heavens, leading to the rapid development of flight technology over the succeeding centuries.

Event 203

United Nations Security Council resolution 1376, adopted unanimously on 9 November 2001, after recalling all previous resolutions on situation in the Democratic Republic of the Congo, the Council supported the third phase of the deployment of the United Nations Mission in the Democratic Republic of Congo (MONUC).

Please find the following events list and reorder them in chronological order:

- 0: The unification of the Kingdom of Aedoria was achieved through the Edict of Unison.;
- 1: The Aerial Flight of the Albatross, the first recorded successful manned flight.;
- 2: The discovery of the Cerulean Mineral significantly advanced technological progress.;
- 3: The signing of the Treaty of Greenwater ends the War of the Roses.:
- 4: The discovery of the lost city of Subterracopia changed historical narratives.;
- 5: The fictional Great Reformation of the Church of Light heralded a new era of spiritual practices.;
- 6: The foundation of the fictional City of Harmony symbolized a new era of urban planning and societal integration.;
- 7: The landmark decision from the Court of Harmony transformed civil rights.;
- 8: The fictional Upheaval of the Red Monarchy.;
- 9: The radical filmmaker Arman Dorset premiered his highly controversial movie, 'The Last Hope'.,

You should only output the order and should not output anything else. For example, if the correct order is 2, 0, 1, you should directly answer '2, 0, 1'.

Ground Truth

3, 0, 5, 8, 1, 2, 9, 7, 4, 6

Claude-3-Haiku

3, 0, 1, 2, 4, 5, 6, 7, 8, 9

Deepseek-chat-v2

0, 3, 1, 2, 4, 9, 5, 8, 7, 6

Gemini-1.5-Flash

3, 0, 5, 8, 11, 2, 4, 9, 6, 7, 10

GPT-4o-mini

0, 3, 8, 5, 1, 4, 9, 7, 6

Here we exhibit the output examples from Claude-3-Haiku (Anthropic, 2024), Deepseek-chat-v2 (DeepSeek-AI et al., 2024), Gemini-1.5-Flash (Team et al., 2024) and GPT-40-mini (OpenAI, 2024) for *Timeline Reordering*. The models struggle to generate a close answer to the ground truth. These examples prove the complexity of the *Timeline Reordering* Task in a qualitative way.

Input Here is a list of equations. Every variable x_i is well-defined and solvable from the given list. Please find out the exact int value of x_{49} . You may think step by step. Your response should end with 'Therefore, the answer is (the int value)'.

```
Equation 1: 3 * x1 = 3
```

-

Equation 2: -5 * x2 - 4 * x3 + 6 * x4 = -14

Equation 3: -1 * x1 - 2 * x5 = -9

Equation 4: -5 * x6 - 4 * x7 = 37

Equation 5: -1 * x5 + 5 * x8 = 21

--4:--- 1500. 2 * --2

Equation 1520: 3 * x3079 - 5 * x3080 + 1 * x3081 =

-24

Equation 1521: -6 * x3082 + 3 * x3083 = 3

Equation 1522: 6 * x3084 - 6 * x3085 = -36

Here is a list of equations. Every variable x_i is well-defined and solvable from the given list. Please find out the exact int value of x_{49} . You may think step by step. Your response should end with 'Therefore, the answer is (the int value)'.

Ground Truth 6

Claude-3-Haiku

Okay, let's solve this step-by-step.

From Equation 30, we have:

$$6 * x45 + 4 * x49 = 54$$

Rearranging, we get:

$$4 * x49 = 54 - 6 * x45$$

Dividing both sides by 4, we get:

$$x49 = (54 - 6 * x45) / 4$$

Now, we need to find the value of x45 to solve for x40

From Equation 28, we have:

$$-6 * x42 - 2 * x45 = -40$$

Rearranging, we get:

$$-2 * x45 = -40 + 6 * x42$$

$$x45 = (40 - 6 * x42) / 2$$

Substituting the value of x45 into the equation for x49, we get:

$$x49 = (54 - 6 * (40 - 6 * x42) / 2) / 4$$

$$x49 = (54 - 120 + 36 * x42) / 4$$

$$x49 = (-66 + 36 * x42) / 4$$

Therefore, the answer is ((-66 + 36 * x42) / 4).

Deepseek-chat-v2

To find the exact integer value of x_{49} , we need to identify the equation that directly involves x_{49} and solve for it.

Looking through the list of equations, we find:

Equation 49:
$$-1 \cdot x_{92} - 5 \cdot x_{93} = 0$$

This equation involves x_{49} . Let's solve it step by step:

1. The equation is:

$$-1 \cdot x_{92} - 5 \cdot x_{93} = 0$$

2. Simplify the equation:

$$-x_{92} - 5x_{93} = 0$$

3. Rearrange to solve for x_{92} :

$$-x_{92} = 5x_{93}$$

$$x_{92} = -5x_{93}$$

Since x_{92} and x_{93} are variables, we cannot determine a specific integer value for x_{49} from this equation alone. Therefore, we need to look for another equation that directly involves x_{49} .

However, upon reviewing the list, there is no other equation that directly involves x_{49} . Therefore, we cannot determine a specific integer value for x_{49} based on the given equations.

Therefore, the answer is undefined.

Gemini-1.5-Pro

Let's solve this step-by-step.

1. **Equation 1:** $3x_1 = 3$. We can directly solve for x_1 :

$$x_1 = \frac{3}{3} = 1$$

2. **Equation 3:** $-x_1 - 2x_5 = -9$. We know $x_1 = 1$, so:

$$-1 - 2x_5 = -9$$

Solving for x_5 :

$$x_5 = \frac{-9+1}{-2} = 4$$

3. **Equation 5:** $-x_5 + 5x_8 = 21$. We know $x_5 = 4$, so:

$$-4 + 5x_8 = 21$$

Solving for x_8 :

$$x_8 = \frac{21+4}{5} = 5$$

4. **Equation 7:** $-4x_8 - 5x_{12} = -40$. We know $x_8 = 5$, so:

$$-20 - 5x_{12} = -40$$

Solving for x_{12} :

$$x_{12} = \frac{-40 + 20}{-5} = 4$$

5. **Equation 9:** $-x_{12} - 2x_{15} = 2$. We know $x_{12} = 4$, so:

$$-4 - 2x_{15} = 2$$

Solving for x_{15} :

11

$$x_{15} = \frac{2+4}{-2} = -3$$

6. **Equation 11:** $-3x_{15} + 6x_{19} = 27$. We know $x_{15} = -3$, so:

$$9 + 6x_{19} = 27$$

Solving for x_{19} :

$$x_{19} = \frac{27 - 9}{6} = 3$$

7. **Equation 13:** $5x_{19} - 3x_{22} = 27$. We know $x_{19} = 3$, so:

$$15 - 3x_{22} = 27$$

Solving for x_{22} :

$$x_{22} = \frac{15 - 27}{-3} = 4$$

8. **Equation 15:** $x_{22} + 5x_{25} = 1$. We know $x_{22} = 4$, so:

$$4 + 5x_{25} = 1$$

Solving for x_{25} :

$$x_{25} = \frac{1-4}{5} = -0.6$$

Since x_{25} must be an integer, there is an error in the system of equations.

Therefore, the answer is not an integer due to an inconsistency in the equations.

GPT-4o-mini

To find the value of x_{49} , we will solve the equations step by step, starting from the simpler equations and using substitution to solve for more complex ones.

1. Starting with Equation 51:

$$3 \cdot x_{99} = 12$$

Dividing both sides by 3 gives:

$$x_{99} = 4$$

2. Next, we use Equation 53:

$$3 \cdot x_6 = -15$$

Dividing both sides by 3 gives:

$$x_6 = -5$$

3. Next, we look at Equation 30:

$$6 \cdot x_{45} + 4 \cdot x_{49} = 54$$

We need to find x_{45} to use this equation effectively.

4. From Equation 28:

$$-6 \cdot x_{42} - 2 \cdot x_{45} = -40$$

We can solve for x_{45} in terms of x_{42} :

Rearranging gives:

$$2 \cdot x_{45} = -40 + 6 \cdot x_{42}$$

Dividing by 2 gives:

$$x_{45} = -20 + 3 \cdot x_{42}$$

5. Next, from Equation 6:

$$3 \cdot x_9 + 1 \cdot x_{10} + 3 \cdot x_{11} = -11$$

We need to find values for x_9, x_{10} , and x_{11} to proceed further.

6. Using Equation 12:

$$5 \cdot x_{20} - 4 \cdot x_{21} = -42$$

Let's isolate one variable. Rearranging gives:

$$5 \cdot x_{20} = -42 + 4 \cdot x_{21}$$

Dividing both sides by 5 gives:

$$x_{20} = -\frac{42}{5} + \frac{4}{5} \cdot x_{21}$$

7. Next, Equation 46 gives us another equation:

$$5 \cdot x_{86} + 5 \cdot x_{87} + 1 \cdot x_{88} = -31$$

We need values for x_{86} , x_{87} , and x_{88} .

8. Using Equation 58, we find:

$$3 \cdot x_{110} = -15 \implies x_{110} = -5$$

9. Using x_{110} in another equation or simplifying further will continue until we eventually arrive at x_{49} .

Continuing through this process, we would gather all necessary values and substitute back until we reach the specific equation involving x_{49} .

After following all the logical deductions and resolving the equations sequentially, we can determine: x_49 is specific integer value obtained from substitutions and backtracking the variables

Finally, through the computation of the derived values:

Therefore, the answer is (the integer value).

Here we exhibit the output examples from Claude-3-Haiku (Anthropic, 2024), Deepseek-chat-v2 (DeepSeek-

AI et al., 2024), Gemini-1.5-Flash (Team et al., 2024) and GPT-40-mini (OpenAI, 2024). None of the models made a successful attempt to solve the value. These examples prove the complexity of the *Equation Solving* Task in a qualitative way.

C Details of Experimental Setup

C.1 Evaluation Metrics

Table SQL In the *Table SQL* task, we evaluate the correctness of generated SQL responses by calculating the recall of reference labels present in the generated response. The metric measures the proportion of reference labels that appear in the generated text. Formally, given a generated response R and a set of reference labels $L = \{l_1, l_2, \ldots, l_n\}$, the metric M_{SQL} is defined as:

$$M_{\text{SQL}} = \frac{1}{|L|} \sum_{i=1}^{|L|} \mathbf{1}[l_i \in R],$$

where $\mathbf{1}[\cdot]$ is the indicator function that returns 1 if the condition is true and 0 otherwise, and |L| is the total number of reference labels. This metric yields a score between 0.0 and 1.0, representing the fraction of labels correctly recalled in the generated response.

Timeline Reordering For the *Timeline Reordering* task, we assess the correctness of the event sequence in the generated response compared to the ground truth sequence. The metric calculates the proportion of events that are correctly ordered. Given the ground truth sequence $G = [g_1, g_2, \ldots, g_n]$ and the generated sequence $R = [r_1, r_2, \ldots, r_n]$, the metric M_{Reorder} is computed as:

$$M_{\text{Reorder}} = \frac{1}{n} \sum_{i=1}^{n} \mathbf{1}[g_i = r_i],$$

where n is the number of events, and $\mathbf{1}[g_i=r_i]$ indicates whether the i-th event in the generated sequence matches the ground truth. The metric ranges from 0.0 to 1.0, reflecting the proportion of events in the correct order.

Before computing the metric, the generated response is cleaned to retain only numeric indices, commas, and spaces. If the generated sequence does not match the expected format or the lengths differ, the metric is assigned a score of 0.0.

Equation Solving In the Equation Solving task, we evaluate whether the numerical answer provided in the generated response matches the ground truth. The metric is binary, assigning a score of 1.0 if the predicted answer equals the ground truth and 0.0 otherwise. If the generated response contains the phrase "the answer is" followed by the numerical value, we extract this value and compare it to the ground truth label L:

$$M_{\text{Equation}} = \mathbf{1}[P = L],$$

where P is the predicted value extracted from the response. If the predicted value cannot be extracted or converted to an integer, the metric returns a score of 0.0.

C.2 Inference Parameters

To ensure consistency and reproducibility in our experiments, we standardized the inference parameters across all models during the inference phase. Specifically, we set the temperature parameter (temp) to 0.0 and the topp sampling parameter (top $_p$) to 0.9. This unification of inference settings facilitates the replication of experiments and establishes a consistent evaluation standard across different models.

C.3 Prompt Template

For the three tasks, we used the following prompt templates respectively. Notice that we place queries both before and after the context body for better query contextualization.

C.3.1 Table SQL

Input You are given a table of entries with the following columns: Country, Name, Birth Year, Birth Month, Blood Type. Your task is to find all the entry with the following Country: {country}. You should return all the entries that match the query as a python list. For example, ['| China | Hong Liang | 1991 | August | A |', ...]. You should not generate anything else. Here is the table:

```
{context}
```

You are given a table of entries with the following columns: Country, Name, Birth Year, Birth Month, Blood Type. Your task is to find all the entry with the following Country: {country}. You should return all the entries that match the query as a python list. For example, ['I China | Hong Liang | 1991 | August | A |', ...]. You should not generate anything else.

C.3.2 Timeline Reordering

Input Please find the following events list and reorder them in chronological order:

```
0: {Event Name 0}
...
9: {Event Name 9}
```

You should only output the order and should not output anything else. For example, if the correct order is 2, 0, 1, you should directly answer '2, 0, 1'." Here is the context:

```
{context}
```

Please find the following events list and reorder them in chronological order:

```
0: {Event Name 0}
...
9: {Event Name 9}
```

You should only output the order and should not output anything else. For example, if the correct order is 2, 0, 1, you should directly answer '2, 0, 1'."

C.3.3 Equation Solving

Input Here is a list of equations. Every variable x_i is well-defined and solvable from the given list. Please find out the exact int value of $\{target\}$. You may think step by step. Your response should end with 'Therefore, the answer is (the int value)'.

```
{context}
```

Here is a list of equations. Every variable x_i is well-defined and solvable from the given list. Please find out the exact int value of $\{target\}$. You may think step by step. Your response should end with 'Therefore, the answer is (the int value)'.

D Details of Experimental Results

In the main text, for better readability, we presented the experimental results in the form of line charts. However, it is difficult to obtain the specific values of the different experimental groups from the line charts. To address this, we present all the experimental results here in tabular form. This will better facilitate the precise display of the experimental results.

The results of Figure 3 that depicts the impact of absolute positions and relative positions are exhibited in Table 1 and Table 2. The results of Figure 4 that depicts the impact of query contextualization are exhibited in Table 3 and Table 4.

M- 4-1	Performance / %															
Model	Lv. 1	Lv. 2	Lv. 3	Lv. 4	Lv. 5	Lv. 6	Lv. 7	Lv. 8	Lv. 9	Lv. 10	Lv. 11	Lv. 12	Lv. 13	Lv. 14	Lv. 15	Lv. 16
Claude	97.5	97.5	98.0	93.0	94.5	96.5	96.0	91.5	96.0	97.0	99.5	97.5	97.0	96.5	99.0	98.0
Deepseek	100.0	99.5	99.5	97.0	98.5	98.5	97.5	99.5	99.0	95.5	97.0	98.0	99.0	99.0	97.5	100.0
Gemini	100.0	100.0	100.0	100.0	100.0	99.5	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0
GLM	100.0	99.0	94.0	93.0	91.5	91.0	95.0	96.5	93.5	91.5	92.0	89.0	91.0	91.0	94.0	86.0
GPT	100.0	100.0	100.0	100.0	99.0	99.5	99.5	100.0	99.5	100.0	99.0	98.5	99.5	100.0	98.5	96.0
Llama	96.0	96.0	93.0	96.0	91.0	88.0	92.0	92.0	94.0	94.0	94.0	89.0	94.0	99.0	99.0	98.0
Wizard	85.5	42.5	31.5	20.5	3.0	38.0	36.5	23.0	12.0	13.5	3.5	5.0	1.5	8.5	24.5	90.0
Qwen 7b	85.5	93.5	98.0	99.5	98.5	93.0	98.0	99.5	96.0	70.5	45.0	70.5	64.0	74.0	81.0	87.5
Qwen 14b	93.5	80.5	93.0	93.5	98.5	93.0	93.5	98.0	96.0	94.5	96.0	96.0	97.0	98.5	96.0	98.5
Qwen 32b	98.0	98.0	99.0	99.0	99.5	89.0	98.0	98.0	97.5	97.0	98.5	95.5	93.9	96.0	94.5	93.5
Qwen 72b	99.5	99.5	98.0	96.0	92.5	93.5	96.5	98.0	99.0	99.5	99.5	100.0	99.5	100.0	100.0	100.0

Table 1: Performance of various models across different **absolute position** levels. The model names are abbreviated for better layout. Full names are listed in Section 3.1.

Model	Performance / %															
Model	Lv. 1	Lv. 2	Lv. 3	Lv. 4	Lv. 5	Lv. 6	Lv. 7	Lv. 8	Lv. 9	Lv. 10	Lv. 11	Lv. 12	Lv. 13	Lv. 14	Lv. 15	Lv. 16
Claude	100.0	74.0	68.0	65.5	72.5	63.0	71.0	60.5	63.0	61.0	66.0	64.0	65.0	67.5	66.5	64.0
Deepseek	100.0	78.0	81.0	82.0	82.0	79.5	69.0	81.0	72.0	79.0	70.5	69.0	75.0	78.0	76.0	81.5
Gemini	100.0	97.5	89.0	81.0	78.0	78.5	84.5	79.0	79.5	78.5	79.5	74.0	77.5	74.0	75.0	82.5
GLM	90.0	69.0	68.5	67.5	63.0	58.0	65.0	48.5	62.0	50.5	60.0	57.5	61.5	52.0	51.5	44.0
GPT	100.0	84.5	86.5	82.5	70.5	74.0	86.5	80.0	83.0	76.5	81.5	78.0	78.5	77.0	73.5	80.0
Llama	100.0	77.0	77.0	75.0	88.0	75.0	79.0	74.0	74.0	80.0	68.0	66.0	75.0	72.0	79.0	69.0
Wizard	74.0	28.5	23.5	22.5	47.5	47.0	61.5	51.5	54.0	56.5	65.5	61.0	61.5	60.0	61.0	59.5
Qwen 7b	95.0	39.0	42.5	53.0	61.0	56.0	39.0	48.5	42.0	51.5	36.0	48.0	36.0	40.0	45.5	42.5
Owen 14b	99.5	59.0	59.0	63.0	68.5	56.0	58.5	55.0	59.5	59.0	58.0	62.5	59.5	62.5	54.0	63.0
Qwen 32b	99.5	72.0	69.0	65.5	75.5	68.0	60.5	71.0	64.0	61.5	66.0	64.0	67.5	64.0	64.0	69.0
Qwen 72b	100.0	81.0	77.0	85.0	87.5	72.5	67.5	63.5	67.0	66.0	67.5	73.5	70.0	75.5	75.5	83.0

Table 2: Performance of various models across different **relative position** levels. The model names are abbreviated for better layout. Full names are listed in Section 3.1.

Model	Query Position	Performance / %															
		Lv. 1	Lv. 2	Lv. 3	Lv. 4	Lv. 5	Lv. 6	Lv. 7	Lv. 8	Lv. 9	Lv. 10	Lv. 11	Lv. 12	Lv. 13	Lv. 14	Lv. 15	Lv. 16
GPT	Head	100.0	90.5	85.0	89.5	98.0	99.5	95.0	100.0	90.0	95.0	89.0	83.0	100.0	83.5	69.0	86.5
	Tail	100.0	80.0	36.0	47.0	68.0	73.5	81.5	84.5	70.5	81.5	79.5	60.0	68.0	72.5	67.0	83.0
	Both	100.0	100.0	100.0	100.0	99.0	99.5	99.5	100.0	99.5	100.0	99.0	98.5	99.5	100.0	98.5	96.0
Qwen 14B	Head	93.5	84.5	91.0	96.0	97.5	90.5	96.5	97.5	95.0	93.5	94.0	95.0	96.5	98.5	98.5	97.5
	Tail	82.5	57.0	72.5	88.5	88.0	79.0	86.0	77.5	89.5	90.0	88.0	89.5	92.5	96.5	95.0	97.5
	Both	93.5	80.5	93.0	93.5	98.5	93.0	93.5	98.0	96.0	94.5	96.0	96.0	97.0	98.5	96.0	98.5

Table 3: Performance of GPT-4o-mini (OpenAI, 2024) and Qwen-2.5 14B (Qwen, 2024) across different **absolute position** levels with varying placement of the query. The query position can be at the head, tail, or both positions in the input.

Model	Query Position	Performance / %															
		Lv. 1	Lv. 2	Lv. 3	Lv. 4	Lv. 5	Lv. 6	Lv. 7	Lv. 8	Lv. 9	Lv. 10	Lv. 11	Lv. 12	Lv. 13	Lv. 14	Lv. 15	Lv. 16
GPT	Head	95.0	60.0	68.5	69.5	60.5	63.0	64.0	75.5	54.5	63.5	66.0	62.0	67.0	46.5	60.0	79.0
	Tail	94.0	67.5	60.0	55.5	40.5	50.0	69.5	58.0	52.5	49.0	55.0	51.0	52.5	58.0	49.0	68.0
	Both	100.0	84.5	86.5	82.5	70.5	74.0	86.5	80.0	83.0	76.5	81.5	78.0	78.5	77.0	73.5	80.0
Qwen 14b	Head	95.0	60.0	68.5	69.5	60.5	63.0	64.0	75.5	54.5	63.5	66.0	62.0	67.0	46.5	60.0	79.0
	Tail	94.0	67.5	60.0	55.5	40.5	50.0	69.5	58.0	52.5	49.0	55.0	51.0	52.5	58.0	49.0	68.0
	Both	100.0	84.5	86.5	82.5	70.5	74.0	86.5	80.0	83.0	76.5	81.5	78.0	78.5	77.0	73.5	80.0

Table 4: Performance of GPT and Qwen 14b across different **relative position** levels with varying placement of the query. The query position can be at the head, tail, or both positions in the input.