

Same Task, More Tokens: the Impact of Input Length on the Reasoning Performance of Large Language Models

Mosh Levy*¹ Alon Jacoby*¹ Yoav Goldberg^{1,2}

¹Bar-Ilan University ²Allen Institute for AI

{moshe0110, alonj4}@gmail.com

Abstract

This paper explores the impact of extending input lengths on the capabilities of Large Language Models (LLMs). Despite LLMs advancements in recent times, their performance consistency across different input lengths is not well understood. We investigate this aspect by introducing a novel QA reasoning framework, specifically designed to assess the impact of input length. We isolate the effect of input length using multiple versions of the same sample, each being extended with padding of different lengths, types and locations. Our findings show a notable degradation in LLMs' reasoning performance at much shorter input lengths than their technical maximum. We show that the degradation trend appears in every version of our dataset, although at different intensities. Additionally, our study reveals that traditional perplexity metrics do not correlate with performance of LLMs' in long input reasoning tasks. We analyse our results and identify failure modes that can serve as useful guides for future research, potentially informing strategies to address the limitations observed in LLMs.

1 Introduction

Recent advancements in Large Language Models (LLMs) show impressive performance across a range of tasks (OpenAI, 2023; Anil et al., 2023; Jiang et al., 2024), including answering correctly complex questions requiring multiple reasoning steps (Kojima et al., 2022; Wei et al., 2022). These models also claim to support increasingly longer inputs. This development underscores the need to examine their performance on the longer inputs they are now technically supporting.

A reasonable assumption is that support for long inputs would transfer across tasks and enable a model adept at solving a task when presented in a short input prompt, to perform the same task when

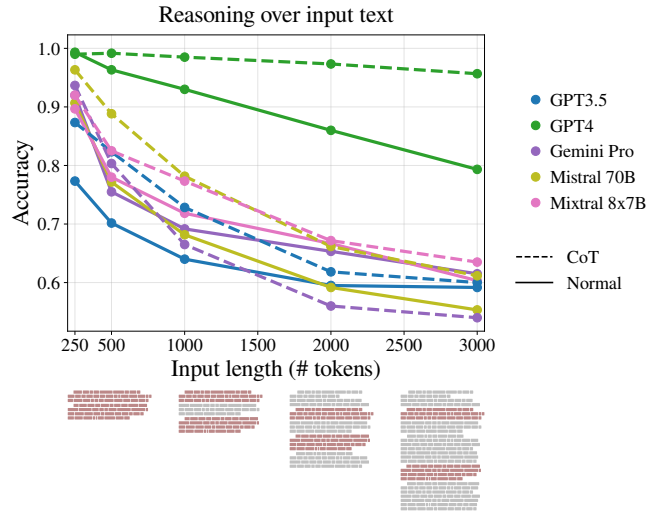


Figure 1: Reasoning performance drops as input grows, across a variety of tasks. Inputs are composed of text containing information relevant to the task (in red), and irrelevant text (grey) which is drawn from various sources and extended incrementally. Two separate text spans are required to answer correctly, and are located randomly in the input. Each point reflects the performance across 600 samples.

it is embedded within a longer prompt. Does this assumption hold? Recent studies that benchmark models over tasks that involve longer inputs, including reasoning tasks, indicate that indeed models often struggle with reasoning over long inputs (Shaham et al., 2023; Li et al., 2023; Bai et al., 2023). However, these studies do not properly control their variables, and vary both the input length and the associated tasks to be performed. This makes it hard to say if the degraded performance is due to the requirement to work with longer input, or due to the task being generally harder.

In this work, we study the effect of increasing the input length on model performance, while keeping other factors as constant as possible.

We employ a methodology to measure model performance trends as a function of input length,

*These authors contributed equally to this work.

by isolating it as a variable, while keeping the underlying task intact (§2).

To that end, we introduce **Flexible LENGTH Question Answering** dataset (FLenQA)¹, a QA dataset for text-based reasoning (§3). For each sample, composed of a True/False question over two pieces of information required to answer it (the context), we create multiple versions of different lengths by embedding the context parts within longer, irrelevant texts. To ensure that models utilize their entire input, the dataset is composed of tasks for which both pieces of information must be reasoned over together in order to correctly answer the question. At the same time, we keep the tasks simple enough such that models answer most of them correctly when the information pieces are presented on their own, with no additional padding.

We show that LLMs quickly degrade in their reasoning capabilities, even on input length of 3000 tokens, which is much shorter than their technical maximum (on average over all tested models, a drop in accuracy from 0.92 to 0.68).

Additionally, we explore the effect of embedding the information pieces in various locations within the context, as well as with two kinds of contexts: similar to the information pieces, or dissimilar to them (§4). We find that regardless of the experimental setting, there are similar trends of degradation.

We also show that next-word prediction performance of models on long inputs is uncorrelated with their performance on downstream tasks of reasoning on long inputs (§5).

Furthermore, we find that while *Chain-of-Thought* (CoT) prompting (Kojima et al., 2022; Wei et al., 2022) increases performance in short inputs, in most models it does not mitigate the degradation of performance when inputs are longer: while CoT prompting increases the accuracy over non-CoT prompting, the amount of increase is roughly consistent across context lengths, and is far from closing the performance drop due to long context (§6). The only exception to that is GPT4², in which the gap between CoT and normal prompting increases as the input is longer.

Finally, we analyse our results and identify several failure modes in model responses (§7). We find that with longer inputs models tend not to follow specific instructions in the input, either providing

no answer, or - in the case of CoT prompting - presenting the final answer before outlining the reasoning steps. We also observe a bias towards answering "false", as well as a decline in the models' ability to incorporate relevant information in their responses, as input length increases.

2 Desired Data Properties

Our goal is to understand how input length affects LLMs reasoning capabilities over text, given that the relevant information remains the same. We thus use question answering tasks that require models to reason over a given text. For the investigation to be applicable to both open and closed models, we chose a behavioral approach that relies on input intervention (Holtzman et al., 2023).

We aim for our data to satisfy the following requirements:

Ensuring models reason over the input. To examine the performance of models on long inputs, we require that the task can only be solved correctly by drawing conclusions from evidence in the text (Huang and Chang, 2022).

1. *Each data sample should contain several relevant text spans that are both necessary and sufficient to correctly solve the task.*
2. *All relevant spans must be consulted jointly to reach a successful solution.* Some tasks, like text summarization, can be solved using a "divide-and-conquer" approach (Gidiotis and Tsoumakas, 2020; Liu et al., 2022; Wolhandler et al., 2022), where each relevant span is individually identified, and then paraphrased and added to the output. We wish to avoid such decomposable tasks, as they do not really require reasoning over long inputs.
3. *To avoid model reliance on parametric knowledge rather than on the text, and to avoid data contamination (Jacovi et al., 2023) the question and supporting relevant spans should consist of novel facts not seen in training.*³

Isolating the length factor. To isolate the effect of length, we impose the following requirements:

1. *The required reasoning should be independent of the length of the sample:* the relevant spans

¹<https://github.com/alonj/Same-Task-More-Tokens>

²we refer to the models gpt-4-1106-preview, gpt-3.5-turbo-1106 as GPT4 and GPT3.5 accordingly.

³Models often answer the question correctly even if none, or only one of the required supporting facts is present in its input. We discuss this further in Appendix A.

should remain the same in all length variations.

2. The *added material* (a.k.a “padding”, text that is added to control the samples’ length) *should not contradict or interfere with the reasoning over the relevant text spans*.
3. The location of each relevant span within the input should be controllable.

Maintaining natural-looking inputs. The input should reflect something a user may naturally use in an LLM prompt. For example, a sequence of unrelated sentences is not natural. In contrast, a sequence of unrelated paragraphs but where each paragraph is cohesive is more natural, as such an input may result from collecting relevant information from multiple sources. To best maintain the naturality of the inputs while changing an input’s length, we require that the input should be cohesive at least at the level of paragraphs.

3 FLenQA

We introduce the **Flexible LENgth Question Answering** dataset (FLenQA), which follows the requirements set in §2.

FLenQA is composed of three reasoning tasks: Monotone Relations (a new task), People In Rooms (a new task) and a simplified version of Ruler (Clark et al., 2021) (§3.2). Each task consists of 100 base instances, from which we create variations of different lengths, different background texts, and different dispersion of facts within the background texts (§3.3).

Each task is completely balanced in its label distribution (“True” and “False”), and we ensure that most base-instances within it will be solved correctly by the LLMs when presented in their unexpanded forms (§3.4).

We release the dataset to support future studies of reasoning and long input performance. Details and statistics of the tasks appear in Appendix B.

3.1 Base instances.

Each base-instance consists of (1) an *optional prefix* (for example introducing the task or supporting facts); (2) *two key paragraphs*, each of which is thematically coherent and starts with a *key sentence* needed for solving the task; and (3) an *optional suffix* (for example, asking a question about the pre-

ceding context).⁴ For each instance, the different parts are joined by newlines and fed to the LLM.

Throughout the text, key paragraphs are typeset in red, the supporting sentences within them in darker red, and the optional prefixes and suffixes in black. The full prompts used for each dataset are in Appendix D.

Deriving the key paragraphs Each task relies on two facts, expressed as simple sentences. Each of these sentences is then expanded to a thematically-coherent paragraph, in order to ensure the naturality requirement. This expansion is performed using GPT-4, which we prompt to extend the sentences without adding new information, followed by a manual verification of the results by the authors.

3.2 The tasks

Monotone relations (MonoRel) Each key sentence is comparing two person names on monotone scale, e.g. “X is larger than Y”, “Y is larger than Z”. The suffix is a True/False question that asks about a relation between two entities that appear in different sentences (they are not explicitly compared in the text). The relations are transitive and monotone in nature.

MonoRel Example:

Julie Baker is younger than Julian Barton.
This is a fact that remains constant, unchanging like the northern star. It’s a truth that is as clear as day that she ...
Samantha Arnold is younger than Julie Baker.
It means that Samantha Arnold has experienced fewer birthdays than Julie Baker. ...
Is Samantha Arnold younger than Julian Barton?

This data is inspired by different monotonic relations describing kinship, introduced by Sinha et al. 2018. We define a new set of relation types in this work. Following the requirements in §2, answering the question requires reasoning over both key sentences. The data is created programmatically by randomly drawing names from Faker python library (Faraglia and Contributors, 2012) and a relation from a list of hand-crafted relations.

⁴The optionality is at the task level, either all instances in the task have a prefix/suffix, or they don’t.

People In Rooms (PIR) In each sample in the task, in one key sentence person is said to be located in a named room (“*X is in the old library*”), and the other key sentence describes the room to have a certain property (“the old library has wooden floors”). The task is then to infer whether the given person is located in a room with the given property.

PIR Example:

John’s living room is marble-floored, a reality that is as intrinsic to the building as its very foundations. The moment ...
 Ethan Washington is in John’s living room, a fact that has become as much a part of the place as the walls and the ceiling. The truth that Ethan Washington is in John’s living ...
 Is Ethan Washington in a marble-floored room?

This dataset is inspired by the bAbI set of tasks (Weston et al., 2016), where reasoning is conducted on paths taken by one or more agents. PIR is a simplification of the task, involving just one agent. The names of people in the task are drawn randomly (Faraglia and Contributors, 2012). Rooms and properties were hand selected to be mutually exclusive (for example, a room is either blue-walled or red-walled), so no ambiguous examples are created.

Simplified Ruletaker We employ the task formulation from Ruletaker (Clark et al., 2021), a benchmark designed for theorem proving within texts that present explicit logical theories in natural language. Each instance consists of a logical rule, two sentences each introducing a fact, and a question over the rule and facts.⁵

Simplified Ruletaker Example:

Facts:
 Erin is furry. Erin is known for his furriness. He has a lot of fur and ...
 Erin is good. Erin was always known for how good he is. His goodness appears on all matters of life ...
 Rule: If X is big and X is good then X is tall.
 Question: can the statement “Erin is tall” be derived from the rule and the facts?

⁵Initial experiments revealed that most LLMs still struggle with instances involving multiple rules or more than two facts. Our Simplified Ruletaker task consists of generated samples that fit these criteria.

3.3 Length Variations

We expand each base instance to input lengths of roughly 250, 500, 1000, 2000, and 3000 tokens.⁶ To extend the inputs to those targets we add background text that is irrelevant to the question (“padding”, §2). For each basic-instance and length pair we create different versions that differ in their source of background text: either *duplicate*, *similar* or *different* than the key paragraphs of the instance. For each of these, we also vary the dispersion of the key-paragraph within the background text.

3.3.1 Background Texts

Duplicate. To evaluate the extreme case where the length changes but the information remains the same, we perform an experiment where the each length text consists of multiple copies of the key paragraph. We duplicate each key paragraphs without any modification to achieve the target length of the input. The two duplicated paragraphs appear in alternating order until the desired sample length is achieved. In this case, of the two sub-tasks of QA reasoning - identifying the key information and reasoning over it, the first sub-task is trivial.

Similar: resampling from the same task. To get background text that is similar to the key paragraphs, we pad using paragraphs sampled from other base instances of the same task. To avoid creating contradictions, we exclude paragraphs that contain entities appearing in the key paragraphs. This padding therefore does not produce adversarial or ambiguous versions of the samples.

Different: Book Corpus. To get background text that differs from the key paragraphs, we use text from the Books Corpus (Zhu et al., 2015). We sample a random (continuous) text from the Book Corpus, and inject each of the key paragraphs within it, while respecting sentence boundaries.

3.3.2 Location of key paragraphs in the text

We consider four distinct ways in which the key paragraphs are dispersed within the background text: in the first three cases the key paragraphs appear adjacent to each other, while in the fourth the key paragraphs are separated by intervening text of various lengths.

(1) *Key paragraphs first:* The key paragraphs appear at the beginning of the text followed by padding;

⁶We consider a sample to be of length N if its token count as measured by the GPT4 tokenizer is in $(N - 70, N + 70)$.

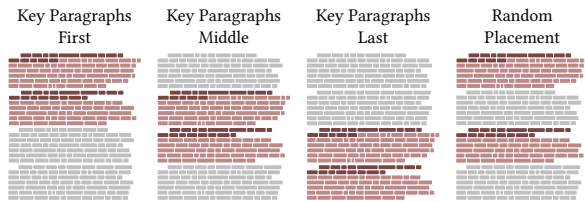


Figure 2: **Inputs construction.** Key sentences (dark red), are expanded to key paragraphs (light red) which are dispersed in controlled locations among padding text (grey) which is irrelevant to the task.

(2) *Key paragraphs middle*: Half of the padding is affixed before and half after the key paragraphs, but not between them (the key paragraphs are exactly in the middle);

(3) *Key paragraphs last*: The key paragraphs appear at the end of the text, with padding prepended before them as a prefix;

(4) *Random placement*: padding is added before, between and after the paragraphs, with random intervals.

A visual representation is provided in Figure 2.

3.4 Base instances are answerable

We estimate the baseline accuracy by evaluating the LLMs on the minimal text of each sample in the dataset that includes only the question and the key paragraphs relevant to it. We found that even when using non-CoT prompting, four out of the five models achieve high accuracy (>0.89). The lowest performing model (GPT3.5) achieve high enough accuracy for degradation to be observable (0.77). Full results can be found in Appendix C.

4 Main Experiments

We report average accuracies over all three tasks, and maintain the same setup (prompt, temperature, etc.) over all input lengths. We evaluate five recent capable LLMs: GPT4, GPT3.5, Gemini-Pro, Mistral 70B and Mixtral 8x7B. See Appendix E for a detailed breakdown of our setup parameters.

4.1 Impact of Length and Location

We start by validating the impact of input length on LLM reasoning performance (Figure 1) in various experimental settings.

No irrelevant paragraphs We first look into the extreme case where only relevant tokens are added (“duplicate padding”). Shi et al. (2023) Demonstrate that appending irrelevant texts to the input

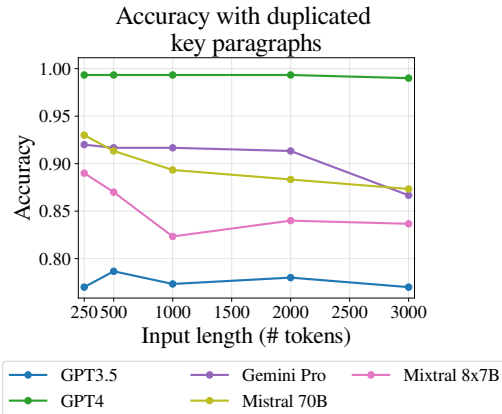


Figure 3: The relevance of padding is a factor, but it is distinct from the effect of length itself. Some models degrade in reasoning performance. Note, both GPT3.5 and GPT4 are less affected by length when the added tokens are relevant. Each point reflects 300 samples.

of a reasoning task (GSM-8K (Cobbe et al., 2021)) reduces model performance substantially. We isolate the effect of relevance by testing a setting in which the padding is duplications of the exact text of the key paragraphs. In this setup, the LLMs are not required to “search” the input to find the key paragraphs, so any bias towards any position becomes irrelevant (Liu et al., 2023b). Also, any difficulty that might be imposed by the distance between the key paragraphs also becomes irrelevant. Hence, we expect that there will be no degradation in performance. Surprisingly, the *Results* shown in Figure 3, reveal that even in this setup length does play a factor, *and accuracy decreases with length for all models.*

Adjacent paragraphs surrounded by irrelevant ones We now move to the more realistic case where the prompt includes the key paragraphs as well as additional irrelevant ones. In the first set of experiments, we keep the key paragraphs adjacent to each other: the LLM just needs to focus and operate on a single area of the input, ignoring the rest. Liu et al. (2023b) Found that in the task of extractive QA, the position of the answer in the text affects the ability of models to answer correctly. We thus experiment with the three scenarios: positioning both key paragraphs at the start, end or middle of the text. In all cases we average over both types of irrelevant padding.

The results in Figure 4 show a significant drop in accuracy as length increase beyond 500 tokens. For most models, adjacency of key paragraphs produces higher accuracy, and when the key para-

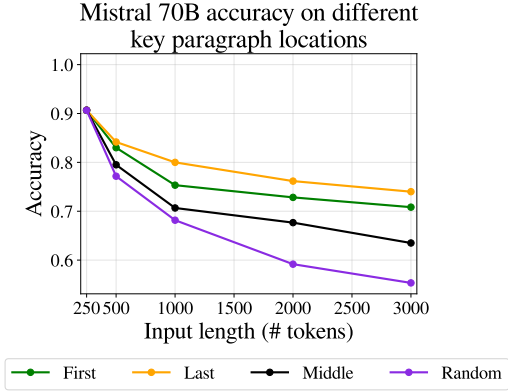


Figure 4: Effect of key paragraphs positions on accuracy. Accuracy decreases as input length grows regardless of where the key paragraphs are placed within the input. Each point reflects 300 samples. Results for all models appear in Appendix E

graphs appear last, accuracy is often highest (suggesting recency bias).

Non-adjacent relevant paragraphs. Finally, we test the scenario in which the relevant facts need to be collected from two non-adjacent locations within the text.

Here, *the results* in Figure 1 show a very large drop in performance as length increases, indicating that reasoning tasks become significantly harder for LLMs when they need to collect evidence from two distinct locations in a large-ish context length.

4.2 Kind of irrelevant material

We now focus only on the non-adjacent key-paragraphs case, and explore the effect of the kind of irrelevant text. We consider two scenarios: when the irrelevant paragraphs are *similar* to the relevant ones (taken from the same task), and when they are *different* (taken from the books corpus).

Our initial expectation was that the setup in which the irrelevant paragraphs are *different* from the relevant ones will be easier for the model, as the irrelevant paragraphs will be easier to discard, aiding focusing on the relevant ones. However, the results (Figure 5) show that is not the case: the drop for the *different* setup is mostly larger than for the *similar* one.

5 Correlation with Next Word Prediction

Perplexity is used as the main benchmark to show that models utilize their entire input (Anil et al., 2023; Jiang et al., 2024). However, it was shown that performance on downstream tasks does not

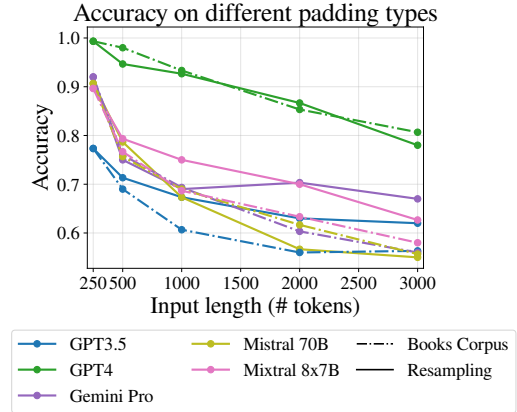


Figure 5: Performance degrade in both types of padding. Books padding impact is much greater in most models. Each point reflects the performance across 300 samples.

necessarily correlate with model perplexity (Liu et al., 2023a; Xia et al., 2022; Tay et al., 2022). Here, we will use the flexibility of our dataset to understand the correlation between perplexity and reasoning accuracy.

In closed models we lack access to full vocabulary token probabilities so model perplexity cannot be measured, therefore we resort to measuring next word accuracy on our data. We prompt models to complete the next word in a given text, and the output is correct if it is an exact match to the true next word. We use the samples in our dataset (without the questions) as the text and compare the results to the reasoning performance on the same samples.

Our method finds similar trends on the next word prediction task to those shown in other works (Anil et al., 2023; Jiang et al., 2024), namely accuracy increases as input is longer. However, as shown in Figure 1, next word accuracy correlates negatively with reasoning on FlenQA⁷.

This implies that measuring next word prediction and, similarly, perplexity, cannot substitute downstream task evaluation on long inputs.

6 Does Chain of Thought Help?

Chain of Thought (CoT) prompting, introduced by Kojima et al. (2022); Wei et al. (2022), is a technique by which the LLM is pushed to produce a text comprising of reasoning steps before concluding the correct answer for a question. Zhou et al. (2022) found that a more specific and optimised instruction ("Let's work this out in a step by step way to be sure we have the right answer.>").

⁷ $\rho_{Pearson} = -0.95, p = 0.01$

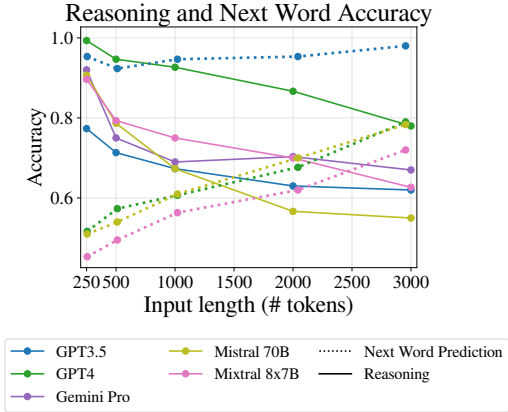


Figure 6: Next word accuracy correlates negatively with the reasoning accuracy on FlenQA. Each point reflects the performance across 300 samples. Gemini-Pro is not included as it answered empty replies to the next word prediction task at any length.

The CoT technique was shown to significantly improve the accuracy on many reasoning-based question-answering setups. Will using it change the trend and allow the LLMs to perform effectively on longer inputs? We experiment with CoT using the elicitation string of Zhou et al. (2022).

The results show (Figure 1) that CoT has different effects on different LLMs, and overall does not mitigate the drop in performance due to length. In most cases (GPT4, Mixtral 8x7B, Mistral 70B and GPT3.5) it improves performance, but only in GPT4 it has an increased effect as length increases, making it a limited mitigation technique. In the case of Gemini-Pro, we see that CoT decrease performance as input length is increased, even though it increase performance on short length.

The full results of the CoT prompting over all tasks and setups can be found in Appendix E.

7 Length-induced Failure modes

We find in the results four *failure modes*:⁸ consistent patterns that correlate with incorrect responses.

Failure to answer All of the samples in the dataset can be answered with either "True" or "False", as instructed in our prompts (Appendix D). However, some of LLMs responded with a refusal answer the question, often preceded by a sentence such as "There is not enough information in the text". **This tendency grows as the input length increases**, indicating a failure to comply to

⁸All failure modes can be measured automatically using the code in our repository.

the instruction that specified a clear choice between "True" and "False". The trend is demonstrated in figure 7, and results over all models in Appendix E.

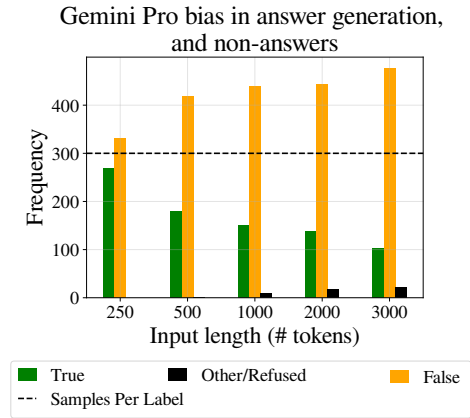


Figure 7: The models exhibit two types of input-length dependent biases: (a) They tend to generate "False" more often than "True", and (b) they ignore instructions and generate answers which do not contain neither.

Label bias As discussed in §3, our dataset is completely balanced in the label distribution. We find that certain LLMs tend to favour one of the labels, typically "false", as the input length grows. Results for all models are in Appendix E.

Answer first, reason later When using Chain-of-Thought prompting, some LLMs were much more likely to output the final true/false answer *before* the expected reasoning steps, as inputs grow longer. In recent work, Kojima et al. 2022 found that when models are elicited to provide the reasoning steps after the answer their performance does not increase (as expected when using autoregressive models that only attend to earlier tokens). This can be viewed as a case of failing to follow prompt instructions (see prompt instructions in Appendix D) as the length increases. In testing, we found that incorrect responses are statistically dependent on the occurrence of answers before the reasoning steps⁹.

Chain-of-Thought lack of coverage All the tasks in FlenQA require the LLM to: (1) locate the relevant texts within the input; and (2) perform the relevant reasoning over them. Ideally, the CoT prompt would elicit the LLM to first locate each of the relevant texts and copy them to the "steps" part,

⁹Corresponding odds-ratio is 3.643 with $p < 0.001$ obtained through Fisher exact test.

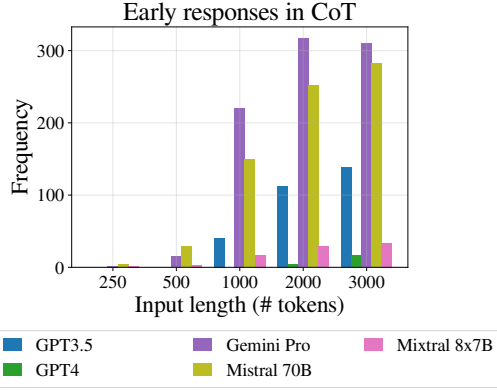


Figure 8: Most of the models tend to generate an answer before the reasoning steps, in a zero-shot CoT prompt setting, and do so more as input length increases.

hence avoiding the effect of long input on reasoning. However, we find that as input length grows, LLMs ability to do this degrades (Figure 9).

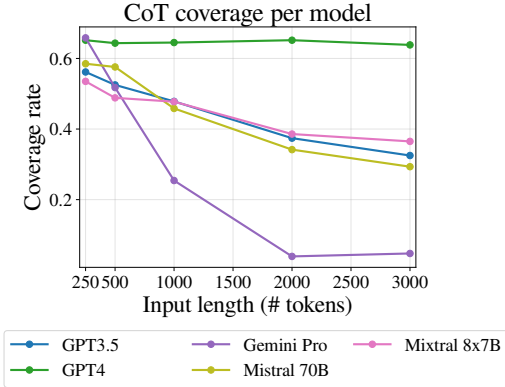


Figure 9: CoT coverage of relevant facts. As input grows, all models fail more often in outputting the task-relevant information at the CoT reasoning steps stage.

We measure this by computing the coverage of the relevant text (the key sentences in each sample) in the models’ “steps” part of the outputs (details in Appendix D.4). We find that in most models, the ability to locate the relevant text within the input *decreases* as the input length gets longer. We found incorrect responses were statistically dependent on the incomplete coverage of the facts¹⁰.

8 Related Work

The evaluation of LLMs on long inputs has followed two distinct pathways: benchmarks of downstream tasks and next word prediction. In the realm of benchmarks, studies proposed datasets of long

input samples that can be used to evaluate models (Shaham et al., 2023, 2022; An et al., 2023b,a). Those datasets are curated over inputs of different, but fixed, length. This approach, while straightforward, limits our ability to inputs of varying lengths, posing a challenge in understanding the true impact of input length on model performance. On the other hand, next word prediction evaluations do offer an insights into how models handle inputs of different lengths (like done in Anil et al. 2023; Jiang et al. 2024). However, the correlation of this task with downstream performance was found not consistent (Liu et al., 2023a; Xia et al., 2022; Tay et al., 2022). In this paper we reproduce this finding with respect to extended length.

This study builds upon prior research that examined different aspects through input intervention, studying the semantic content (theme) of a task (Dasgupta et al., 2022), prompting strategies (Kojima et al., 2022; Yao et al., 2023; Jin et al., 2024) and various properties of the QA task (Levy et al., 2023). Our investigation focuses on input length, isolating it, to reveal its impact on performance.

9 Discussion

We study the effect of input length on reasoning performance of current Large Language Models (LLMs). Our findings reveal a significant drop in performance with longer inputs, occurring well before reaching the models’ maximum input-length capacity. Our experiments relied on FLenQA, a dataset we constructed that allows to isolate the length factor, by adjusting the parts in the input that are irrelevant to the task. We show that regardless of how we adjust the samples, there is still a strong effect of length on reasoning performance.

Finally, we identified specific failure modes, including difficulties in following extended instructions and biases towards less relevant information. Our analysis reveals specific failings, providing possible directions for future studies to address and rectify the weaknesses found in LLMs.

In conclusion, our work indicates that evaluating a model’s performance based on a single input length does not provide a full picture, and more nuanced evaluation is required. We argue that for a model to be considered capable at long range, it must maintain its performance at any length it technically supports.

¹⁰Corresponding odds-ratio is 3.138 with $p < 0.001$ obtained through Fisher exact test.

Limitations

Because of the nature of behavioral testing, the observed drop in performance with varying input lengths remains unexplained; because of lack of access to many of the models, we suspect this direction will continue to be limited. Secondly, our approach aimed to create a universally applicable test across different LLMs, leading to the selection of tasks that cater to the lowest common denominator. This approach potentially overlooks the nuanced performance differences in more complex reasoning tasks (e.g 5 key paragraphs), where, for instance, stronger models might exhibit performance degradation at shorter input lengths compared to what our findings suggest. Additionally, we focused on a subset of reasoning task types which may differ behaviourally from other types. Finally, our study did not test the distance between key paragraphs, leaving an aspect of LLM performance unexplored that we leave for future research.

Acknowledgements

References

- Chen An, Shansan Gong, Ming Zhong, Mukai Li, Jun Zhang, Lingpeng Kong, and Xipeng Qiu. 2023a. [L-eval: Instituting standardized evaluation for long context language models](#). *ArXiv*, abs/2307.11088.
- Chenxin An, Shansan Gong, Ming Zhong, Mukai Li, Jun Zhang, Lingpeng Kong, and Xipeng Qiu. 2023b. [L-eval: Instituting standardized evaluation for long context language models](#).
- Rohan Anil, Sebastian Borgeaud, Yonghui Wu, and Gemini Team Google. 2023. [Gemini: A family of highly capable multimodal models](#). *ArXiv*, abs/2312.11805.
- Yushi Bai, Xin Lv, Jiajie Zhang, Hongchang Lyu, Jiankai Tang, Zhidian Huang, Zhengxiao Du, Xiao Liu, Aohan Zeng, Lei Hou, et al. 2023. Longbench: A bilingual, multitask benchmark for long context understanding. *arXiv preprint arXiv:2308.14508*.
- Jifan Chen and Greg Durrett. 2019. Understanding dataset design choices for multi-hop reasoning. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*.
- Peter Clark, Oyvind Tafjord, and Kyle Richardson. 2021. Transformers as soft reasoners over language. In *Proceedings of the Twenty-Ninth International Conference on International Joint Conferences on Artificial Intelligence*, pages 3882–3890.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.
- Ishita Dasgupta, Andrew K Lampinen, Stephanie CY Chan, Antonia Creswell, Dharshan Kumaran, James L McClelland, and Felix Hill. 2022. Language models show human-like content effects on reasoning. *arXiv preprint arXiv:2207.07051*.
- Daniele Faraglia and Other Contributors. 2012. [Faker](#).
- Alexios Gidiotis and Grigorios Tsoumakas. 2020. A divide-and-conquer approach to the summarization of long documents. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 28:3029–3040.
- Ari Holtzman, Peter West, and Luke Zettlemoyer. 2023. Generative models as a complex systems science: How can we make sense of large language model behavior? *arXiv preprint arXiv:2308.00189*.
- Jie Huang and Kevin Chen-Chuan Chang. 2022. Towards reasoning in large language models: A survey. *arXiv preprint arXiv:2212.10403*.
- Alon Jacovi, Avi Caciularu, Omer Goldman, and Yoav Goldberg. 2023. [Stop uploading test data in plain text: Practical strategies for mitigating data contamination by evaluation benchmarks](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 5075–5084, Singapore. Association for Computational Linguistics.
- 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, L lio Renard Lavaud, Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le Scao, Th ophile Gervet, Thibaut Lavril, Thomas Wang, Timoth e Lacroix, and William El Sayed. 2024. [Mixtral of experts](#).
- Mingyu Jin, Qinkai Yu, Haiyan Zhao, Wenyue Hua, Yanda Meng, Yongfeng Zhang, Mengnan Du, et al. 2024. The impact of reasoning step length on large language models. *arXiv preprint arXiv:2401.04925*.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. *Advances in neural information processing systems*, 35:22199–22213.
- Mosh Levy, Shauli Ravfogel, and Yoav Goldberg. 2023. Guiding llm to fool itself: Automatically manipulating machine reading comprehension shortcut triggers. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 8495–8505.

- Jiaqi Li, Mengmeng Wang, Zilong Zheng, and Muhan Zhang. 2023. [Loogle: Can long-context language models understand long contexts?](#) *ArXiv*, abs/2311.04939.
- Hong Liu, Sang Michael Xie, Zhiyuan Li, and Tengyu Ma. 2023a. Same pre-training loss, better downstream: Implicit bias matters for language models. In *International Conference on Machine Learning*, pages 22188–22214. PMLR.
- Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2023b. Lost in the middle: How language models use long contexts. *arXiv preprint arXiv:2307.03172*.
- Yang Liu, Chenguang Zhu, and Michael Zeng. 2022. End-to-end segmentation-based news summarization. In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 544–554.
- Sewon Min, Eric Wallace, Sameer Singh, Matt Gardner, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2019. [Compositional questions do not necessitate multi-hop reasoning](#). In *Annual Meeting of the Association for Computational Linguistics*.
- OpenAI. 2023. [Gpt-4 technical report](#).
- Oscar Sainz, Jon Campos, Iker García-Ferrero, Julen Etxaniz, Oier Lopez de Lacalle, and Eneko Agirre. 2023. Nlp evaluation in trouble: On the need to measure llm data contamination for each benchmark. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 10776–10787.
- Uri Shaham, Maor Ivgi, Avia Efrat, Jonathan Berant, and Omer Levy. 2023. Zeroscrolls: A zero-shot benchmark for long text understanding. *arXiv preprint arXiv:2305.14196*.
- Uri Shaham, Elad Segal, Maor Ivgi, Avia Efrat, Ori Yoran, Adi Haviv, Ankit Gupta, Wenhan Xiong, Mor Geva, Jonathan Berant, et al. 2022. Scrolls: Standardized comparison over long language sequences. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 12007–12021.
- Freda Shi, Xinyun Chen, Kanishka Misra, Nathan Scales, David Dohan, Ed H Chi, Nathanael Schärli, and Denny Zhou. 2023. Large language models can be easily distracted by irrelevant context. In *International Conference on Machine Learning*, pages 31210–31227. PMLR.
- Koustuv Sinha, Shagun Sodhani, William L. Hamilton, and Joelle Pineau. 2018. [Compositional language understanding with text-based relational reasoning](#). *ArXiv*, abs/1811.02959.
- Yi Tay, Mostafa Dehghani, Samira Abnar, Hyung Won Chung, William Fedus, Jinfeng Rao, Sharan Narang, Vinh Q Tran, Dani Yogatama, and Donald Metzler. 2022. Scaling laws vs model architectures: How does inductive bias influence scaling? *arXiv preprint arXiv:2207.10551*.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 35:24824–24837.
- Jason Weston, Antoine Bordes, Sumit Chopra, Alexander M Rush, Bart Van Merriënboer, Armand Joulin, and Tomas Mikolov. 2016. Towards ai-complete question answering: A set of prerequisite toy tasks. In *4th International Conference on Learning Representations, ICLR 2016*.
- Ruben Wolhandler, Arie Cattan, Ori Ernst, and Ido Dagan. 2022. How “multi” is multi-document summarization? In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 5761–5769.
- Mengzhou Xia, Mikel Artetxe, Chunting Zhou, Xi Victoria Lin, Ramakanth Pasunuru, Danqi Chen, Luke Zettlemoyer, and Ves Stoyanov. 2022. Training trajectories of language models across scales. *arXiv preprint arXiv:2212.09803*.
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L Griffiths, Yuan Cao, and Karthik Narasimhan. 2023. Tree of thoughts: Deliberate problem solving with large language models. *arXiv preprint arXiv:2305.10601*.
- Yongchao Zhou, Andrei Ioan Muresanu, Ziwen Han, Keiran Paster, Silviu Pitis, Harris Chan, and Jimmy Ba. 2022. Large language models are human-level prompt engineers. *arXiv preprint arXiv:2211.01910*.
- Yukun Zhu, Ryan Kiros, Rich Zemel, Ruslan Salakhutdinov, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. 2015. Aligning books and movies: Towards story-like visual explanations by watching movies and reading books. In *The IEEE International Conference on Computer Vision (ICCV)*.

A Training Contamination in Reasoning Tasks

Data contamination is a major concern when evaluating models (Sainz et al., 2023; Jacovi et al., 2023). Ensuring that a task requires reasoning across multiple text spans is a stronger requirement than a task that requires multi hop reasoning (§2). Evaluating models on questions they answer using parametric knowledge prevents us from assessing their reasoning capabilities. Datasets originating from internet sources are especially vulnerable to contamination, thereby undermining the evaluation of a model’s capacity to reason over novel facts.

Furthermore, it was shown that small models can answer existing reasoning dataset when given one

some of the parts that were claimed to be required for the task (Chen and Durrett, 2019; Min et al., 2019). This raises question on how well those dataset can be used to evaluate the effect of length, where we require that the reasoning will be done on multiple parts of the text.

We conclude that parametric knowledge should be accounted for when evaluating text-based reasoning capabilities. In this work we introduced FlenQA, which is composed of novel generated data to make sure that reasoning over the input is required.

B Datasets

Each task of the following contains 100 base instances. In each sample, there are two paragraph-length texts (*key paragraphs*). To achieve paragraphs of similar length, we edit them by truncating sentences beyond a specific length, resulting in an average paragraph length of 125 tokens.

B.1 Ruletaker

The key paragraphs in the task are as evidence for the reasoning task, a rule and a question. In the original data (Clark et al., 2021), the samples contain different number of reasoning steps. In this study, we generate new, simpler samples of the task: each sample is composed of only two facts and one logical rule. The samples we generate are of similar flavor to those that exist in the original Ruletaker data, but are generated with new statements, rules and facts. The key paragraphs and the padding appear as the facts of each sample.

Padding Type	Target Input Length	Mean Number Tokens
Books	250	249.8
	500	508.78
	1000	1009.56
	2000	2009.64
	3000	3008.38
Same	250	249.8
	500	503.535
	1000	1004.41
	2000	2005.51
	3000	3005.125

Figure 10: Summary of statistics of the Ruletaker* task data.

B.2 MonoRel

The key paragraphs in the task act as evidence for the reasoning task, and a question. Both key paragraphs describe a monotonic relation between two

people, where one person is shared between both. The key paragraphs are embedded in padding text to create a text mixture.

Padding Type	Target Input Length	Mean Number Tokens
Books	250	238.06
	500	490.84
	1000	991.41
	2000	1990.34
	3000	2990.95
Same	250	238.06
	500	491.69
	1000	991.43
	2000	1991.31
	3000	2991.44

Figure 11: Summary of statistics of the MonoRel task data.

B.3 People in Rooms (PIR)

One key paragraph describes the location of an individual, and the other describes some attribute of that location. The key paragraphs are embedded in padding text to create a text mixture.

Padding Type	Target Input Length	Mean Number Tokens
Books	250	305.36
	500	491.85
	1000	989.91
	2000	1992.00
	3000	2988.67
Same	250	305.36
	500	484.63
	1000	985.82
	2000	1985.04
	3000	2984.80

Figure 12: Summary of statistics of the People In Rooms (PIR) task data.

C Base Instances Result

Model	Prompt	MonoRel	PIR	Ruletaker*
GPT3.5 Turbo	Direct	0.77	0.81	0.74
	CoT	0.86	0.88	0.88
GPT4 Turbo	Direct	1.00	1.00	0.98
	CoT	1.00	1.00	0.97
Gemini Pro	Direct	0.84	1.00	0.92
	CoT	0.88	0.96	0.97
Mistral 70B	Direct	0.99	1.00	0.73
	CoT	1.00	1.00	0.89
Mixtral 8x7B	Direct	0.92	0.97	0.80
	CoT	0.86	0.97	0.93

Table 1: **Minimal length accuracy.** The evaluated models have high accuracy on the tasks in our dataset when evaluated on the minimal text (250 tokens). CoT improve performance across almost all tasks and models.

D Full Evaluation Setup

D.1 Prompts

Rulemaker prompt - Normal:

Answer whether the statement {statement} can be derived from the rule and the facts. Answer with either "True" or "False".

Rule: {rule}

Facts: {facts + padding}

Answer with either "True" or "False".

Rulemaker prompt - CoT:

Answer whether the statement {statement} can be derived from the rule and the facts.

Show your steps then answer with either "True" or "False".

Rule: {rule}

Facts: {facts + padding}

Answer with either "True" or "False". Let's work this out in a step by step way to be sure we have the right answer.

PIR prompt - Normal:

{facts + padding}

True/False Question: {question}

Answer only True or False.

PIR prompt - CoT:

Show your steps then answer with 'true' or 'false'.

{facts + padding}

True/False Question: {question}

Let's work this out in a step by step way to be sure we have the right answer.

MonoRel prompt - Normal:

Here are some facts. Answer the exact following question based on the text:

{question} Answer the question as it appears exactly. {facts + padding}

{question}

Answer only True or False.

MonoRel prompt - CoT:

Here are some facts. Answer the exact following question based on the text: {question} Answer the question as it appears exactly.

Show your steps then answer with 'true' or 'false'.

{facts + padding}

{question}

Let's work this out in a step by step way to be sure we have the right answer. Show your work and finally answer with 'true' or 'false'. The final step should include the exact text of the question and the answer.

D.2 Parameters

All models were evaluated with a temperature of 0 and "top p" of 0 where available to make results as reproducible as possible. Additionally, We configured Gemini Pro to ignore safety guardrails ("HARM_CATEGORY" configurations) to overcome its blank answers in some samples.

D.3 Locating the answer in models' replies

To identify the models' answers in their responses, we searched for the occurrences of "false" or "true," disregarding case sensitivity. In cases where these words appeared multiple times, only the last instance was considered relevant. We tested the reliability of this method by manually examining a random sample of 100 responses and confirmed its accuracy in all instances.

D.4 Evaluating the coverage of key facts in CoT

Coverage of the key facts that are relevant to the reasoning task in CoT outputs, was done by searching for (case-insensitive) match of the key sentences in the key paragraphs, within the output of each model. Full coverage means that both key sentences from the input appear in the CoT output. We verified the reliability of this method manually on a sample of 100 responses.

E Full results

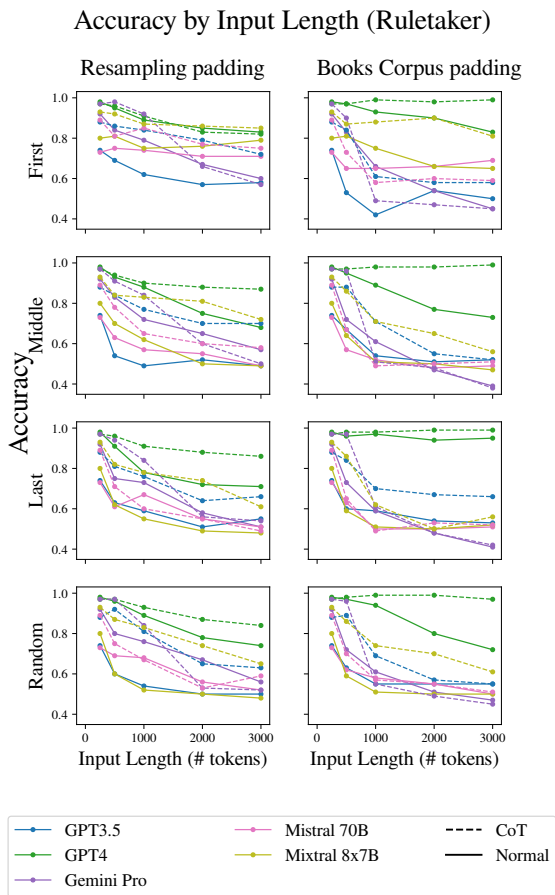


Figure 13: Full results for the Ruletaker dataset.

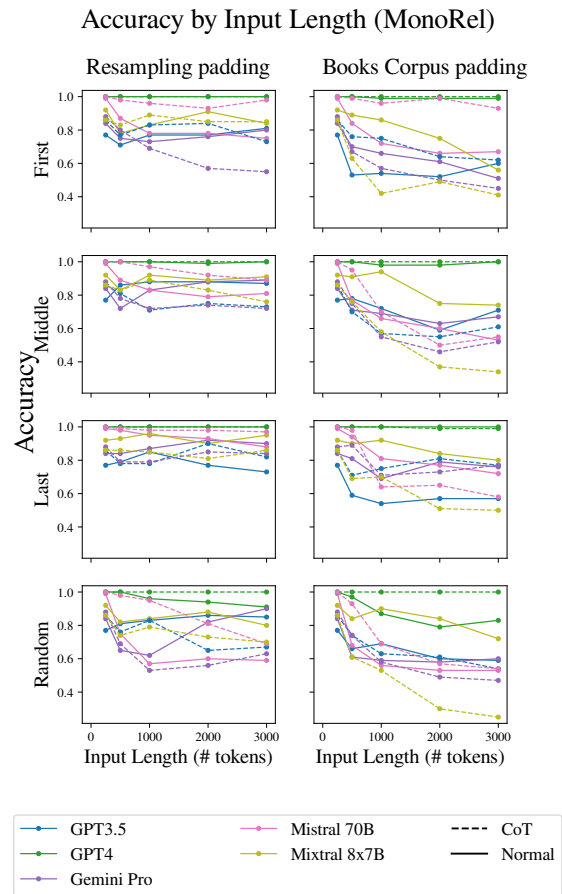


Figure 14: Full results for the MonoRel dataset.

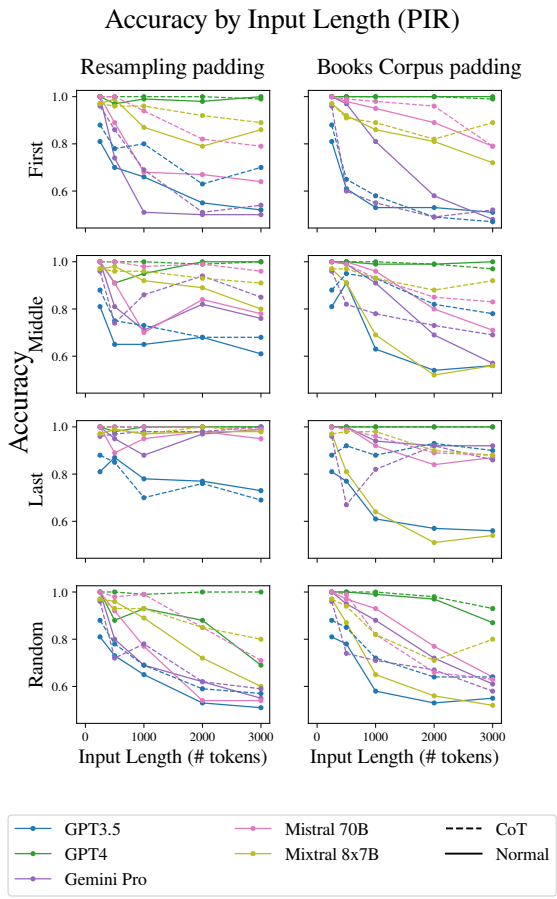


Figure 15: Full results for the People In Rooms (PIR) dataset.

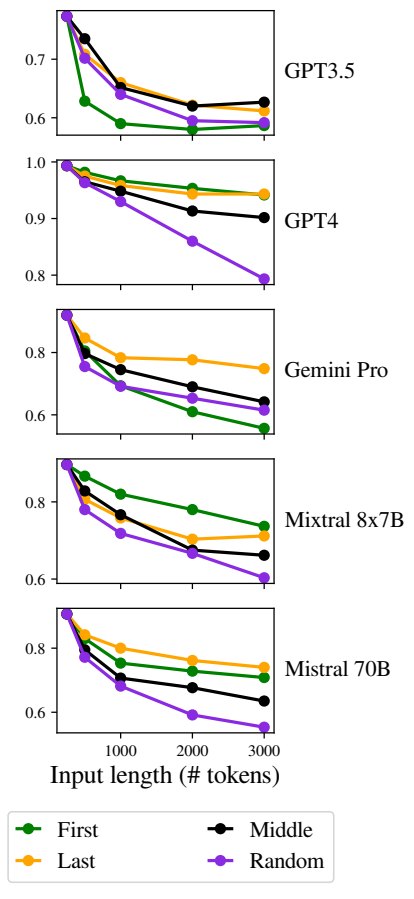


Figure 16: Differences in accuracy between different positions of key paragraphs in input. Averaged over both types of irrelevant padding: similar (resampling from the data) and dissimilar (Books corpus) padding.

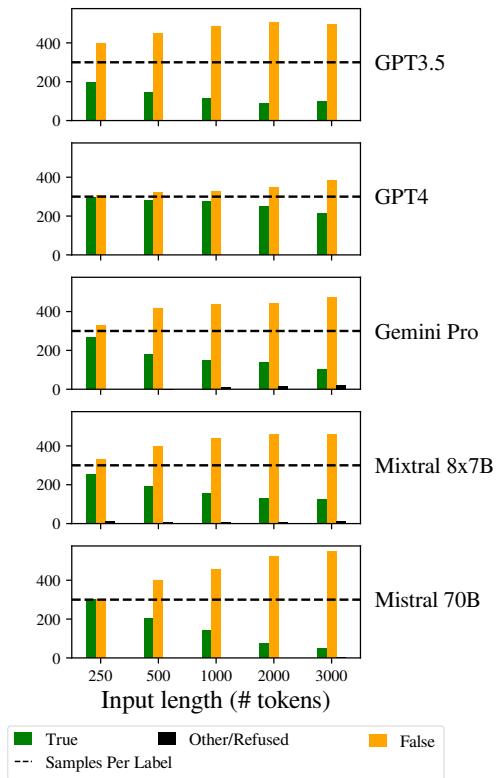


Figure 17: **Biases in answer generation and non-answers.** Frequency of responses with True, False, or neither, per model.