AbsenceBench: Language Models Can't Tell What's Missing

Harvey Yiyun Fu*, Aryan Shrivastava¹, Jared Moore²
Peter West², Chenhao Tan¹, Ari Holtzman¹

¹University of Chicago ²Stanford University

Abstract

Large language models (LLMs) are increasingly capable of processing long inputs and locating specific information within them, as evidenced by their performance on the Needle in a Haystack (NIAH) test. However, while models excel at recalling surprising information, they still struggle to identify *clearly omitted* information. We introduce AbsenceBench to assesses LLMs' capacity to detect missing information across three domains: numerical sequences, poetry, and GitHub pull requests. AbsenceBench asks models to identify which pieces of a document were deliberately removed, given access to both the original and edited contexts. Despite the apparent straightforwardness of these tasks, our experiments reveal that even state-of-the-art models like Claude-3.7-Sonnet achieve only 69.6% F1-score with a modest average context length of 5K tokens. Our analysis suggests this poor performance stems from a fundamental limitation: Transformer attention mechanisms cannot easily attend to "gaps" in documents since these absences don't correspond to any specific keys that can be attended to. Overall, our results and analysis provide a case study of the close proximity of tasks where models are already superhuman (NIAH) and tasks where models breakdown unexpectedly (AbsenceBench).

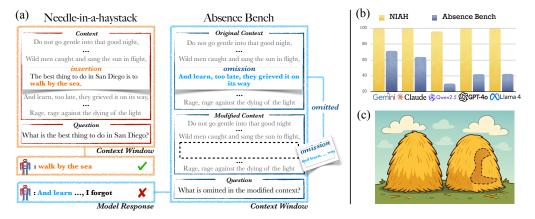


Figure 1: (a) An overview of the difference between the Needle-in-a-haystack (NIAH) test setting and AbsenceBench task setting. AbsenceBench is asking models to identify omitted pieces of content. (b) Performance of 5 SoTA LLMs on AbsenceBench is significantly lower than on the NIAH test, measured by F1-score. (c) an illustration of our task setting using the "haystack" metaphor, generated by ChatGPT.

^{*}Corresponding author(s): harveyfu@uchicago.edu

1 Introduction

Recent Large Language Models (LLMs) have shown exceptional abilities across a wide range of long-context tasks [Bai et al., 2024, Zhang et al., 2024, *inter alia*]. In particular, the Needle-in-a-Haystack (NIAH) test [Kamradt, 2023, Hsieh et al., 2024] has been used to evaluate whether models can find small bits of surprising information in extremely long inputs, with recent models pushing the frontier of NIAH into the millions of tokens. We ask the converse question: **Can LLMs spot information that has clearly been left out?**

To answer this question, we present AbsenceBench: a new benchmark designed to evaluate the abilities of LLMs in locating conspicuously missing information from long inputs. Instead of asking LLMs to find off-topic information (the 'needle' in NIAH), LLMs are prompted to identify and recall intentionally **omitted** information. Cutting-edge, closed-source models perform poorly on AbsenceBench, despite the apparent similarity to NIAH and the fact that AbsenceBench is simple to describe and entirely unambiguous.

AbsenceBench includes three domains: poetry, numerical sequences, and GitHub pull requests (PRs). The average context length of AbsenceBench is 5K, significantly shorter than most long-context benchmarks—we thus call it a medium-context benchmark. We benchmark a total of 14 LLMs on AbsenceBench, including cutting-edge models such as GPT-4 [OpenAI et al., 2024b], Claude-3.7-Sonnet [Anthropic, 2024], and Gemini-2.5-flash [Google, 2025]², as well as inference-time compute models such as o3-mini [OpenAI, 2025], Grok-3-mini [xAI, 2025] and DeepSeek-R1 [DeepSeek-AI et al., 2025]. We further study the impact of context length and the rate of omission. Our results suggest that: (1) Longer context length makes the task harder especially under the poetry domain. (2) Using inference-time compute boost models' performance by only 7.9% at the cost of generating an average of extra 8K thinking tokens—nearly 3x the average document length! (3) Counter-intuitively, a *lower* rate of omission in the task set-up leads to *worse* model performance.

AbsenceBench is also significantly more difficult for LLMs than the NIAH test. To illustrate this, we compare the performance of three LLMs on both our proposed task setting and the original NIAH test setting, observing a a massive 56.9% drop in F1-score on average. Why is AbsenceBench so much more difficult than the NIAH test? We hypothesize that the Transformer [Vaswani et al., 2017] attention may have trouble addressing "gaps" in the document, leading us to experiment with adding placeholder strings where information is omitted. This boosts the performance by a dramatic 35.7% on average (see §4.2).

What explains these results? While LLMs often generate deceptively human-like responses, and may even have human-like cognitive-biases [Koo et al., 2024], the contrast in success between NIAH and AbsenceBench suggests that models may have failed for completely different reasons than humans do. This has practical implications: if LLMs cannot properly notice when information is missing, can we rely on LLM-as-a-Judge [Zheng et al., 2023] for a wide range of tasks where recognizing missing information matters? We hope that AbsenceBench can serve as a starting point for examining LLMs' ability to tell what is missing, and as a case-study of an LLM-specific cognitive bias.

Overall, our contributions are as follows:

- We release a new medium-length-context benchmark for evaluating LLMs' abilities in locating intentionally-omitted information across three diverse domains.
- We evaluate 14 popular LLMs, including those with inference-time compute, and show that our benchmark is challenging even for cutting-edge language models.
- We show that while the NIAH test is essentially solved for very long contexts, AbsenceBench tasks have low performance on medium length contexts, suggesting LLMs find it harder to identify omissions than insertions.
- We analyze the effect of inference-time compute, finding that it offers only a modest performance improvement at a high cost—an average chain-of-thought length that is nearly 3x the size of the average document length.
- Explicitly marking omissions improves models' performance by 35.7% on average. This suggests that AbsenceBench may be caused by inherent weaknesses within Transformer-style self-attention mechanisms.

²Claude-3.7-Sonnet and Gemini-2.5-flash have the option to leverage inference-time compute and we evaluate them both with and without this feature.

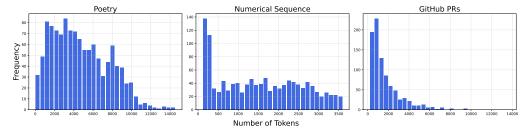


Figure 2: The three domains in AbsenceBench test models' abilities across a variety of document lengths and omission probabilities. Frequency reports the number of tasks in the domain within a given range of document lengths. The average context length across all tasks in AbsenceBench is 5K tokens. On the document level, the average document length is 2.7K, while it is 4.7K for poetry, 1.5K for numerical sequences, and 1.7K for Github pull requests. We use the GPT-4 Tokenizer⁴to measure document and context lengths.

Domain	Original	Modified
Poetry	To me, I give these weedy rhymes In memory of early times	And so, to you, who always were In memory of early times
Numerical	117, 121, 125, 129, 133, 137	$117, 125, 129, 133\dots$
GitHub PRs	+ \$replacements = [+ '[' => '\[', + '<' => '<', + '>' => '>', +];	+ \$replacements = [+ '[' => '\[', + '>' => '>', +];
Absence Token	To me, I give these weedy rhymes In memory of early times	And so, to you, who always were <missing line=""> In memory of early times</missing>

Table 1: Examples of tasks in AbsenceBench by domain. Models are given both the original and modified documents and are asked to identify which elements are missing (in blue). The benchmark includes the Poetry, Numerical Sequences, and Github PRs domains, while the Absence Token domain includes an *omission token* that is used for analysis in §4.2

2 AbsenceBench

AbsenceBench³ is built on a simple idea: LLMs have trouble keeping track of what's missing.

2.1 Task Definition

We formulate all tasks in AbsenceBench into a straightforward controlled generation framework: given an original document composed of n elements $D_{\text{orig}} = \{e_1, \dots, e_n\}$, we intentionally omit p percent of the elements $D_{\text{omit}} = \{e_{r1}, e_{r2}, \dots, e_{rk}\}$, with $k = p \cdot n$, and produce a modified version of the document $D_{\text{modified}} = D_{\text{orig}} \backslash D_{\text{omit}}$. We present both versions of the document D_{orig} and D_{modified} to the LLMs and then prompt them to generate the exact set of omitted elements D_{omit} . We use "document length" to denote the token count of the original document, and "context length" to indicate the token count of the entire input context.

2.2 Dataset Construction

³Our code is available at https://github.com/harvey-fin/absence-bench

⁴https://github.com/openai/tiktoken?tab=readme-ov-file#-tiktoken

```
System Prompt
You are helping a student practice memorizing poems. The student will
recite a poem, but they may have missed some lines. Your task is to
identify exactly which lines are missing from their recitation.
List only the missing lines, nothing else.
User Message
Here is the complete original poem:
{original poem}
Now, here is my recitation which may be missing some lines:
{modified poem}
What lines did I miss? Please list only the missing lines, nothing else.
```

Table 2: Default prompt template used for evaluating language models under the poetry domain

AbsenceBench covers three distinct domains: poetry, numerical sequences, and GitHub pull requests (PRs). Poetry and GitHub PRs are realistic data directly collected from open sources, while numerical sequences are synthetic. We choose p=0.1 for all domains. Overall, AbsenceBench contains 4302 instances in total, with an average context length of 5K tokens. We plot the distribution of document lengths for each domain in Figure 2. We include justifications for our construction parameter choices in Appendix D.

Poetry We use poems from the Gutenberg Poetry Corpus [Parrish, 2018]. For each poem, we omit at a line level and use the newline character as the delimiter between different lines. In order to create diversity in document length, we truncate the poems such that the number of lines for each poem are uniformly distributed from 100 to 1000. We include justification over the choices of dataset construction parameters in Appendix

Numerical Sequences LLMs demonstrate impressive abilities in solving mathematical tasks [Frieder et al., 2023], and have seen numbers frequently in the pre-training data. Nonetheless, it is unclear whether LLMs are able to reliably keep track of numerical sequences. We generate a total of 1200 numerical sequences. Each sequence consists of n decimal numbers $\{a^{(1)}, a^{(2)}, \ldots, a^{(n)}\}$ that are arranged in a particular order $\mathcal{L} \in \{\text{ascending, descending, random}\}$, with a step size $s \in \{1, 4, 7, 13\}$ between two consecutive numbers. The first number $a^{(1)}$ is randomly chosen from 0 to 9999. We omit each number from the set with a certain probability, and delimit numbers with a newline character.

GitHub PRs GitHub is one of the largest open-source code repositories and is frequently used as a high-quality data source for (LLMs) such as Code Llama [Rozière et al., 2024]. We retrieve the PRs from the top 20 GitHub repositories⁵ with the most PRs using GitHub API, and only keep those PRs with a diff that includes 10 to 200 updated lines (i.e., a line that starts with "+" or "-"). We format this task similarly to the poetry domain: omissions are applied at a line level, and we use the newline character as the delimiter. In particular, we only omit the updated lines that are unique within each PR's diff. This domain has practical application to current LLM-usage: an LLM that resolve and verify merge conflicts must inherently be able to detect omissions in file diffs.

3 Experiments

3.1 Experiment Setup

Models We evaluate a total of 14 LLMs, including seven models that leverage inference-time compute (Gemini-2.5-flash [Google, 2025], Claude-3.7-Sonnet [Anthropic, 2024], o3-mini [OpenAI, 2025], DeepSeek R1 [DeepSeek-AI et al., 2025], Grok-3-mini-Beta [xAI, 2025], Qwen3-235B [Qwen, 2025a], QwQ-32B [Qwen, 2025b]), three OpenAI models (GPT4o[OpenAI et al., 2024b], GPT-4.1, GPT-4.1-mini), and four open-weights models (Llama-4-Maverick, Llama-3.3-70B-Instruct [Meta, 2025], Qwen2.5-72B-Instruct [Yang et al., 2024], Mixtral-8x7B-Instruct [Jiang et al., 2024]). We run both Gemini-2.5-flash and Claude-3.7-Sonnet with and without inference-time compute (i.e.,

⁵https://top1000repos.com/based-on-pr

Models	Poetry	Numerical Sequences	GitHub PRs	Average
Gemini-2.5-flash*	87.3	95.4	30.9	71.2
Claude-3.7-Sonnet*	72.7	96.0	40.0	69.6
Claude-3.7-Sonnet	73.5	91.4	35.7	66.9
Gemini-2.5-flash	79.3	85.2	26.2	63.6
o3-mini*	65.0	78.1	38.9	60.7
GPT-4.1	54.3	57.5	36.2	49.3
Grok-3-mini-Beta*	40.7	56.3	36.4	44.5
GPT-4o	38.4	48.1	39.4	42.0
QwQ-32B*	32.1	57.7	31.6	40.5
Llama-4-Maverick	32.8	58.7	29.0	40.2
GPT-4.1-mini	30.2	45.0	31.3	35.5
Llama-3.3-70B-Instruct	25.3	37.7	28.7	30.6
Qwen2.5-72B-Instruct	19.0	45.4	26.8	30.4
DeepSeek-R1*	38.7	29.5	23.1	30.4
Qwen3-235B*	26.1	18.5	24.6	23.1
Llama-3.1-8B-Instruct	9.4	17.1	21.6	16.0
Mixtral-8x7B-Instruct	4.9	21.9	17.3	14.7

Table 3: Micro F1-score (%) of 15 LLMs evaluated on all three domains of Absence Bench. * indicates the model uses inference-time compute during evaluation. GitHub PRs represent the most challenging domain: models that perform well in other domains often struggle there.

"thinking mode"). All selected models have a claimed context length of over 32K tokens. All model inferences are obtained through API requests (see Appendix B for details).

Evaluation Metric Our evaluation is two-fold: (1) we use exact match to check whether each element (e.g., a line of a poem) is present in a model's response, and (2) we use **micro F1-score** to evaluate whether models generate the correct set of case-insensitive elements. Note that we use a different evaluation metric than the recall accuracy metric that is frequently used in the NIAH test. In our task setting, a model could achieve a perfect recall score simply by copy-pasting the entire original context. However, this would result in many false positives cases (i.e., when model generates non-omitted elements) that recall accuracy does not account for.

Prompt Templates & In-context Learning Table 2 presents an example of the prompt template used for evaluation in the poetry domain. See Appendix A for detailed prompts under all domains. Due to compute and context length constraints, we include limited perturbation studies on the prompt template and in-context learning examples in this paper and leave it for future work. See Appendix C for a perturbation study on the prompt template and §7 for further discussion.

3.2 Main Results

We show the evaluation results of 14 LLMs on AbsenceBench in Table 3. Gemini-2.5-flash (thinking)⁶ outperforms the rest of the models on poetry, numerical sequences, and overall by a significant margin, followed by Claude-3.7-Sonnet (thinking). Most open-weights models are generally weaker in performance and struggle to reach 40% F1-score except for QwQ-32B and Llama-4-Maverick. GitHub PRs is a challenging domain for all models: the highest score is only 40.0% by Claude-3.7-Sonnet (thinking). It is important to highlight that while Mixtral-8x7B achieves a perfect score on the NIAH test at an equivalent context length, it attains only a 14.7% F1-score on the AbsenceBench, indicating substantial space for improvement among smaller-scale models. Overall, AbsenceBench presents a surprisingly challenging task to cutting-edge language models.

Inference-time compute is helpful but costly We consider seven models that include inference-time compute capabilities. We plot the distribution of *thinking token ratio*: the ratio between the total thinking tokens generated by seven inference-time compute models and the original document length under all three domains in Figure 3. The figure shows that reasoning models typically generate thinking tokens several times greater than the document length, which further suggests that the models

⁶We use (thinking) to denote the model that use inference-time compute during evaluation

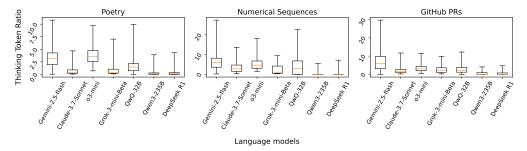


Figure 3: Reasoning models often generate an order of magnitude more text than input document. Distribution of the *thinking token ratio* (number of generated thinking tokens divided by number of tokens in the original document) for four inference-time compute models under each domain. We set the parameters of the boxplot to capture 99% of the distribution. The outliers are hidden for better clarity (see Figure 8 for the full distribution).

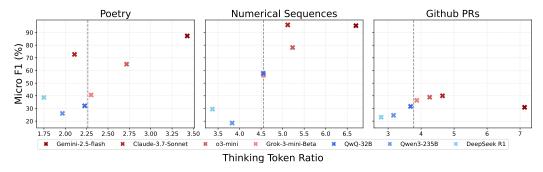


Figure 4: Closed-source models (reds) perform better than open-weights models (blues) on AbsenceBench, while generating more thinking tokens. Each plot shows the average F1-score (x-axis) and the average thinking token ratio (y-axis). The grey line presents a visual boundary.

may attempt to reconstruct the original document using thinking token as a way to identify omissions. Notably, Gemini-2.5-flash shows the highest thinking token ratio of 5.7 on average across three domains, partly accounting for its advantage in benchmark performance. Additionally, we compare the performance of Gemini-2.5-flash and Claude-3.7-Sonnet under conditions with and without the use of "thinking mode". We observe that inference-time compute leads to a modest performance improvement of 7.9%, but it comes at the cost of producing an additional 8K tokens for intermediate reasoning steps on average—significantly longer than our average context length of 5K tokens.

Closed-source models outperform open-weights models We plot the average F1-score and thinking token ratio of four closed-source models and three open-weights models in Figure 4. We observe that closed-source models generally achieve higher average F1-scores than open-weight models across all domains, with the exception of QwQ-32B, which outperforms Grok-3-mini-Beta in numerical sequences and Gemini-2.5-flash in GitHub PRs. On the other hand, the higher performance of closed-source models come at the cost of generating 42.0% more thinking tokens than open-weights models on average across all domains.

Locating omissions are harder than insertions The increased difficulty of AbsenceBench relative to the original NIAH test could potentially be attributed to differences in the task design and the higher frequency of document modifications (insertions or omissions). To investigate this, we run Claude-3.7-Sonnet, GPT-4.1-mini, and Llama-4-Maverick under an "insertion bench" setting: rather than omitting elements from the context, we insert "needles" at an equivalent frequency. We choose the Harry Potter dataset⁷ as the needles and two realistic domains (poetry and GitHub PRs) as the haystack. All three models achieve nearly 99.5% F1-score on the poetry domain, and at least 86.2% F1-score on Github PRs (see Appendix E for full results). By comparison, models

⁷https://www.kaggle.com/datasets/gulsahdemiryurek/harry-potter-dataset

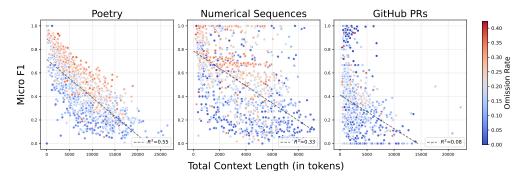


Figure 5: **GPT-4.1-mini performs worse on longer tasks in Poetry,** but the relationship is not clear in Numerical Sequences and Github PRs. Each plot shows the F1-score (y-axis) and the total context length (x-axis). Dark blue represents a lower and dark red represents a higher percentage of omission (number of omitted lines divided by total number of lines across each of three domains). Each dot on the graph represents the performance on a single instance. The dashed lines represent the least squares fit, with \mathbb{R}^2 indicating the strength of correlation between the two axes.

experience a substantial average performance drop of 56.9% under the AbsenceBench task setting, suggesting that the observed difficulty likely arises from the distinction between omissions and insertions. A possible explanation for this observation is that the Transformer attention mechanism might struggle to effectively handle "gaps" when processing documents containing omissions. We carry out additional experiments to assess this hypothesis in Section 4.2.

4 Analysis

We study the effect of perturbing the context length as well as the percentage of omissions on the models' performance. We additionally perform an error analysis and present a simple strategy that effectively addresses the difficulty in AbsenceBench. For the sake of cost, we limit our extended analysis to three LLMs: Claude-3.7-Sonnet, GPT-4.1-mini, and Llama-4-Maverick.

4.1 Perturbation Studies

We construct a perturbed dataset with a higher variance of omission rates in order to lower the impact of various factors on the difficulty of AbsenceBench tasks. We randomly sample the probability of omission to be $p \sim \mathrm{U}(0,0.5)$ for each document, and then randomly omit each element from that document with probability p. Figure 5 illustrates the relationship between the performance of GPT-4.1-mini and total context length. The analysis of Claude-3.7-Sonnet and Llama-4-Maverick are shown in Figure 6 and 7 in Appendix F.

Context length has a major impact on performance For contexts shorter than 5K tokens, the performance varies significantly with changes in context length. This finding contrasts with observations from the NIAH test, where context length affects performance only when contexts exceed 100K tokens. We perform a linear regression analysis to assess the correlation between F1-score and total context length. Results across three LLMs indicate a stronger average correlation ($R^2 = 0.55$) within the poetry domain compared to numerical sequences (0.19) and GitHub PRs (0.08).

LLMs perform worse with lower omission rate In Figure 5, we observe a notable pattern (as indicated by color) that shows a positive correlation between model's performance and the omission rate, especially in the poetry and numerical sequences domain. In other words, models actually perform *worse* when required to recall *fewer lines*. While this may initially appear counterintuitive, we hypothesize that a higher omission rate increases the likelihood of consecutive omissions occurring. Consequently, language models may be able generate text consistently until encountering a non-omitted element that disrupts their generation process. On the other hand, language models encounter greater difficulty with contexts containing fewer omissions, as numerous missing spans may interrupt the continuity between segments that must be generated.

Models	Poetry	Numerical Sequences	GitHub PRs	Average	
Placeholder: None (baseline)					
Claude-3.7-Sonnet	73.5	91.4	35.7	66.9	
GPT-4.1-mini	30.2	45.0	31.3	35.5	
Llama-4-Maverick	32.8	58.7	29.0	40.2	
Placeholder: <missing line=""></missing>					
Claude-3.7-Sonnet	$87.4_{(+18.9\%)}$	$97.7_{(+6.9\%)}$	$64.9_{(+81.8\%)}$	$83.3_{(+24.5\%)}$	
GPT-4.1-mini	$50.4_{(+66.9\%)}$	$59.2_{(+31.5\%)}$	46.8(+49.5%)	$52.1_{(+46.8\%)}$	
Llama-4-Maverick	$60.7_{(+85.0\%)}$	$78.9_{(+34.4\%)}$	$40.0_{(+37.9\%)}$	59.8 _(+48.9%)	
Placeholder:					
Claude-3.7-Sonnet	$85.5_{(+16.3\%)}$	$97.2_{(+6.3\%)}$	$61.3_{(+71.7\%)}$	$81.3_{(+21.5\%)}$	
GPT-4.1-mini	$45.9_{(+52.0\%)}$	$57.3_{(+27.3\%)}$	$36.5_{(+16.6\%)}$	$46.5_{(+31.2\%)}$	
Llama-4-Maverick	$53.9_{(+64.3\%)}$	$73.0_{(+24.4\%)}$	$35.2_{(+21.4\%)}$	$54.0_{(+34.4\%)}$	

Table 4: Micro F1-score of three LLMs evaluated on all three domains of Absence Bench using 2 different placeholders: "<missing line>" and "______" (10 consecutive underlines), along with the percentage increase in F1-score compared to the none-placeholder baseline.

4.2 Marking omissions with placeholders

To assess whether the Transformer's attention mechanism is the reason for AbsenceBench's difficulty, we explicitly mark omissions in the modified context using *placeholders*, special tokens that explicitly signal omitted segments. Specifically, we test across all domains (see Table 4 for full results) and consider two different placeholders: "<missing line>" and 10 consecutive underlines. Across all three domains, using "<missing line>" as the placeholder performs better, with an average boost in performance of 41.9%. Notably, Llama-4-Maverick achieves the highest improvement of 54.5% overall. The improved performance is comparable to that of o3-mini, yet o3-mini generates 9.1K extra thinking tokens on average.

Why does having a placeholder help so much? We hypothesize that AbsenceBench is so difficult for models because, when a segment of text is omitted, there is no position to attend to that corresponds directly to that omission. Transformer-based LLMs are incredibly good at finding the correct piece of information to attend to from a document, but what happens when the information is a conspicuous absence? What should the Transformer attend to? The fact that including a placeholder improves performance so drastically suggests that Transformer self-attention struggles to anchor on information gaps. This is a key area for future architecture and inference-time research, as omissions are common—from missing paperwork to comments deliberately left out of a recommendation letter.

5 Discussion

Our experiments highlight a failure of cutting-edge models to consistently succeed at the simple task of detecting absences, despite very strong performance on existing benchmarks. This suggests that popular evaluations do not effectively cover the capability of absence recognition, although it is fundamental for applications such as LLM-as-a-judge in which models must recognize which rubric elements are not covered. Our results are especially surprising given the simplicity of our evaluation: we only ask about surface-form absence (rather than absence on a higher semantic level), and directly provide models with both the original and modified documents for correct comparison. Indeed, our tasks could be solved by a simple program in linear time, yet deployed LLMs with hundreds of billions of parameters consistently struggle on them.

We propose an initial hypothesis explaining this behavior: identifying presence is simpler than absence with the attention mechanisms underlying Transformers [Vaswani et al., 2017]. Information included in a document can be directly attended to, while the absence of information cannot. We provide initial evidence of this hypothesis in §4.2. Rather than simply deleting elements from the original document, we include an absence "placeholder" here (e.g. "<missing line>") to provide a span of text for the transformer to attend to, indicating absence. This causes consistent, significant improvements for models, suggesting that the lack of a sequence to attend to is at least one aspect of this problem.

One important question raised by our work is how to make models more effective at handling and identifying absence. The strong performance of Gemini-2.5-flash in 2 of our 3 domains suggests that inference-time compute and reasoning mechanisms may help, although this may not hold in general. However there is a cost: we find that Gemini uses an extremely large number of reasoning tokens in many cases (Figure 3), in many cases exceeding the length of the original document by an order of magnitude. This could be the result of naive solutions, such as regenerating the full document multiple times. While this would work for the simple versions of absence studied here, it would likely fail on more complex notions of absence that look beyond missing surface forms.

We suggest a few directions where future work can go. Reasoning may improve absence detection, but must be studied on more complex, *semantic* notions of absence to be sufficiently validated—AbsenceBench is merely a necessary bar. New architectures, with mechanisms significantly different from attention, may be required to address the deeper problem. Such work should be supported by a more mechanistic understanding of why existing models often fail on these tasks, and how attention is connected to this question. We hope that AbsenceBench will provide a natural playground for such mechanistic studies. This will also be key in guiding how existing models are used—if they fail at these simple notions of absence, it is not clear that models will be trustworthy for more complex tasks that require this capability, such as model-as-a-judge evaluations using rubrics.

Regardless of how this problem might be solved, our work supports the notion that LLMs show *novel* intelligence, not well explained by analogy to humans. LLMs fail at this simple notion of absence while performing at superhuman levels on tasks that are quite similar (e.g. Needle-in-a-haystack) and many tasks that seem much more complex and challenging (e.g. math problem solving). The "shape" of LLM intelligence is simply not that similar to humans. While improving LLM performance on absence should be one goal, this also poses an opportunity to consider: how might we map the axes of variation in LLM intelligence?

Our findings suggest that models are not effective at understanding the conspicuous absence of information, an often overlooked yet foundational component of comprehension. Standard benchmarks overwhelmingly focus on whether models can identify and reason about what is present, which does not appear to tightly correlate with a model's ability to identify absent information. However use-cases such as LLM-as-a-judge, and tasks such as grading, legal reasoning, and misinformation detection, often hinge on recognizing what is missing. The gap in performance between NIAH and AbsenceBench underscores how misleading current evaluations might be if they ignore absence. Evaluators and developers should thus augment rubric-based assessments with systematic tests for absence sensitivity, both at surface and semantic levels. More broadly, absence may be a useful conceptual lens through which to understand model failure. Rather than treating hallucinations, misjudgments, or oversights as separate phenomena, they may all be related to a common limitation: models' weak grasp of what is not there. As a result, better understanding and diagnosing absence failure may reveal general-purpose principles for more robust and trustworthy LLM behavior.

6 Related Work

Long-context Language Models The context window of language models has significantly increased due to recent efforts in training long-context language models with optimized data mixture [Grattafiori et al., 2024, Gao et al., 2025, AI et al., 2025], exploring novel model architectures [Gu and Dao, 2024, Dao and Gu, 2024, Fu et al., 2023, Peng et al., 2023, Bertsch et al., 2023], developing encoding mechanisms [Su et al., 2023, Yen et al., 2024], attention algorithms that reduce memory footprints [Dao et al., 2022, Dao, 2023, Liu et al., 2023] as well as length extrapolation techniques [Press et al., 2022, Sun et al., 2022, Chen et al., 2023].

Long-context Benchmarks Many existing benchmarks evaluate language models' long-context abilities along multiple dimensions. ZeroSCROLLS [Shaham et al., 2023] focus on long-context information gathering tasks. Long-bench [Bai et al., 2024] introduces a bilingual benchmark and studies context-length variations. L-Eval [An et al., 2023] aim for a more diverse and high-quality collection of datasets. Infinite-Bench [Zhang et al., 2024] extend the context-length to over 100K tokens. HELMET [Yen et al., 2025] presents a holistic evaluation beyond synthetic data. Other previous work focus on specific tasks, such as retrieval-augmented generation [Lee et al., 2024], question answering [Kočiský et al., 2017, Wang et al., 2025], in-context learning [Agarwal et al., 2024, Anil et al., 2024, Xu et al., 2024], coreference resolution [Vodrahalli et al., 2024], and document

summarization[Chang et al., 2024, Kim et al., 2024]. Compared to those canonical tasks designed for benchmarking long-context understanding abilities, AbsenceBench features the task of handling long contexts at a finer granularity. Similar to the NIAH test [Kamradt, 2023, Yen et al., 2025] that evaluates the ability of LLMs to locate and retrieve crucial pieces of information, we formulate AbsenceBench into a much more straightforward task of identifying omissions in context, which turns out to be more challenging than identify the presence of information.

7 Limitations

While AbsenceBench exposes fundamental challenges in how current LLMs process omitted content, we note several limitations.

Surface-form omission only. We chose to keep the tasks exceedingly simple; even though our benchmark is solvable by simple heuristics models fail to use these heuristics. Our benchmark exclusively evaluates surface-form omission, where the removed content is a simple deletion. This simplification was intentional and helps make the evaluation unambiguous and trivially automatic. However, it does not capture more subtle or semantic notions of absence—such as omitted reasoning steps or missing evidence—which are critical in many real-world applications like rubric-based grading or code review.

Structured task format. AbsenceBench presents models with both the original and modified documents side-by-side and asks for the omitted elements directly. This setup may not reflect naturalistic settings where users do not provide such explicit comparisons. As a result, our evaluation may overestimate how well models would perform under more realistic situations with omitted information. Still, given the poor performance on AbsenceBench, current models are clearly not performant. However, if a new generation of models solve AbsenceBench, that will not be sufficient evidence to say that models can handle omission elegantly.

Evaluation scope. Our benchmark covers three domains (poetry, numerical sequences, and GitHub pull requests) with medium-length contexts (5K tokens). While this breadth offers early insight into the difficulty of absence detection, our findings may not generalize to other modalities (e.g., vision, audio), domains (e.g., legal, medical), or longer contexts. However, we expect these different settings to be more difficult, not less.

No prompt or instruction tuning. We only explored two different prompts and did not explore prompt tuning or in-context examples, due to the overwhelming cost of running long-context models, especially the most performant ones that make use of inference-time compute. Note that many inference-time models consistently use more the three times more "thinking tokens" than the original length of the document they are scanning for omissions. Thus, it remains an open question whether more elaborate prompting strategies, could significantly improve model performance, something we hope to examine in future work.

Statistical significance. Finally, we did not compute error bars or significance testing for our evaluation across runs due to API cost constraints. Although our findings are consistent across tasks and models, additional replication would help quantify variability and confirm robustness. We especially would have liked to run a document perturbation study, to see whether small variations in the document or prompt affect performance.

8 Conclusion

We introduce AbsenceBench, a benchmark that tests LLMs' ability to detect omitted information—an ability distinct from the well-studied task of recalling present content. Despite the benchmark's simplicity, models perform poorly, with surprising trends: fewer omissions make the task harder, and inference-time compute models often generate three times as many "thinking tokens" than the document length itself. Explicitly inserting placeholders where content is missing substantially improves performance, supporting the hypothesis that attention struggles to represent absence. These results reveal a core limitation in current models and motivate future work on absence-aware architectures and evaluation.

9 Acknowledgement

We thank OpenAI for API credits granted via their Researcher Access Program. We thank Aswathy Ajith, Todd Nief, and Chenghao Yang for their thoughtful feedback and suggestions.

References

- Rishabh Agarwal, Avi Singh, Lei M. Zhang, Bernd Bohnet, Luis Rosias, Stephanie Chan, Biao Zhang, Ankesh Anand, Zaheer Abbas, Azade Nova, John D. Co-Reyes, Eric Chu, Feryal Behbahani, Aleksandra Faust, and Hugo Larochelle. Many-shot in-context learning, 2024. URL https://arxiv.org/abs/2404.11018.
- 01. AI, :, Alex Young, Bei Chen, Chao Li, Chengen Huang, Ge Zhang, Guanwei Zhang, Guoyin Wang, Heng Li, Jiangcheng Zhu, Jianqun Chen, Jing Chang, Kaidong Yu, Peng Liu, Qiang Liu, Shawn Yue, Senbin Yang, Shiming Yang, Wen Xie, Wenhao Huang, Xiaohui Hu, Xiaoyi Ren, Xinyao Niu, Pengcheng Nie, Yanpeng Li, Yuchi Xu, Yudong Liu, Yue Wang, Yuxuan Cai, Zhenyu Gu, Zhiyuan Liu, and Zonghong Dai. Yi: Open foundation models by 01.ai, 2025. URL https://arxiv.org/abs/2403.04652.
- Chenxin An, Shansan Gong, Ming Zhong, Xingjian Zhao, Mukai Li, Jun Zhang, Lingpeng Kong, and Xipeng Qiu. L-eval: Instituting standardized evaluation for long context language models, 2023. URL https://arxiv.org/abs/2307.11088.
- Cem Anil, Esin DURMUS, Nina Rimsky, Mrinank Sharma, Joe Benton, Sandipan Kundu, Joshua Batson, Meg Tong, Jesse Mu, Daniel J Ford, Francesco Mosconi, Rajashree Agrawal, Naomi Bashkansky Rylan Schaeffer, Samuel Svenningsen, Mike Lambert, Ansh Radhakrishnan, Carson Denison, Evan J Hubinger, Yuntao Bai, Trenton Bricken, Timothy Maxwell, Nicholas Schiefer, James Sullyand Alex Tamkin, Tamera Lanham, Karina Nguyen, Tomasz Korbak, Jared Kaplan, Deep Ganguli, Samuel R. Bowman, Ethan Perez, Roger Baker Grosse, and David Duvenaud. Many-shot jailbreaking, 2024. URL https://www.anthropic.com/research/many-shot-jailbreaking.
- Anthropic. Claude 3 haiku: our fastest model yet, 2024. URL https://www.anthropic.com/news/claude-3-haiku.
- 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. Longbench: A bilingual, multitask benchmark for long context understanding, 2024. URL https://arxiv.org/abs/2308.14508.
- Amanda Bertsch, Uri Alon, Graham Neubig, and Matthew R. Gormley. Unlimiformer: Long-range transformers with unlimited length input, 2023. URL https://arxiv.org/abs/2305.01625.
- Yapei Chang, Kyle Lo, Tanya Goyal, and Mohit Iyyer. Booookscore: A systematic exploration of book-length summarization in the era of llms, 2024. URL https://arxiv.org/abs/2310.00785.
- Shouyuan Chen, Sherman Wong, Liangjian Chen, and Yuandong Tian. Extending context window of large language models via positional interpolation, 2023. URL https://arxiv.org/abs/2306.15595.
- Tri Dao. Flashattention-2: Faster attention with better parallelism and work partitioning, 2023. URL https://arxiv.org/abs/2307.08691.
- Tri Dao and Albert Gu. Transformers are ssms: Generalized models and efficient algorithms through structured state space duality, 2024. URL https://arxiv.org/abs/2405.21060.
- Tri Dao, Daniel Y. Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. Flashattention: Fast and memory-efficient exact attention with io-awareness, 2022. URL https://arxiv.org/abs/2205.14135.

DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L. Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Meng Li, Miaojun Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen, Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, R. J. Chen, R. L. Jin, Ruyi Chen, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng Zhou, Shuting Pan, S. S. Li, Shuang Zhou, Shaoqing Wu, Shengfeng Ye, Tao Yun, Tian Pei, Tianyu Sun, T. Wang, Wangding Zeng, Wanjia Zhao, Wen Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, W. L. Xiao, Wei An, Xiaodong Liu, Xiaohan Wang, Xiaokang Chen, Xiaotao Nie, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin, X. Q. Li, Xiangyue Jin, Xiaojin Shen, Xiaosha Chen, Xiaowen Sun, Xiaoxiang Wang, Xinnan Song, Xinyi Zhou, Xianzu Wang, Xinxia Shan, Y. K. Li, Y. Q. Wang, Y. X. Wei, Yang Zhang, Yanhong Xu, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Wang, Yi Yu, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yuan Ou, Yuduan Wang, Yue Gong, Yuheng Zou, Yujia He, Yunfan Xiong, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Y. X. Zhu, Yanhong Xu, Yanping Huang, Yaohui Li, Yi Zheng, Yuchen Zhu, Yunxian Ma, Ying Tang, Yukun Zha, Yuting Yan, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhenda Xie, Zhengyan Zhang, Zhewen Hao, Zhicheng Ma, Zhigang Yan, Zhiyu Wu, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Zizheng Pan, Zhen Huang, Zhipeng Xu, Zhongyu Zhang, and Zhen Zhang. Deepseek-r1: Incentivizing reasoning capability in Ilms via reinforcement learning, 2025. URL https://arxiv.org/abs/2501.12948.

Simon Frieder, Luca Pinchetti, Alexis Chevalier, Ryan-Rhys Griffiths, Tommaso Salvatori, Thomas Lukasiewicz, Philipp Christian Petersen, and Julius Berner. Mathematical capabilities of chatgpt, 2023. URL https://arxiv.org/abs/2301.13867.

Daniel Y. Fu, Elliot L. Epstein, Eric Nguyen, Armin W. Thomas, Michael Zhang, Tri Dao, Atri Rudra, and Christopher Ré. Simple hardware-efficient long convolutions for sequence modeling, 2023. URL https://arxiv.org/abs/2302.06646.

Tianyu Gao, Alexander Wettig, Howard Yen, and Danqi Chen. How to train long-context language models (effectively), 2025. URL https://arxiv.org/abs/2410.02660.

Google. Gemini 2.5: Our most intelligent ai model, 2025. URL https://blog.google/technology/google-deepmind/gemini-model-thinking-updates-march-2025/.

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, Danny Wyatt, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Francisco Guzmán, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Govind Thattai, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jack Zhang, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua

Saxe, Junteng Jia, Kalyan Vasuden Alwala, Karthik Prasad, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Kushal Lakhotia, Lauren Rantala-Yeary, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Maria Tsimpoukelli, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Ning Zhang, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohan Maheswari, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Raparthy, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vítor Albiero, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaofang Wang, Xiaoqing Ellen Tan, Xide Xia, Xinfeng Xie, Xuchao Jia, Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aayushi Srivastava, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Amos Teo, Anam Yunus, Andrei Lupu, Andrew Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Dong, Annie Franco, Anuj Goyal, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Ce Liu, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Cynthia Gao, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Eric-Tuan Le, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Filippos Kokkinos, Firat Ozgenel, Francesco Caggioni, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hakan Inan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Hongyuan Zhan, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Ilias Leontiadis, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Janice Lam, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kiran Jagadeesh, Kun Huang, Kunal Chawla, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabsa, Manav Avalani, Manish Bhatt, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Miao Liu, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikhil Mehta, Nikolay Pavlovich Laptey, Ning Dong, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy,

- Raghu Nayani, Rahul Mitra, Rangaprabhu Parthasarathy, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Russ Howes, Ruty Rinott, Sachin Mehta, Sachin Siby, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Mahajan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shishir Patil, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Summer Deng, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Koehler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaojian Wu, Xiaolan Wang, Xilun Wu, Xinbo Gao, Yaniv Kleinman, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yu Zhao, Yuchen Hao, Yundi Qian, Yunlu Li, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, Zhiwei Zhao, and Zhiyu Ma. The llama 3 herd of models, 2024. URL https://arxiv.org/abs/2407.21783.
- Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces, 2024. URL https://arxiv.org/abs/2312.00752.
- Cheng-Ping Hsieh, Simeng Sun, Samuel Kriman, Shantanu Acharya, Dima Rekesh, Fei Jia, Yang Zhang, and Boris Ginsburg. Ruler: What's the real context size of your long-context language models?, 2024. URL https://arxiv.org/abs/2404.06654.
- Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, Gianna Lengyel, Guillaume Bour, Guillaume Lample, Lélio Renard Lavaud, Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le Scao, Théophile Gervet, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. Mixtral of experts, 2024. URL https://arxiv.org/abs/2401.04088.
- Gregory Kamradt. Needle in a haystack pressure testing llms, 2023. URL https://github.com/gkamradt/LLMTest_NeedleInAHaystack.
- Yekyung Kim, Yapei Chang, Marzena Karpinska, Aparna Garimella, Varun Manjunatha, Kyle Lo, Tanya Goyal, and Mohit Iyyer. Fables: Evaluating faithfulness and content selection in book-length summarization, 2024. URL https://arxiv.org/abs/2404.01261.
- Ryan Koo, Minhwa Lee, Vipul Raheja, Jong Inn Park, Zae Myung Kim, and Dongyeop Kang. Benchmarking cognitive biases in large language models as evaluators, 2024. URL https://arxiv.org/abs/2309.17012.
- Tomáš Kočiský, Jonathan Schwarz, Phil Blunsom, Chris Dyer, Karl Moritz Hermann, Gábor Melis, and Edward Grefenstette. The narrativeqa reading comprehension challenge, 2017. URL https://arxiv.org/abs/1712.07040.
- Jinhyuk Lee, Anthony Chen, Zhuyun Dai, Dheeru Dua, Devendra Singh Sachan, Michael Boratko, Yi Luan, Sébastien M. R. Arnold, Vincent Perot, Siddharth Dalmia, Hexiang Hu, Xudong Lin, Panupong Pasupat, Aida Amini, Jeremy R. Cole, Sebastian Riedel, Iftekhar Naim, Ming-Wei Chang, and Kelvin Guu. Can long-context language models subsume retrieval, rag, sql, and more?, 2024. URL https://arxiv.org/abs/2406.13121.
- Hao Liu, Matei Zaharia, and Pieter Abbeel. Ring attention with blockwise transformers for near-infinite context, 2023. URL https://arxiv.org/abs/2310.01889.
- Meta. The llama 4 herd: The beginning of a new era of natively multimodal ai innovation, 2025. URL https://ai.meta.com/blog/llama-4-multimodal-intelligence/.
- OpenAI. Openai o3-mini, pushing the frontier of cost-effective reasoning., 2025. URL https://openai.com/index/openai-o3-mini/.

OpenAI, :, Aaron Hurst, Adam Lerer, Adam P. Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, Aleksander Madry, Alex Baker-Whitcomb, Alex Beutel, Alex Borzunov, Alex Carney, Alex Chow, Alex Kirillov, Alex Nichol, Alex Paino, Alex Renzin, Alex Tachard Passos, Alexander Kirillov, Alexi Christakis, Alexis Conneau, Ali Kamali, Allan Jabri, Allison Moyer, Allison Tam, Amadou Crookes, Amin Tootoochian, Amin Tootoonchian, Ananya Kumar, Andrea Vallone, Andrej Karpathy, Andrew Braunstein, Andrew Cann, Andrew Codispoti, Andrew Galu, Andrew Kondrich, Andrew Tulloch, Andrey Mishchenko, Angela Baek, Angela Jiang, Antoine Pelisse, Antonia Woodford, Anuj Gosalia, Arka Dhar, Ashley Pantuliano, Avi Nayak, Avital Oliver, Barret Zoph, Behrooz Ghorbani, Ben Leimberger, Ben Rossen, Ben Sokolowsky, Ben Wang, Benjamin Zweig, Beth Hoover, Blake Samic, Bob McGrew, Bobby Spero, Bogo Giertler, Bowen Cheng, Brad Lightcap, Brandon Walkin, Brendan Quinn, Brian Guarraci, Brian Hsu, Bright Kellogg, Brydon Eastman, Camillo Lugaresi, Carroll Wainwright, Cary Bassin, Cary Hudson, Casey Chu, Chad Nelson, Chak Li, Chan Jun Shern, Channing Conger, Charlotte Barette, Chelsea Voss, Chen Ding, Cheng Lu, Chong Zhang, Chris Beaumont, Chris Hallacy, Chris Koch, Christian Gibson, Christina Kim, Christine Choi, Christine McLeavey, Christopher Hesse, Claudia Fischer, Clemens Winter, Coley Czarnecki, Colin Jarvis, Colin Wei, Constantin Koumouzelis, Dane Sherburn, Daniel Kappler, Daniel Levin, Daniel Levy, David Carr, David Farhi, David Mely, David Robinson, David Sasaki, Denny Jin, Dev Valladares, Dimitris Tsipras, Doug Li, Duc Phong Nguyen, Duncan Findlay, Edede Oiwoh, Edmund Wong, Ehsan Asdar, Elizabeth Proehl, Elizabeth Yang, Eric Antonow, Eric Kramer, Eric Peterson, Eric Sigler, Eric Wallace, Eugene Brevdo, Evan Mays, Farzad Khorasani, Felipe Petroski Such, Filippo Raso, Francis Zhang, Fred von Lohmann, Freddie Sulit, Gabriel Goh, Gene Oden, Geoff Salmon, Giulio Starace, Greg Brockman, Hadi Salman, Haiming Bao, Haitang Hu, Hannah Wong, Haoyu Wang, Heather Schmidt, Heather Whitney, Heewoo Jun, Hendrik Kirchner, Henrique Ponde de Oliveira Pinto, Hongyu Ren, Huiwen Chang, Hyung Won Chung, Ian Kivlichan, Ian O'Connell, Ian O'Connell, Ian Osband, Ian Silber, Ian Sohl, Ibrahim Okuyucu, Ikai Lan, Ilya Kostrikov, Ilya Sutskever, Ingmar Kanitscheider, Ishaan Gulrajani, Jacob Coxon, Jacob Menick, Jakub Pachocki, James Aung, James Betker, James Crooks, James Lennon, Jamie Kiros, Jan Leike, Jane Park, Jason Kwon, Jason Phang, Jason Teplitz, Jason Wei, Jason Wolfe, Jay Chen, Jeff Harris, Jenia Varavva, Jessica Gan Lee, Jessica Shieh, Ji Lin, Jiahui Yu, Jiayi Weng, Jie Tang, Jieqi Yu, Joanne Jang, Joaquin Quinonero Candela, Joe Beutler, Joe Landers, Joel Parish, Johannes Heidecke, John Schulman, Jonathan Lachman, Jonathan McKay, Jonathan Uesato, Jonathan Ward, Jong Wook Kim, Joost Huizinga, Jordan Sitkin, Jos Kraaijeveld, Josh Gross, Josh Kaplan, Josh Snyder, Joshua Achiam, Joy Jiao, Joyce Lee, Juntang Zhuang, Justyn Harriman, Kai Fricke, Kai Hayashi, Karan Singhal, Katy Shi, Kavin Karthik, Kayla Wood, Kendra Rimbach, Kenny Hsu, Kenny Nguyen, Keren Gu-Lemberg, Kevin Button, Kevin Liu, Kiel Howe, Krithika Muthukumar, Kyle Luther, Lama Ahmad, Larry Kai, Lauren Itow, Lauren Workman, Leher Pathak, Leo Chen, Li Jing, Lia Guy, Liam Fedus, Liang Zhou, Lien Mamitsuka, Lilian Weng, Lindsay McCallum, Lindsey Held, Long Ouyang, Louis Feuvrier, Lu Zhang, Lukas Kondraciuk, Lukasz Kaiser, Luke Hewitt, Luke Metz, Lyric Doshi, Mada Aflak, Maddie Simens, Madelaine Boyd, Madeleine Thompson, Marat Dukhan, Mark Chen, Mark Gray, Mark Hudnall, Marvin Zhang, Marwan Aljubeh, Mateusz Litwin, Matthew Zeng, Max Johnson, Maya Shetty, Mayank Gupta, Meghan Shah, Mehmet Yatbaz, Meng Jia Yang, Mengchao Zhong, Mia Glaese, Mianna Chen, Michael Janner, Michael Lampe, Michael Petrov, Michael Wu, Michele Wang, Michelle Fradin, Michelle Pokrass, Miguel Castro, Miguel Oom Temudo de Castro, Mikhail Pavlov, Miles Brundage, Miles Wang, Minal Khan, Mira Murati, Mo Bavarian, Molly Lin, Murat Yesildal, Nacho Soto, Natalia Gimelshein, Natalie Cone, Natalie Staudacher, Natalie Summers, Natan LaFontaine, Neil Chowdhury, Nick Ryder, Nick Stathas, Nick Turley, Nik Tezak, Niko Felix, Nithanth Kudige, Nitish Keskar, Noah Deutsch, Noel Bundick, Nora Puckett, Ofir Nachum, Ola Okelola, Oleg Boiko, Oleg Murk, Oliver Jaffe, Olivia Watkins, Olivier Godement, Owen Campbell-Moore, Patrick Chao, Paul McMillan, Pavel Belov, Peng Su, Peter Bak, Peter Bakkum, Peter Deng, Peter Dolan, Peter Hoeschele, Peter Welinder, Phil Tillet, Philip Pronin, Philippe Tillet, Prafulla Dhariwal, Qiming Yuan, Rachel Dias, Rachel Lim, Rahul Arora, Rajan Troll, Randall Lin, Rapha Gontijo Lopes, Raul Puri, Reah Miyara, Reimar Leike, Renaud Gaubert, Reza Zamani, Ricky Wang, Rob Donnelly, Rob Honsby, Rocky Smith, Rohan Sahai, Rohit Ramchandani, Romain Huet, Rory Carmichael, Rowan Zellers, Roy Chen, Ruby Chen, Ruslan Nigmatullin, Ryan Cheu, Saachi Jain, Sam Altman, Sam Schoenholz, Sam Toizer, Samuel Miserendino, Sandhini Agarwal, Sara Culver, Scott Ethersmith, Scott Gray, Sean Grove, Sean Metzger, Shamez Hermani, Shantanu Jain, Shengjia Zhao, Sherwin Wu, Shino Jomoto, Shirong Wu, Shuaiqi, Xia, Sonia Phene, Spencer

Papay, Srinivas Narayanan, Steve Coffey, Steve Lee, Stewart Hall, Suchir Balaji, Tal Broda, Tal Stramer, Tao Xu, Tarun Gogineni, Taya Christianson, Ted Sanders, Tejal Patwardhan, Thomas Cunninghman, Thomas Degry, Thomas Dimson, Thomas Raoux, Thomas Shadwell, Tianhao Zheng, Todd Underwood, Todor Markov, Toki Sherbakov, Tom Rubin, Tom Stasi, Tomer Kaftan, Tristan Heywood, Troy Peterson, Tyce Walters, Tyna Eloundou, Valerie Qi, Veit Moeller, Vinnie Monaco, Vishal Kuo, Vlad Fomenko, Wayne Chang, Weiyi Zheng, Wenda Zhou, Wesam Manassra, Will Sheu, Wojciech Zaremba, Yash Patil, Yilei Qian, Yongjik Kim, Youlong Cheng, Yu Zhang, Yuchen He, Yuchen Zhang, Yujia Jin, Yunxing Dai, and Yury Malkov. Gpt-4o system card, 2024a. URL https://arxiv.org/abs/2410.21276.

OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O'Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. Gpt-4 technical report, 2024b. URL https://arxiv.org/abs/2303.08774.

Allison Parrish. gutenberg-poetry-corpus: A corpus of poetry from project gutenberg, 2018.

- Bo Peng, Eric Alcaide, Quentin Anthony, Alon Albalak, Samuel Arcadinho, Stella Biderman, Huanqi Cao, Xin Cheng, Michael Chung, Matteo Grella, Kranthi Kiran GV, Xuzheng He, Haowen Hou, Jiaju Lin, Przemyslaw Kazienko, Jan Kocon, Jiaming Kong, Bartlomiej Koptyra, Hayden Lau, Krishna Sri Ipsit Mantri, Ferdinand Mom, Atsushi Saito, Guangyu Song, Xiangru Tang, Bolun Wang, Johan S. Wind, Stanislaw Wozniak, Ruichong Zhang, Zhenyuan Zhang, Qihang Zhao, Peng Zhou, Qinghua Zhou, Jian Zhu, and Rui-Jie Zhu. Rwkv: Reinventing rnns for the transformer era, 2023. URL https://arxiv.org/abs/2305.13048.
- Ofir Press, Noah A. Smith, and Mike Lewis. Train short, test long: Attention with linear biases enables input length extrapolation, 2022. URL https://arxiv.org/abs/2108.12409.
- Qwen. Qwen3 technical report, 2025a. URL https://github.com/QwenLM/Qwen3/blob/main/Qwen3_Technical_Report.pdf.
- Qwen. Qwq-32b: Embracing the power of reinforcement learning, 2025b. URL https://qwenlm.github.io/blog/qwq-32b/.
- Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Romain Sauvestre, Tal Remez, Jérémy Rapin, Artyom Kozhevnikov, Ivan Evtimov, Joanna Bitton, Manish Bhatt, Cristian Canton Ferrer, Aaron Grattafiori, Wenhan Xiong, Alexandre Défossez, Jade Copet, Faisal Azhar, Hugo Touvron, Louis Martin, Nicolas Usunier, Thomas Scialom, and Gabriel Synnaeve. Code llama: Open foundation models for code, 2024. URL https://arxiv.org/abs/2308.12950.
- Uri Shaham, Maor Ivgi, Avia Efrat, Jonathan Berant, and Omer Levy. ZeroSCROLLS: A zero-shot benchmark for long text understanding. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 7977–7989, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.findings-emnlp.536. URL https://aclanthology.org/2023.findings-emnlp.536/.
- Jianlin Su, Yu Lu, Shengfeng Pan, Ahmed Murtadha, Bo Wen, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding, 2023. URL https://arxiv.org/abs/2104.09864.
- Yutao Sun, Li Dong, Barun Patra, Shuming Ma, Shaohan Huang, Alon Benhaim, Vishrav Chaudhary, Xia Song, and Furu Wei. A length-extrapolatable transformer, 2022. URL https://arxiv.org/abs/2212.10554.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- Kiran Vodrahalli, Santiago Ontanon, Nilesh Tripuraneni, Kelvin Xu, Sanil Jain, Rakesh Shivanna, Jeffrey Hui, Nishanth Dikkala, Mehran Kazemi, Bahare Fatemi, Rohan Anil, Ethan Dyer, Siamak Shakeri, Roopali Vij, Harsh Mehta, Vinay Ramasesh, Quoc Le, Ed Chi, Yifeng Lu, Orhan Firat, Angeliki Lazaridou, Jean-Baptiste Lespiau, Nithya Attaluri, and Kate Olszewska. Michelangelo: Long context evaluations beyond haystacks via latent structure queries, 2024. URL https://arxiv.org/abs/2409.12640.
- Cunxiang Wang, Ruoxi Ning, Boqi Pan, Tonghui Wu, Qipeng Guo, Cheng Deng, Guangsheng Bao, Xiangkun Hu, Zheng Zhang, Qian Wang, and Yue Zhang. Novelqa: Benchmarking question answering on documents exceeding 200k tokens, 2025. URL https://arxiv.org/abs/2403.12766.
- xAI. Grok 3 beta the age of reasoning agents, 2025. URL https://x.ai/news/grok-3.
- Xiaoyue Xu, Qinyuan Ye, and Xiang Ren. Stress-testing long-context language models with lifelong icl and task haystack, 2024. URL https://arxiv.org/abs/2407.16695.
- An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang,

- Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. Qwen2.5 technical report. *arXiv preprint arXiv:2412.15115*, 2024.
- Howard Yen, Tianyu Gao, and Danqi Chen. Long-context language modeling with parallel context encoding, 2024. URL https://arxiv.org/abs/2402.16617.
- Howard Yen, Tianyu Gao, Minmin Hou, Ke Ding, Daniel Fleischer, Peter Izsak, Moshe Wasserblat, and Danqi Chen. Helmet: How to evaluate long-context language models effectively and thoroughly, 2025. URL https://arxiv.org/abs/2410.02694.
- 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. ∞bench: Extending long context evaluation beyond 100k tokens, 2024. URL https://arxiv.org/abs/2402.13718.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. Judging Ilm-as-a-judge with mt-bench and chatbot arena, 2023. URL https://arxiv.org/abs/2306.05685.

A Prompt Templates

We show the detailed prompt templates for language model evaluation across each domain included in AbsenceBench in Table 6 that are used to obtain the evaluation results in Table 3. Additionally, Table 7 presents prompt templates with adapted instructions used for assessing the NIAH test under the poetry and GitHub PRs domains.

```
Poetry (default)

System Prompt
You are helping a student practice memorizing poems. The student will recite a poem, but they may have missed some lines. Your task is to identify exactly which lines are missing from their recitation.
List only the missing lines, nothing else.
User Message
Here is the complete original poem:
{original poem}
Now, here is my recitation which may be missing some lines:
{modified poem}
What lines did I miss? Please list only the missing lines, nothing else.
```

Numerical Sequences (default)

```
You are helping a student practice reciting sequences. The student will recite a sequence, but they may have missed some numbers. Your task is to identify exactly which numbers are missing from their recitation.

List only the missing numbers, nothing else

User Message

Here is a sequence of numbers:

{original sequence}

Now, here is my recitation of the sequence which may be missing some numbers:

{modified sequence}

What numbers did I miss? Please list only the missing numbers, nothing else.
```

GitHub PRs (default)

```
System Prompt
You are helping a software developer determine if their merge of a pull request was successful. The developer had to edit the commit history and just wants to make sure that they have not changed what will be merged. They will list the changed lines. Your job is to figure out if they have missed any insertions or deletions from the original merge. Only pay attention to the insertions and deletions (ignore the context of the diff).

User Message
Here is the complete original diff:
{original diff}
And here is the merge diff after the developer fixed the commit history:
{modified diff}
What changed lines (insertions or deletions) present in the original diff are missing in the merge diff (if any)?
List only the missing changed lines, nothing else.
```

Table 5: Detailed prompt templates used for evaluating language models under each domain of AbsenceBench

Poetry (post-instruction)

System Prompt

You are helping a student practice memorizing poems. The student will recite a poem, but they may have missed some lines. Your task is to identify exactly which lines are missing from their recitation. List only the missing lines, nothing else.

User Message

{original poem}

{modified poem}

The original poem is followed by my recitation of the poem which may be missing some lines.
What lines did I miss? Please list only the missing lines, nothing else.

Numerical Sequences (post-instruction)

System Prompt

You are helping a student practice reciting sequences. The student will recite a sequence, but they may have missed some numbers. Your task is to identify exactly which numbers are missing from their recitation. List only the missing numbers, nothing else

User Message

{original sequence}

{modified sequence}

The original sequence is followed by my recitation of the sequence which may be missing some numbers.

What numbers did I miss? Please list only the missing numbers, nothing else.

GitHub PRs (post-instruction)

System Prompt

You are helping a software developer determine if their merge of a pull request was successful. The developer had to edit the commit history and just wants to make sure that they have not changed what will be merged. They will list the changed lines. Your job is to figure out if they have missed any insertions or deletions from the original merge. Only pay attention to the insertions and deletions (ignore the context of the diff).

User Message

{original diff}

{modified diff}

The original diff is followed by a merge diff after the developer fixed the commit history $\$

What changed lines (insertions or deletions) present in the original diff are missing in the merge diff (if any)?
List only the missing changed lines, nothing else.

Table 6: Perturbed prompt templates where the task instructions are positioned after the original and modified documents.

Poetry (NIAH)

System Prompt

You are helping a student practice memorizing poems. The student will recite a poem, but they may have added some random lines that related to Harry Potter characters. Your task is to identify exactly which lines are not in the original poem. List only the missing lines, nothing else.

User Message

Here is the complete original poem:

{original poem}

Now, here is my recitation with some extra lines that is related to Harry Potter novel series:

{modified poem}

What lines did I miss? Please list only the extra lines, nothing else.

GitHub PRs (NIAH)

System Prompt

You are helping a software developer determine if their merge of a pull request was successful. The developer had to edit the commit history and accidently added some random lines related to Harry Potter characters. They will list the changed lines. Your job is to figure out if they have added any insertions from the original merge. Only pay attention to the insertions.

User Message

Here is the complete original diff:

{original diff}

And here is the merge diff after the developer fixed the commit history: $\{ modified \ diff \}$

What changed lines (insertions or deletions) present in the original diff are missing in the merge diff (if any)?
List only the missing changed lines, nothing else.

Table 7: Modified prompt templates used for evaluating language models under the NIAH test setting within the Poetry and GitHub PRs domains.

B Inference Details

We evaluate a total of 14 LLMs on AbsenceBench. All model inferences are obtained through API requests. We show the detailed information of each model's maximum context length, API provider and reference in Table 8. We fix the generation temperature at the default value of 1.0 across all API providers. For models with inference-time compute, we adopt the default thinking budget set by these API providers: Google Gemini (2.5-flash; dynamic with an upper bound of 24576 tokens), OpenAI (03-mini; medium reasoning effort), TogetherAI (Qwen3 models and Deepseek R1). We also manually set thinking budget for Claude-3.7-Sonnet to be 10K tokens and Grok-3-mini-beta to be "low reasoning effort" due to cost considerations.

Models	Context Length	API Provider	API Reference
Gemini-2.5-flash [Google, 2025]	1M	Google	models/gemini-2.5-flash-preview-04-17
Claude-3.7-Sonnet [Anthropic, 2024]	200K	Anthropic	claude-3-7-sonnet-latest
o3-mini [OpenAI, 2025]	200K	OpenAI	o3-mini
GPT-4.1 [OpenAI et al., 2024b]	1M	OpenAI	gpt-4.1
GPT-4.1-mini [OpenAI et al., 2024b]	1M	OpenAI	gpt-4.1-mini
GPT-4o [OpenAI et al., 2024a]	128K	OpenAI	gpt-4o
Grok-3-mini-Beta [xAI, 2025]	131K	xAI	grok-3-mini-beta
Llama-4-Maverick [Meta, 2025]	1M	TogetherAI	meta-llama/Llama-4-Maverick-17B-128E-Instruct-FP8
Llama-3.3-70B-Instruct [Grattafiori et al., 2024]	131K	TogetherAI	meta-llama/Llama-3.3-70B-Instruct-Turbo
Llama-3.1-8B-Instruct [Grattafiori et al., 2024]	128K	TogetherAI	meta-llama/Meta-Llama-3.1-8B-Instruct-Turbo
Qwen3-235B [Qwen, 2025a]	41K	TogetherAI	Qwen/Qwen3-235B-A22B-fp8-tput
Qwen2.5-72B-Instruct [Yang et al., 2024]	32K	TogetherAI	Qwen/Qwen2.5-72B-Instruct-Turbo
QwQ-32B [Qwen, 2025b]	32K	TogetherAI	Qwen/QwQ-32B
DeepSeek-R1 [DeepSeek-AI et al., 2025]	128K	TogetherAI	deepseek-ai/DeepSeek-R1
Mixtral-8x7B-Instruct [Jiang et al., 2024]	32K	TogetherAI	mistralai/Mixtral-8x7B-Instruct-v0.1

Table 8: Detailed information of models evaluated on AbsenceBench

C Alternative Prompt Template

Alongside the default prompt template (see Table 2), we curate a perturbed template where the task instructions are positioned after the original and modified documents. We refer to this perturbed prompt template as "post-instruction". Table 9 shows that the performance of three LLMs prompted with both prompt templates. All three models exhibit reduced performance with the post-instruction template under both the poetry and numerical sequences domains, suggesting that LLMs more effectively detect omissions when instructions are presented beforehand. Meanwhile, LLMs remain robust in the GitHub PRs domain. All detailed prompt templates are included in Appendix A.

Models	Poetry		Numerical Sequences		GitHub PRs	
	default	post-instruction	default	post-instruction	default	post-instruction
Claude-3.7-Sonnet GPT-4.1-mini	78.8 45.4	66.1 39.2	94.3	93.7 42.5	35.1	34.7 35.4
Llama-4-Maverick	48.7	36.6	71.5	63.6	29.6	30.5

Table 9: Micro F1-score (%) of three LLMs evaluated on the AbsenceBench with two different prompt templates. All three models exhibit reduced performance with the *post-instruction* template in the poetry and numerical sequences domains, but remain robust in the GitHub PRs domain.

D Dataset Construction Parameters

Our main considerations for the selection of dataset construction parameters were cost, balancing difficulty, and ensuring a valid evaluation, all described below:

- Cost and Practicality: our poetry source data includes poems of widely varying lengths, with some extending beyond 5000 lines (145 out of 1187). Without truncation, the benchmark would be technically not applicable for LLMs with a shorter context-length limit.
- Utility and Difficulty Balance: according to our perturbation studies in Section 4, further extending the context length would add to the difficulty of AbsenceBench, but is not necessary since the benchmark is already highly challenging to leading LLMs at these medium context lengths (see Table 2). We believe balancing the utility and difficulty of the benchmark is important, leading us to the current medium-context.
- Validity: the selected length range ensures that instances are neither trivially short nor excessively
 long, thus avoiding scenarios where context length would become the dominant confounder in
 model accuracy. On the other hand, the benchmark only focuses on surface-level omissions.
 Truncating the instances would not bring any semantic concern to the validity of the benchmark.

E Comparison to the NIAH test

We extend the analysis of Section 3.2 on the comparison between the task design of AbsenceBench and the NIAH test. In table 10, we show 3 LLMs achieve near perfect results in the NIAH test under the poetry domain, as well as substantially improved results within the GitHub PRs domain.

Models	Poetry	Poetry (NIAH)	GitHub PRs	GitHub PRs (NIAH)
Claude-3.7-Sonnet	73.5	99.52	35.7	97.1
GPT-4.1-mini	30.2	99.5	31.3	92.1
Llama-4-Maverick	32.8	99.47	29.0	86.2

Table 10: Micro F1-score (%) of three LLMs evaluated on the AbsenceBench and the NIAH test setting under the Poetry and GitHub PRs domains.

F Additional Results

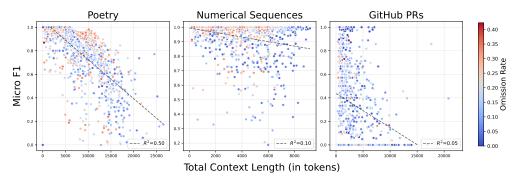


Figure 6: Micro-F1 score of **Claude-3.7-Sonnet** (y-axis) as a function of the total context length (x-axis) as well as the percentage of omission (color)

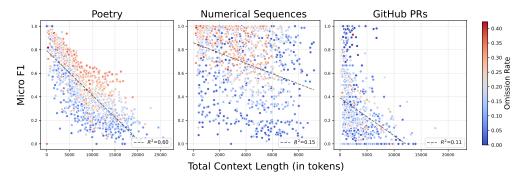


Figure 7: Micro-F1 score of **Llama-4-Maverick** (y-axis) as a function of the total context length (x-axis) as well as the percentage of omission (color)

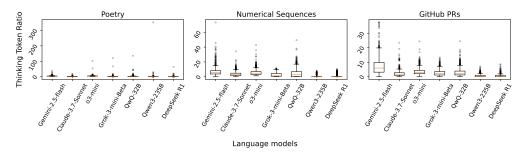


Figure 8: The thinking token ratio (number of generated thinking tokens divided by number of tokens in the original context) for four inference-time compute models under each domain.

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: We do not over-claim our key points and insights, which are all justified and supported by experiment results.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the
 contributions made in the paper and important assumptions and limitations. A No or
 NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals
 are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: We include a dedicated Limitation Section 7 to discuss the limitations of the work.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory Assumptions and Proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [NA]

Justification: Does not apply. The paper does not contain theoretical proofs. All results and experiments are empirical.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and crossreferenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: We fully disclose all the details needed to reproduce the main results.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
- (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
- (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
- (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
- (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: We release the data, which is already hosted, and our code to be submitted as supplemental material.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: We specify all details to reproduce the results.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental
 material.

7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [NA]

Justification: Due to the cost constraints of performing experiments (through API), we are not able to perform extensive studies, but we mentioned it as one of the limitations.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).

- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error
 of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments Compute Resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: We provide full details on experiment resources.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code Of Ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?

Answer: [Yes]

Justification: The research conducted in the paper is conform with the NeurIPS Code of Ethics

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a
 deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader Impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: Our benchmark allows us to measure a necessary (but insufficient) task to ensure that LLMs will not miss important omissions when deployed in potentially sensitive situations. However, such technologies are inevitably dual use—LLMs that are better at spotting omissions are also better at surveillance. Since our benchmark is mainly a tool for analysis, we believe will have an overall positive impact, by creating common knowledge about the capabilities of LLMs, both closed and open source.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [Yes]

Justification: We ensure that our benchmark will not pose any threats of being misused and unsafe.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with
 necessary safeguards to allow for controlled use of the model, for example by requiring
 that users adhere to usage guidelines or restrictions to access the model or implementing
 safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: All sources of assets are properly credited and cited.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.

- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New Assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [Yes]

Justification: The new assets are well-documented

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and Research with Human Subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: We do not crowdsource in this paper.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects Guidelines:

 The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.

- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.