

On Many-Shot In-Context Learning for Long-Context Evaluation

Anonymous ACL submission

Abstract

Many-shot in-context learning (ICL) has emerged as a unique setup to both utilize and test the ability of large language models to handle long context. This paper delves into long-context language model (LCLM) evaluation through many-shot ICL. We first ask: what types of ICL tasks benefit from additional demonstrations, and how effective are they in evaluating LCLMs? We find that classification and summarization tasks show performance improvements with additional demonstrations, while translation and reasoning tasks do not exhibit clear trends. Next, we investigate the extent to which different tasks necessitate retrieval versus global context understanding. We develop metrics to categorize ICL tasks into two groups: (i) similar-sample learning (SSL): tasks where retrieval of the most similar examples is sufficient for good performance, and (ii) all-sample learning (ASL): tasks that necessitate a deeper comprehension of all examples in the prompt. Lastly, we introduce a new many-shot ICL benchmark, **MANYICLBENCH**, to characterize model’s ability on both fronts and benchmark 12 LCLMs using MANYICLBENCH. We find that while state-of-the-art models demonstrate good performance up to 64k tokens in SSL tasks, many models experience significant performance drops at only 16k tokens in ASL tasks.

1 Introduction

Long-context language models (LCLMs) have revolutionized the way users interact with language models by extending the context size from 2K to 128K or even 1M tokens (Team et al., 2023; GLM et al., 2024; Dubey et al., 2024). This unlocks challenging applications, such as long- and multi-document summarization, multi-turn dialogue, and code repository comprehension. Despite the recent progress in building LCLMs, existing benchmarks primarily evaluate these models’ retrieval capabili-

ties (Liu et al., 2023; Hsieh et al., 2024). From synthetic tasks such as Needle-in-A-Haystack (NIAH) (Kamradt, 2023) and RULER benchmark (Hsieh et al., 2024) to real-world challenges like long-novel QA (Karpinska et al., 2024), the majority of benchmarks assess how well LCLMs retrieve specific pieces of information from extensive contexts. As a result, **evaluating models’ global understanding of the full context remains lacking**.

To fill the gap, Li et al. (2024) introduce Long-ICLBench, which uses many-shot ICL classification tasks to evaluate models’ long-context performance, arguing that these tasks require the comprehension of the entire input. A few other works have also explored many-shot ICL for long-context models (Agarwal et al., 2024; Bertsch et al., 2024). However, they have mainly relied on classification tasks (Li et al., 2024; Bertsch et al., 2024), which are insufficient to distinguish which skills LCLMs require to perform well in many-shot ICL classification tasks. Agarwal et al. (2024) study non-classification ICL tasks but only limit its study to Gemini 1.5 Pro. In this work, we want to conduct a comprehensive study on many-shot ICL across a wide range of models, with a goal of identifying tasks that **benefit from additional demonstrations** and explore their utility in evaluating long-context models. Moreover, we seek to determine the extent to which these tasks rely on **retrieving most similar examples versus learning from all samples**.

RQ1: Which tasks benefit from many-shot ICL? First, we investigate ICL tasks that are used in prior work, including classification, summarization, and reasoning, under many-shot settings with context lengths from 1k to 128k (Agarwal et al., 2024). We find that classification and summarization tasks show *strong positive correlation between context lengths and model performance*. Our findings indicate that translation and reasoning tasks such as ARC (Clark et al., 2018) and FLORES-

200 (NLLB Team, 2022) do not gain much performance with an increasing number of demonstrations. Science and symbolic reasoning tasks exhibit inconsistent trends between context lengths and model performance. This variance in performance is mainly attributed to the specific nature of tasks, where more demonstrations do not boost the models’ task understanding. Interestingly, math tasks benefit from additional demonstrations only when step-by-step solutions (or chain-of-thoughts) are derived and using strong LCLMs.

RQ2: To what extent does each task require learning from a limited number of samples versus learning from more samples with broader context from LCLMs? We use the ratio between the performance change of removing dissimilar examples and the change of removing similar examples. A high ratio means a more pronounced drop in performance upon removing similar examples, which indicates the task’s heavy reliance on retrieving and learning its prediction with more similar examples. Our analysis indicates that existing many-shot ICL *classification* tasks (Li et al., 2024) predominantly assess models’ skills to learn from *similar examples* rather than all examples, overshadowing the model’s ability to understand all samples. This leads us to categorize tasks into similar-sample learning (SSL) and all-sample learning groups (ASL).

Following the categorization, we propose a new many-shot ICL benchmark, **MANYICLBENCH**, designed to evaluate long-context models and advocate for the inclusion of many-shot ICL tasks as effective evaluation candidates. Importantly, on MANYICLBENCH, models are tested to either retrieve and learn from the most similar demonstrations or assimilate and learn from all demonstrations to enhance their understanding of the task (Lin and Lee, 2024; Bertsch et al., 2024). Therefore, MANYICLBENCH *evaluates both retrieval skills and global context understanding*, thus providing a holistic assessment of long-context models’ capabilities.

Unlike NIAH, which focuses on retrieving a statement given a query, MANYICLBENCH redefines retrieval to test a model’s ability to identify and leverage similar examples during in-context learning. Additionally, the benchmark retains a focus on global context understanding, as explored in prior works (Zhang et al., 2024; Karpinska et al., 2024), by requiring models to reason over and synthesize information from all demonstrations to gen-

erate accurate outputs.

In summary, we make contributions as below:

- Investigate whether ICL tasks benefit from additional demonstrations and assess their suitability for evaluating LCLMs with a context length up to 128k tokens.
- Develop methods to characterize the primary skills evaluated by ICL tasks, where we focus on distinguishing between similar-sample learning and all-sample learning skills.
- Construct a many-shot ICL benchmark, MANYICLBENCH, designed for evaluating LCLMs on both retrieval and global context understanding, while excluding irrelevant datasets previously used in LCLM evaluation.
- Benchmark 12 widely-used state-of-the-art LCLMs on MANYICLBENCH to assess their performance comprehensively.

2 Related Work

2.1 Long-Context Language Models and Evaluation

As large language models grow in scale, there is an increasing demand for handling tasks that require extended contexts. Tasks such as long document summarization (Kryściński et al., 2022), conversations with long-context memory (Xu et al., 2021), and repository-level code completion (Zhang et al., 2023) have garnered significant interest. Advances in efficient attention mechanisms, such as flash attention (Dao et al., 2022) and grouped query attention (Ainslie et al., 2023), alongside the development of GPUs with larger memory capacities, have enabled LLMs to be trained on extended contexts. Techniques like position interpolation (Chen et al., 2023; Peng et al., 2023) and context compression (Chevalier et al., 2023; Mohtashami and Jaggi, 2023; Jiang et al., 2024) have further extended the context window size to up to 1 million tokens.

Despite these advancements, the NLP community still seeks a universal and effective method for evaluating long-context models. One prominent task is Needle-in-a-Haystack (Kamradt, 2023), which requires models to retrieve the most relevant document from a large set of documents. Currently, most evaluation benchmarks focus on synthetic tasks that primarily assess the retrieval capabilities of long-context models (Hsieh et al., 2024; Kamradt, 2023; Lee et al., 2024; Lei et al., 2024). Only

a few benchmarks, such as Karpinska et al. (2024) and Zhang et al. (2024), emphasize the model’s ability to comprehend the global context. For example, Karpinska et al. (2024) manually curated a set of challenging questions based on various novels to evaluate global context understanding. It creates a realistic long-context benchmark emphasizing retrieval and global context understanding.

2.2 Many-shot ICL with LCLMs

As the context length of a model grows, the number of demonstrations that can be utilized in ICL also increases. Studies by Li et al. (2024), Bertsch et al. (2024), and Agarwal et al. (2024) have examined various properties of ICL under the many-shot setting. Bertsch et al. (2024) explore whether models are merely performing retrieval tasks or genuinely understanding the tasks during many-shot ICL classification. Similarly, Agarwal et al. (2024) analyzes the performance of tasks beyond classification in the many-shot context, using Gemini-Pro, and finds that additional demonstrations generally enhance task performance. Furthermore, Li et al. (2024) propose a long-context evaluation benchmark LongICLBench comprising many-shot ICL classification tasks, noting that current long-context models still face challenges in this area. None of the prior works has studied what skill each ICL task measures LCLMs for. LongICLBench mostly focuses on classification tasks, which may only evaluates the retrieval ability of LCLMs. Unlike previous studies, our work provides a more comprehensive analysis of many-shot ICL across a diverse set of tasks and multiple models. We introduce novel metrics to measure retrieval skills required for each task. We identify a set of ICL tasks suitable for evaluation and present a refined long-context evaluation benchmark with fine-grained categorization based on required similar-sample learning and all-sample learning skills. More related work on ICL can be found in Appendix A.

3 Experiment Setting

We select 12 models ranging from 3.8B to 123B parameters and our evaluation includes 12 datasets with 21 subtasks, spanning classification, summarization, reasoning, and translation domains. For each task, we randomly sample 200 data points from the test set, using the full test set if it contains fewer than 200 samples.

For each task, we construct prompts for different

context window sizes by incrementally adding new demonstrations from the training set to the prompt of the shorter context window size and duplicate training examples if they are insufficient to fill the context window. To ensure a fair comparison, we randomize the order of the demonstrations and consistently use the same set of examples across all context sizes. For simplicity, we apply greedy decoding across all models and conduct each experiment using three different random seeds. For prompt construction, we only include demonstrations and provide minimal task instruction. See Table 4 for more details.

3.1 Datasets

We include five datasets for **classification** tasks: BANKING77, GoEmotions, DialogRE, TREC, and CLINC150. For the **summarization** task, we use XLSUM, and for **translation**, we use FLORES-200. Additionally, we incorporate four datasets for **reasoning** tasks: MATH, BBH, and GPQA, and ARC. More details about each dataset can be found in Table 1 and B.

For the MATH, BBH, GPQA, and ARC tasks, we use accuracy as the evaluation metric. Macro F1-score is employed as the metric for all classification tasks. Rouge-L (Lin, 2004) is used for the XLSUM summarization task. ChrF (Popović, 2015) is applied for translation evaluation.

3.2 Models

The list of models we use in our experiment is: Llama-3.1 8B and 70B (Dubey et al., 2024), GLM-4-9B-Chat (GLM et al., 2024), Mistral Nemo (12B) and Large (123B) (Mistral AI, 2024), Qwen2 7B and 72B (Yang et al., 2024), Phi-3 mini (3.8B), small (7B), and medium (14B) (Abdin et al., 2024), and Jamba 1.5 Mini (12B/52B) (Team et al., 2024), and Gemini-1.5-Pro (Team, 2024). We use the instruction-tuned version of all the models. For models with more than 50B, we run the quantized version of the models.¹ More details about each model can be found in Appendix C.

4 Preliminary Study on many-shot ICL

In this section, we explore the extent to which many-shot ICL enhances model performance across different task types. Previous work has either focused on only classification tasks (Bertsch et al., 2024) or studied only one specific model

¹In Appendix D, we show that quantized and unquantized models roughly exhibit the same trend.

Dataset	Task Category	Avg. Tokens / Shot	Max # of Shots	# of Tasks
BANKING77	Intent Classification	13.13	5386	1
GoEmotions	Emotion Classification	15.85	5480	1
DialogRE	Relation Classification	233.27	395	1
TREC	Question Classification	11.25	6272	1
CLINC150	Intent Classification	8.95	7252	1
MATH	Math reasoning	[185.52, 407.90]	[286, 653]	4
GSM8K	Math reasoning	55.78	784	1
BBH	Reasoning	[48.27, 243.01]	[406, 2660]	4
GPQA	MQ - Science	[183.55, 367.02]	[314, 580]	1
ARC	MQ - Science	[61.54, 61.54]	[1997, 2301]	2
XLSUM	New Summarization	621.32	220	1
FLORES-200	Translation	[63.63, 101.74]	[570, 1965]	3

Table 1: Dataset information. GPT-4o tokenizer is used to calculate # of tokens. Max # of shots is the number of shots can be fitted into the 128k context window. For datasets that have multiple subtasks, we list the range for each value. We have 21 tasks in total.

(Agarwal et al., 2024). In contrast, our analysis provides a comprehensive evaluation of many-shot ICL across both classification and generation tasks using ten open-weights LCLMs, excluding Mistral-Large and Gemini-1.5-Pro in this section. We collect tasks from previous work (Bertsch et al., 2024; Agarwal et al., 2024; Li et al., 2024) from six categories: classification, translation, summarization, math reasoning, science reasoning, and symbolic reasoning. The results, illustrated in Figure 1, include aggregated model performance across task types and the correlation coefficients between context lengths and performance from 1k to 64k. We also plot models’ performance on individual tasks in Appendix H and present more analysis on task categories in Appendix D.

Classification performance steadily improves with more shots. Figure 1a demonstrates a consistent performance increase across all models as more demonstrations are added for classification tasks. This trend indicates a strong positive correlation between context length and performance, which is illustrated in Figure 1b. Given that some classification tasks often involve extensive label spaces, e.g., CLINC150 has 150 classes, additional demonstrations provide models with exposure to more classes and thus enhance their ability to perform accurately. This is consistent with prior research findings (Bertsch et al., 2024).

Summarization shows gradual performance gains only. On summarization, most models exhibit a high correlation between context length and performance. However, there is a noticeable slowdown in the performance gains as the number of demonstrations increases. This suggests that

while additional context may improve performance, it does so at a diminishing rate, particularly for smaller models like Llama-3.1-8B that struggle to leverage longer contexts effectively.

Inconsistent trends in science and symbolic tasks. For science and symbolic reasoning tasks, the performance trends are less predictable, with some models displaying minimal changes when seeing additional examples, while others benefit. *This variability suggests that not all tasks lend themselves to the advantages of many-shot ICL equally.* Ideally, for every task, additional demonstrations should either improve performance or, at the very least, not harm it. A model with robust long-context capabilities should exhibit a non-decreasing performance trend as the context length increases. Given the inconsistent performance on non-classification tasks and even decreasing performance on some reasoning tasks, in the next section, we further investigate what aspects these datasets evaluate and identify a set of tasks useful for evaluating important skills of LCLMs.

5 Task Categorization: Similar-Sample Learning vs. All-Sample Learning

To understand what skill each ICL task primarily requires from LCLMs, in this section, we measure the sample learning ratio of each task and divide them into *similar-sample learning* vs. *all-sample learning* tasks. We exclude translation tasks as tokenization can be inconsistent across languages and to avoid difference in multilingual capabilities affecting the results. Similarly, we omit GoEmotions because of its subjective nature.

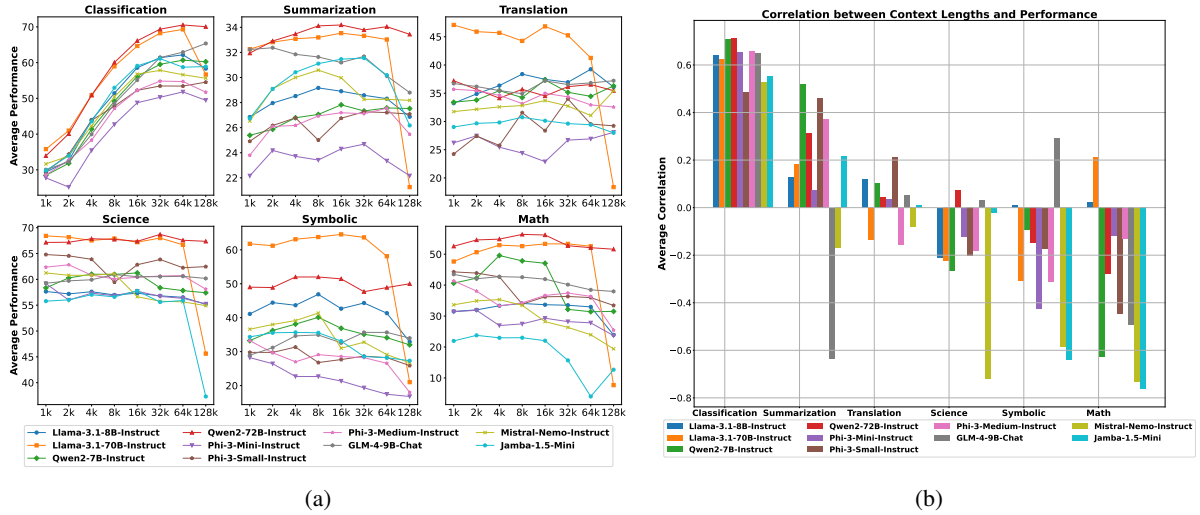


Figure 1: (a) Aggregated performance of models over datasets in different categories of tasks. (b) Average pearson correlation coefficient between context lengths (1k to 64k) and the corresponding performance.

5.1 Sample Learning Ratio

To identify similar-sample learning tasks, we propose a simple metric, **sample learning ratio (SLR)**, to assess whether tasks predominantly rely on models to retrieve relevant examples during many-shot ICL. Concretely, for each ICL task, we create two variants of the original demonstrations at each context size ranging from 1k to 64k by removing the 10% most similar and the 10% least similar examples. At context length l , the model’s performance on these variants is then evaluated, and we have $\text{Perf}_{\text{most}}^{(l)}$ for removing similar examples and $\text{Perf}_{\text{least}}^{(l)}$ for removing dissimilar examples. Here we use BM25 retriever to calculate the similarity. We then average the ratios between $\text{Perf}_{\text{least}}^{(l)}$ and $\text{Perf}_{\text{most}}^{(l)}$ for $l = 1k$ to $l = 64k$ as:

$$\text{SLR} = \frac{1}{7} \sum_{l=1k}^{64k} \frac{\text{Perf}_{\text{least}}^{(l)}}{\text{Perf}_{\text{most}}^{(l)}} \quad (1)$$

Intuitively, if a model predominantly relies on retrieval for a task, removing most similar examples will result in a more pronounced performance drop compared to removing dissimilar ones, which causes the ratio to be larger than 1. Conversely, if there is minimal difference between the two, it means the model does not retrieve similar examples to perform the task, and the ratio will be close to 1.

Classification tasks require retrieval of similar examples. As shown in Figure 2, *all classification tasks exhibit high SLR across the six models*. The BBH geometric shapes task also shows a high SLR, indicating that tasks like BANKING77,

CLINC150, and TREC50 demand strong retrieval capabilities from the models. DialogRE has a relatively lower SLR, suggesting it requires moderate retrieval skills. Among the symbolic tasks, BBH-geometric_shapes is the only reasoning task that has a high SLR. This task involves determining the geometric shape given a full SVG path element, making it similar to a classification task. The high SLR of classification tasks can possibly explain the largest positive correlation between performance and context lengths, as displayed in Figure 1b.

All-sample learning tasks. In Figure 2, tasks such as the math problems and summarization, Dyck languages, translation error detection from BBH, and GPQA with explanations all have a low SLR. This means that *they necessitate a greater degree of global context understanding rather than relying on the retrieval of relevant examples*. These tasks often involve complex reasoning challenges, for which models may lack pretraining skills to solve perfectly, underscoring the need for additional demonstrations or deeper task comprehension.

Additional analysis of SLR. To ensure that performance loss is not caused by the absence of certain labels, we conduct an additional experiment in which the top k% most similar examples are replaced with the most dissimilar examples that share the same labels. In addition to BM25, we also use the sentence transformer (Reimers and Gurevych, 2019) as the retriever. The result of this experiment exhibits the same trend: classification tasks show a higher SLR than non-classification tasks. More details can be found Appendix F.

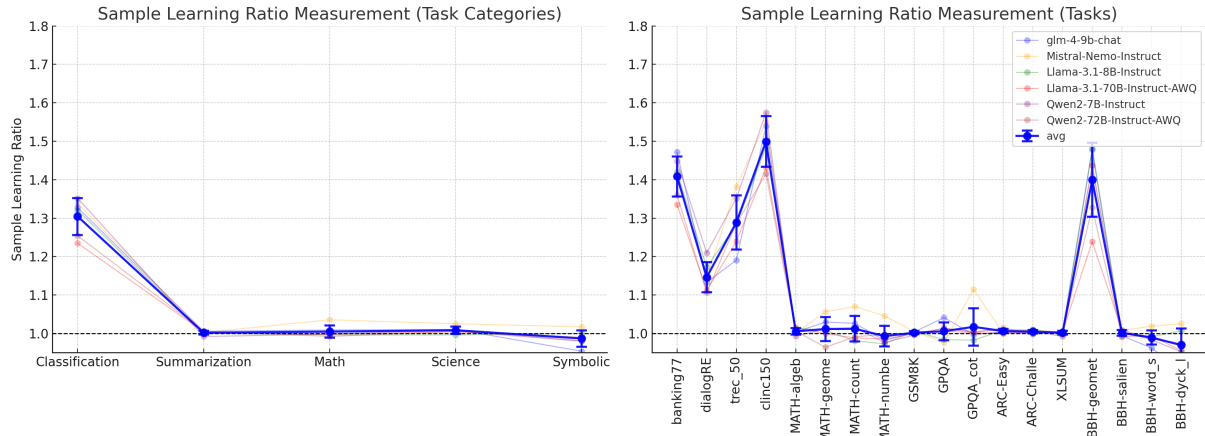


Figure 2: Sample Learning Ratio on different categories of tasks from 1k to 64k tokens. The ratio of 1 indicates models are not doing retrieval during ICL. Classification is the only category of tasks that has a very high ratio, which means classification tasks requires similar-sample retrieval during ICL. The rest of tasks is close to 1, and models’ performance on these tasks does not rely on retrieving similar examples.

6 ManyICLBench: A Many-shot ICL Benchmark to measure retrieval skill and global context understanding

In this section, we present a new long-context benchmark MANYICLBENCH, designed to evaluate LCLMs’ retrieval skills and global context understanding capabilities using the ICL setup. Based on the results from Section 5, we group tasks into two types:

- **5 SSL Tasks:** BANKING77, dialogRE, TREC50, CLINC150, and the geometric shape task from BBH.

- **11 ASL Tasks:** all math tasks, summarization task, GPQA with explanations, ARC_challenge, and all BBH tasks except geometric shapes.

We exclude ARC_Easy and GPQA since ARC_Challenge and GPQA_CoT are already included in the benchmark. Evaluation results of popular LCLMs are summarized in Table 2.

Most models struggle at retrieving examples after 32k length. Up to a context length of 16k, all models demonstrate a steady performance increase, indicating effective retrieval from shorter contexts. However, performance begins to decline after reaching 32k tokens, particularly for the Mistral family and Jamba models. After 64k, the Llama 3.1 family and the mini and medium versions of Phi-3 exhibit a notable downgrade in performance. In contrast, the Qwen-2 family maintains robust performance, with minimal degradation from 64k to 128k. Remarkably, only GLM-4 continues to improve in retrieval performance beyond 64k, indicating its impressive retrieval capabilities within

a very long context window, while larger models such as Mistral-Large and Llama-3.1-70B exhibit the most significant performance losses as context length increases, suggesting that size alone does not ensure superior long-context retrieval ability.

Challenges in ASL tasks. ASL tasks prove to be more challenging, with many models struggling even at short context lengths like 2k or 4k. Only the Llama 3.1 family, Qwen2 family, GLM-4, and Gemini models effectively leverage many demonstrations up to 16k. At 32k, only the Llama 3.1 and Gemini models sustain performance. As context length extends from 32k to 128k, all models experience performance degradation, highlighting that current architectures still struggle to grasp global context and utilize demonstrations effectively. Notably, Qwen2-72B, GLM-4, and Gemini are the only models that do not experience significant performance drops in this category.

The paradox of model size. Despite the common assumption that larger models possess greater capabilities, our findings illustrate that larger models can experience more substantial performance losses compared to smaller models if not trained adequately on long-context data. For instance, Mistral-Large (123B) shows optimal performance from 1k to 32k but experiences a dramatic drop beyond 32k, which is worse than Phi-3-Mini (3.8B). A similar trend is observed with Llama-3.1-70B at 128k. Both underscore the importance of targeted training for long-context tasks. Although the large models we test are quantized, we believe the trend will hold even for the full-precision models.

SSL Tasks	1k	2k	4k	8k	16k	32k	64k	128k	AVG.	AVG.L.
GLM-4-9b-Chat	31.63	34.99	46.37	57.27	63.61	68.34	72.16	72.93	55.91	71.14
Mistral-Nemo-Instruct	33.44	35.45	48.17	57.95	65.38	65.49	63.61	61.73	53.90	63.61
Mistral-Large-Instruct-AWQ	49.15	51.23	60.78	71.95	77.10	79.45	77.77	61.89	66.16	73.04
Llama-3.1-8B-Instruct-AWQ	32.13	34.63	45.76	57.39	66.18	70.02	70.55	65.85	55.31	68.81
Llama-3.1-70B-Instruct-AWQ	38.75	42.87	53.98	66.07	73.12	76.56	78.48	65.56	61.92	73.53
Qwen2-7B-Instruct-AWQ	30.18	34.03	44.40	54.85	62.92	65.91	66.94	66.38	53.20	66.41
Qwen2-72B-Instruct-AWQ	36.41	41.89	54.24	65.33	73.39	76.53	77.51	77.47	62.85	77.17
Phi-3-Mini-Instruct	30.27	30.90	38.09	48.14	53.58	57.29	56.83	48.72	45.48	54.28
Phi-3-Medium-Instruct	31.73	33.55	39.10	49.83	58.29	61.17	60.63	45.32	47.45	55.70
Phi-3-Small-Instruct	31.48	36.27	46.20	54.34	59.63	59.73	60.20	48.97	49.60	56.30
Jamba-1.5-Mini	32.10	36.91	48.61	60.29	66.05	68.33	66.02	65.17	55.44	66.51
Gemini-1.5-Pro	36.40	47.31	58.01	65.49	71.43	74.22	72.43	72.42	62.21	73.03
ASL Tasks	1k	2k	4k	8k	16k	32k	64k	128k	AVG.	AVG.L.
GLM-4-9b-Chat	40.51	40.28	42.04	42.78	40.70	40.46	38.85	39.13	40.59	39.48
Mistral-Nemo-Instruct	38.25	39.07	39.28	38.99	33.06	32.83	30.46	27.11	34.88	30.13
Mistral-Large-Instruct-AWQ	61.47	61.10	61.23	60.87	60.86	58.84	50.01	16.69	53.88	41.85
Llama-3.1-8B-Instruct	37.31	38.84	41.25	40.79	39.83	39.77	39.12	34.41	38.92	37.77
Llama-3.1-70B-Instruct-AWQ	53.32	54.84	55.76	55.87	56.42	56.34	54.42	18.58	50.69	43.12
Qwen2-7B-Instruct	39.52	41.96	45.17	45.39	45.50	37.29	36.97	33.99	40.72	36.09
Qwen2-72B-Instruct-AWQ	48.01	49.24	50.32	50.70	50.97	48.20	47.98	48.16	49.20	48.11
Phi-3-Mini-Instruct	33.54	32.97	29.80	29.75	30.12	28.78	28.06	25.76	29.85	27.53
Phi-3-Medium-Instruct	41.59	40.91	34.85	35.63	36.91	36.84	36.38	28.31	36.43	33.84
Phi-3-Small-Instruct	41.61	41.61	41.61	35.58	37.17	37.73	36.91	35.33	38.44	36.65
Jamba-1.5-Mini	31.96	33.08	32.97	32.70	31.66	28.82	27.14	25.87	30.53	27.28
Gemini-1.5-Pro	57.87	63.39	64.15	66.78	68.02	67.78	66.14	66.42	65.07	66.78

Table 2: Model performance on SSL and ASL tasks. AVG. is the average model performance of all context lengths. AVG.L. is the average model performance of 32k, 64k and 128k. Red indicates performance improvement compared to 1k. Blue indicates performance downgrade compared to 1k. A darker color means higher improvement or downgrade. BOLD number means the largest number of a column. Many models start downgrading their performance after 32k on SSL tasks. On global context ASL tasks, many models start struggling even before 16k.

Llama 3.1 performance and training limitations.

The Llama 3.1 models initially capitalize on additional demonstrations effectively up to 64k but suffer significant performance declines at 128k. This pattern aligns with trends observed in other long-context evaluation benchmarks (Hsieh et al., 2024). We suspect that these performance drops are linked to insufficient training with long-context data during the supervised fine-tuning (SFT) stage. According to Dubey et al. (2024), the average token count for long-context datasets is around 38k, indicating limited exposure for models to effectively learn from data points at 128k lengths.

Gemini is robust. Similar to other open-weight models on SSL tasks, Gemini-1.5-Pro begins to show performance degradation beyond 32k. However, it is one of only three models (alongside Qwen-2-72B and GLM-Chat-9B) that demonstrate impressive retrieval capabilities beyond 64k and maintain performance at 128k. On ASL tasks, Gemini-1.5-Pro significantly outperforms other open-weight models, showcasing its ability to grasp

the global context and effectively utilize all the demonstrations.

Additional analysis on all-sample learning. We want to explore whether ASL tasks genuinely benefit from additional demonstrations and assess models’ global context understanding skills. We compare the performance of models with unique demonstrations versus duplicated examples on ASL tasks and find that some tasks do not benefit from additional unique demonstrations, which means models do not utilize all demonstrations during ICL. More analysis can be found in Appendix G.

6.1 Error Analysis

To better understand how LCLMs fail on ASL tasks, we analyze the performance of Llama-3.1-70B-Instruct and compare it with Qwen-2-72B-Instruct. We focus on non-multiple-choice tasks, including four subtasks from MATH, two symbolic tasks from BBH (word_sorting and dyck_languages), and XLSUM. All the error analysis examples can

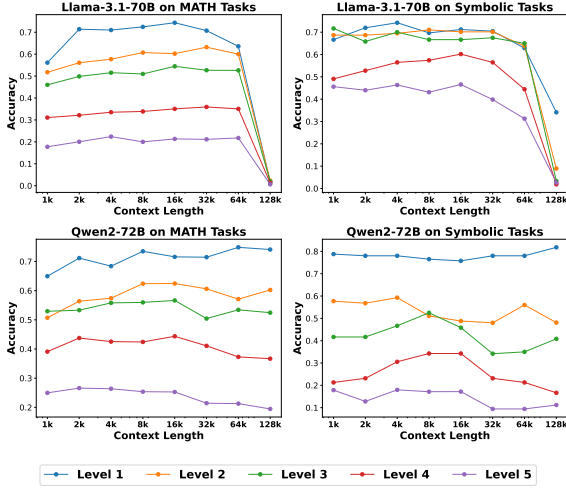


Figure 3: Llama-3.1-70B-Instruct and Qwen2-72B-Instruct aggregated performance on four MATH tasks and two symbolic tasks from Section 6. For the MATH tasks, we use the dataset’s difficulty labels. For the symbolic tasks, we approximate difficulty based on word or symbol length.

be found in Appendix I.

LCLMs’ reasoning ability downgrades at long context lengths. As shown in Figure 3, when the context length increases, Llama-3.1-70B-Instruct tends to lose performance more significantly than Qwen-2-72B-Instruct. Across both easier and more difficult tasks, Qwen-2-72B-Instruct maintains relatively stable performance. A notable observation is that for both models, more difficult tasks degrade at an earlier context length, while easier tasks show a drop only at a later point. For instance, on symbolic tasks, Llama-3.1-70B-Instruct starts losing accuracy on level-4 and level-5 tasks before the simpler level-1 task.

Llama-3.1-70B-Instruct loses reasoning capability at 128k length. On MATH tasks at 128k tokens, Llama-3.1-70B-Instruct often produces repetitive reasoning and gets stuck in loops. Its chain of thought becomes oversimplified, frequently relying on basic sentence structures. In Table 5, the solution at 128k tokens contains many repeated segments, whereas Qwen-2-72B-Instruct still generates valid, consistent reasoning. Furthermore, Llama-3.1-70B-Instruct’s mathematical accuracy suffers: in Table 6, although it correctly calculates $\binom{9}{2} = 36$ at shorter lengths, it incorrectly outputs 42 at 128k tokens.

On symbolic tasks, as test examples become more complex with additional symbols and words,

Llama-3.1-70B-Instruct fails to solve problems it could handle at shorter contexts. While Qwen-2-72B-Instruct still performs well on easier examples at 128k tokens, Llama-3.1-70B-Instruct shows a noticeable drop.

Weaker instruction-following at 128k length.

For summarization tasks, the summaries generated by Llama-3.1-70B-Instruct at 128k tokens become longer and often contain details not present in the reference. In Table 14, Qwen-2-72B-Instruct maintains concise, one-sentence outputs at both 16k and 128k tokens, but Llama-3.1-70B-Instruct adds extraneous information and deviates from the requested one-sentence format. Figure 15 also illustrates a significant increase in output length for Llama-3.1-70B-Instruct at 128k tokens.

Overly long generation at 128k length.

Figure 16 shows that Llama-3.1-70B-Instruct’s chain-of-thought for MATH tasks grows excessively at 128k tokens, partly due to repetitive loops. This trend indicates that beyond a certain context length, the model’s reasoning becomes unfocused and verbose.

7 Conclusion

We investigated many-shot in-context learning (ICL) across various tasks using different open-weight models, assessing their suitability for evaluating long-context language models (LCLMs). Our findings indicate that classification and summarization tasks consistently benefit from additional demonstrations, while other tasks do not. To identify a set of tasks suitable for long-context evaluation, we introduced the concept of SLR to assess the retrieval demands of different tasks. This analysis revealed that classification tasks predominantly rely on the model’s retrieval capabilities. Based on these insights, we categorized tasks into two distinct groups: SSL tasks and ASL tasks. Furthermore, we introduced a novel many-shot ICL benchmark, **ManyICLBench**, designed to evaluate both retrieval and global context understanding skills of LCLMs. Benchmarking open-weight LCLMs on ManyICLBench revealed that most models struggle with ASL tasks at lengths beyond 16k tokens. In contrast, performance on SSL tasks tends to decline after 32k tokens.

8 Limitation

Our study focuses solely on in-context learning setup, and does not address other setups where long context abilities are important, e.g., retrieval-augmented generation. In addition, we restrict our evaluation and benchmark to English ICL tasks, as most existing long-context benchmarks are English-centric. Expanding the benchmark to include multilingual tasks is an important direction for future work. Lastly, given the rapid development and iteration of LCLMs, our findings may not fully generalize to newer models that may be trained to handle long context differently from existing ones. Nevertheless, we believe our findings are useful for future efforts on building models with strong long context understanding.

References

Marah Abidin, Jyoti Aneja, Hany Awadalla, Ahmed Awadallah, Ammar Ahmad Awan, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat Behl, Alon Benhaim, Misha Bilenko, Johan Bjorck, Sébastien Bubeck, Martin Cai, Qin Cai, Vishrav Chaudhary, Dong Chen, Dongdong Chen, Weizhu Chen, Yen-Chun Chen, Yi-Ling Chen, Hao Cheng, Parul Chopra, Xiyang Dai, Matthew Dixon, Ronen Eldan, Victor Fragoso, Jianfeng Gao, Mei Gao, Min Gao, Amit Garg, Allie Del Giorno, Abhishek Goswami, Suriya Gunasekar, Emman Haider, Junheng Hao, Russell J. Hewett, Wenxiang Hu, Jamie Huynh, Dan Iter, Sam Ade Jacobs, Mojan Javaheripi, Xin Jin, Nikos Karampatziakis, Piero Kauffmann, Mahoud Khademi, Dongwoo Kim, Young Jin Kim, Lev Kurilenko, James R. Lee, Yin Tat Lee, Yuanzhi Li, Yunsheng Li, Chen Liang, Lars Liden, Xihui Lin, Zeqi Lin, Ce Liu, Liyuan Liu, Mengchen Liu, Weishung Liu, Xiaodong Liu, Chong Luo, Piyush Madan, Ali Mahmoudzadeh, David Majercak, Matt Mazzola, Caio César Teodoro Mendes, Arindam Mitra, Hardik Modi, Anh Nguyen, Brandon Norick, Barun Patra, Daniel Perez-Becker, Thomas Portet, Reid Pryzant, Heyang Qin, Marko Radmilac, Liliang Ren, Gustavo de Rosa, Corby Rosset, Sambudha Roy, Olatunji Ruwase, Olli Saarikivi, Amin Saied, Adil Salim, Michael Santacrose, Shital Shah, Ning Shang, Hiteshi Sharma, Yelong Shen, Swadheen Shukla, Xia Song, Masahiro Tanaka, Andrea Tupini, Praneetha Vaddamanu, Chunyu Wang, Guanhua Wang, Lijuan Wang, Shuohang Wang, Xin Wang, Yu Wang, Rachel Ward, Wen Wen, Philipp Witte, Haiping Wu, Xiaoxia Wu, Michael Wyatt, Bin Xiao, Can Xu, Jiahang Xu, Weijian Xu, Jilong Xue, Sonali Yadav, Fan Yang, Jianwei Yang, Yifan Yang, Ziyi Yang, Donghan Yu, Lu Yuan, Chenruidong Zhang, Cyril Zhang, Jianwen Zhang, Li Lyna Zhang, Yi Zhang, Yue Zhang, Yunan Zhang, and Xiren Zhou. 2024. [Phi-3 technical report: A highly capable language model locally on your phone](#). *Preprint*, arXiv:2404.14219.

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. 2024. [Many-shot in-context learning](#). *Preprint*, arXiv:2404.11018.

Joshua Ainslie, James Lee-Thorp, Michiel de Jong, Yury Zemlyanskiy, Federico Lebrón, and Sumit Sanghai. 2023. [Gqa: Training generalized multi-query transformer models from multi-head checkpoints](#). *Preprint*, arXiv:2305.13245.

Yushi Bai, Xin Lv, Jiajie Zhang, Yuze He, Ji Qi, Lei Hou, Jie Tang, Yuxiao Dong, and Juanzi Li. 2024. [Longalign: A recipe for long context alignment of large language models](#). *Preprint*, arXiv:2401.18058.

Amanda Bertsch, Maor Ivgi, Uri Alon, Jonathan Berant, Matthew R. Gormley, and Graham Neubig. 2024. [In-context learning with long-context models: An in-depth exploration](#). *Preprint*, arXiv:2405.00200.

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). *Preprint*, arXiv:2005.14165.

Iñigo Casanueva, Tadas Temčinas, Daniela Gerz, Matthew Henderson, and Ivan Vulić. 2020. [Efficient intent detection with dual sentence encoders](#). *Preprint*, arXiv:2003.04807.

Shouyuan Chen, Sherman Wong, Liangjian Chen, and Yuandong Tian. 2023. [Extending context window of large language models via positional interpolation](#). *Preprint*, arXiv:2306.15595.

Alexis Chevalier, Alexander Wettig, Anirudh Ajith, and Danqi Chen. 2023. [Adapting language models to compress contexts](#). *Preprint*, arXiv:2305.14788.

Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv:1803.05457v1*.

Tri Dao, Daniel Y. Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. 2022. [Flashattention: Fast and memory-efficient exact attention with io-awareness](#). *Preprint*, arXiv:2205.14135.

Dorottya Demszky, Dana Movshovitz-Attias, Jeongwoo Ko, Alan Cowen, Gaurav Nemade, and Sujith Ravi. 2020. [Goemotions: A dataset of fine-grained emotions](#). *Preprint*, arXiv:2005.00547.

707	Yiran Ding, Li Lyna Zhang, Chengruidong Zhang, Yuanyuan Xu, Ning Shang, Jiahang Xu, Fan Yang, and Mao Yang. 2024. Longrope: Extending llm context window beyond 2 million tokens . <i>Preprint</i> , arXiv:2402.13753.	763
708		764
709		765
710		766
711		
712	Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. <i>arXiv preprint arXiv:2407.21783</i> .	767
713		768
714		769
715		770
716		
717	Team GLM, :, Aohan Zeng, Bin Xu, Bowen Wang, Chenhui Zhang, Da Yin, Dan Zhang, Diego Rojas, Guanyu Feng, Hanlin Zhao, Hanyu Lai, Hao Yu, Hongning Wang, Jiadao Sun, Jiajie Zhang, Jiale Cheng, Jiayi Gui, Jie Tang, Jing Zhang, Jingyu Sun, Juanzi Li, Lei Zhao, Lindong Wu, Lucen Zhong, Mingdao Liu, Minlie Huang, Peng Zhang, Qinkai Zheng, Rui Lu, Shuaiqi Duan, Shudan Zhang, Shulin Cao, Shuxun Yang, Weng Lam Tam, Wenyi Zhao, Xiao Liu, Xiao Xia, Xiaohan Zhang, Xiaotao Gu, Xin Lv, Xinghan Liu, Xinyi Liu, Xinyue Yang, Xixuan Song, Xunkai Zhang, Yifan An, Yifan Xu, Yilin Niu, Yuantao Yang, Yueyan Li, Yushi Bai, Yuxiao Dong, Zehan Qi, Zhaoyu Wang, Zhen Yang, Zhengxiao Du, Zhenyu Hou, and Zihan Wang. 2024. Chatglm: A family of large language models from glm-130b to glm-4 all tools . <i>Preprint</i> , arXiv:2406.12793.	771
718		772
719		773
720		774
721		775
722		776
723		777
724		778
725		779
726		780
727		781
728		
729		782
730		783
731		784
732		785
733		786
734		787
735		788
736		789
737		790
738		
739	Tahmid Hasan, Abhik Bhattacharjee, Md. Saiful Islam, Kazi Mubasshir, Yuan-Fang Li, Yong-Bin Kang, M. Sohel Rahman, and Rifat Shahriyar. 2021. XLsum: Large-scale multilingual abstractive summarization for 44 languages . In <i>Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021</i> , pages 4693–4703, Online. Association for Computational Linguistics.	791
740		792
741		793
742		794
743		
744	Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021. Measuring mathematical problem solving with the math dataset. <i>NeurIPS</i> .	795
745		796
746		797
747		798
748		
749	Eduard Hovy, Laurie Gerber, Ulf Hermjakob, Chin-Yew Lin, and Deepak Ravichandran. 2001. Toward semantics-based answer pinpointing . In <i>Proceedings of the First International Conference on Human Language Technology Research</i> .	799
750		800
751		801
752	Cheng-Ping Hsieh, Simeng Sun, Samuel Kriman, Shantanu Acharya, Dima Rekesh, Fei Jia, Yang Zhang, and Boris Ginsburg. 2024. Ruler: What’s the real context size of your long-context language models? <i>Preprint</i> , arXiv:2404.06654.	802
753		803
754		804
755		805
756		
757	Huiqiang Jiang, Qianhui Wu, Xufang Luo, Dongsheng Li, Chin-Yew Lin, Yuqing Yang, and Lili Qiu. 2024. Longllmlingua: Accelerating and enhancing llms in long context scenarios via prompt compression . <i>Preprint</i> , arXiv:2310.06839.	806
758		807
759		808
760		
761		809
762	Gregory Kamradt. 2023. Needle in a haystack - pressure testing llms .	810
		811
		812
		813
		814
		815
	Marzena Karpinska, Katherine Thai, Kyle Lo, Tanya Goyal, and Mohit Iyyer. 2024. One thousand and one pairs: A "novel" challenge for long-context language models . <i>Preprint</i> , arXiv:2406.16264.	816
		817
		818
	Wojciech Kryściński, Nazneen Rajani, Divyansh Agarwal, Caiming Xiong, and Dragomir Radev. 2022. Booksum: A collection of datasets for long-form narrative summarization . <i>Preprint</i> , arXiv:2105.08209.	
	Stefan Larson, Anish Mahendran, Joseph J. Peper, Christopher Clarke, Andrew Lee, Parker Hill, Jonathan K. Kummerfeld, Kevin Leach, Michael A. Laurenzano, Lingjia Tang, and Jason Mars. 2019. An evaluation dataset for intent classification and out-of-scope prediction . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 1311–1316, Hong Kong, China. Association for Computational Linguistics.	
	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, Iftexhar Naim, Ming-Wei Chang, and Kelvin Guu. 2024. Can long-context language models subsume retrieval, rag, sql, and more? <i>Preprint</i> , arXiv:2406.13121.	
	Fangyu Lei, Qian Liu, Yiming Huang, Shizhu He, Jun Zhao, and Kang Liu. 2024. S3eval: A synthetic, scalable, systematic evaluation suite for large language models . <i>Preprint</i> , arXiv:2310.15147.	
	Tianle Li, Ge Zhang, Quy Duc Do, Xiang Yue, and Wenhui Chen. 2024. Long-context llms struggle with long in-context learning . <i>Preprint</i> , arXiv:2404.02060.	
	Xin Li and Dan Roth. 2002. Learning question classifiers . In <i>COLING 2002: The 19th International Conference on Computational Linguistics</i> .	
	Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries . In <i>Text Summarization Branches Out</i> , pages 74–81, Barcelona, Spain. Association for Computational Linguistics.	
	Ziqian Lin and Kangwook Lee. 2024. Dual operating modes of in-context learning . <i>Preprint</i> , arXiv:2402.18819.	
	Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2023. Lost in the middle: How language models use long contexts . <i>Preprint</i> , arXiv:2307.03172.	
	Mistral AI. 2024. Mistral nemo. https://mistral.ai/news/mistral-nemo/ . Accessed: 6 September 2024.	
	Amirkeivan Mohtashami and Martin Jaggi. 2023. Landmark attention: Random-access infinite context length for transformers . <i>Preprint</i> , arXiv:2305.16300.	

James Cross Onur Çelebi Maha Elbayad Kenneth	Millican, et al. 2023. Gemini: a family of	876
Heafield Kevin Heffernan Elahe Kalbassi Janice	highly capable multimodal models. <i>arXiv preprint</i>	877
Lam Daniel Licht Jean Maillard Anna Sun Skyler	<i>arXiv:2312.11805</i> .	878
Wang Guillaume Wenzek Al Youngblood Bapi Akula		
Loic Barrault Gabriel Mejia Gonzalez Prangthip	Jamba Team, Barak Lenz, Alan Araz, Amir Bergman,	879
Hansanti John Hoffman Semarley Jarrett Kaushik	Avshalom Manevich, Barak Peleg, Ben Aviram, Chen	880
Ram Sadagopan Dirk Rowe Shannon Spruit Chau	Almagor, Clara Fridman, Dan Padnos, Daniel Gissin,	881
Tran Pierre Andrews Necip Fazil Ayan Shruti Bhos-	Daniel Jannai, Dor Muhlgay, Dor Zimberg, Edden M	882
ale Sergey Edunov Angela Fan Cynthia Gao Vedanuj	Gerber, Elad Dolev, Eran Krakovsky, Erez Safahi,	883
Goswami Francisco Guzmán Philipp Koehn Alexan-	Erez Schwartz, Gal Cohen, Gal Shachaf, Haim	884
dre Mourachko Christophe Ropers Safiyyah Saleem	Rozenblum, Hofit Bata, Ido Blass, Inbal Magar, Itay	885
Holger Schwenk Jeff Wang NLLB Team, Marta R.	Dalmedigos, Jhonathan Osin, Julie Fadlon, Maria	886
Costa-jussà. 2022. No language left behind: Scaling	Rozman, Matan Danos, Michael Gokhman, Mor Zus-	887
human-centered machine translation.	man, Naama Gidron, Nir Ratner, Noam Gat, Noam	888
	Rozen, Oded Fried, Ohad Leshno, Omer Antverg,	889
Bowen Peng, Jeffrey Quesnelle, Honglu Fan, and En-	Omri Abend, Opher Lieber, Or Dagan, Orit Cohavi,	890
rico Shippole. 2023. Yarn: Efficient context win-	Raz Alon, Ro'i Belson, Roi Cohen, Rom Gilad, Ro-	891
dow extension of large language models . <i>Preprint</i> ,	man Glozman, Shahar Lev, Shaked Meirom, Tal Del-	892
<i>arXiv:2309.00071</i> .	bari, Tal Ness, Tomer Asida, Tom Ben Gal, Tom	893
	Braude, Uriya Pumerantz, Yehoshua Cohen, Yonatan	894
Maja Popović. 2015. chrF: character n-gram F-score	Belinkov, Yuval Globerson, Yuval Peleg Levy, and	895
for automatic MT evaluation . In <i>Proceedings of the</i>	Yoav Shoham. 2024. Jamba-1.5: Hybrid transformer-	896
<i>Tenth Workshop on Statistical Machine Translation</i> ,	mamba models at scale . <i>Preprint</i> , <i>arXiv:2408.12570</i> .	897
pages 392–395, Lisbon, Portugal. Association for		
Computational Linguistics.	Sang Michael Xie, Aditi Raghunathan, Percy Liang,	898
	and Tengyu Ma. 2022. An explanation of in-context	899
Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano	learning as implicit bayesian inference . <i>Preprint</i> ,	900
Ermon, Christopher D. Manning, and Chelsea Finn.	<i>arXiv:2111.02080</i> .	901
2024. Direct preference optimization: Your lan-		
guage model is secretly a reward model . <i>Preprint</i> ,	Jing Xu, Arthur Szlam, and Jason Weston. 2021. Be-	902
<i>arXiv:2305.18290</i> .	yond goldfish memory: Long-term open-domain con-	903
	versation . <i>Preprint</i> , <i>arXiv:2107.07567</i> .	904
Nils Reimers and Iryna Gurevych. 2019. Sentence-bert:		
Sentence embeddings using siamese bert-networks .	An Yang, Baosong Yang, Binyuan Hui, Bo Zheng,	905
In <i>Proceedings of the 2019 Conference on Empirical</i>	Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan	906
<i>Methods in Natural Language Processing</i> . Associa-	Li, Dayiheng Liu, Fei Huang, Guanting Dong, Hao-	907
tion for Computational Linguistics.	ran Wei, Huan Lin, Jialong Tang, Jialin Wang,	908
	Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin	909
David Rein, Betty Li Hou, Asa Cooper Stickland,	Ma, Jianxin Yang, Jin Xu, Jingren Zhou, Jinze Bai,	910
Jackson Petty, Richard Yuanzhe Pang, Julien Di-	Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Ke-	911
rani, Julian Michael, and Samuel R. Bowman. 2023.	qin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni,	912
Gpqa: A graduate-level google-proof qa benchmark .	Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize	913
<i>Preprint</i> , <i>arXiv:2311.12022</i> .	Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan,	914
	Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge,	915
Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao,	Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren,	916
Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch,	Xinyu Zhang, Xipin Wei, Xuancheng Ren, Xuejing	917
Adam R Brown, Adam Santoro, Aditya Gupta,	Liu, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan,	918
Adrià Garriga-Alonso, et al. 2022. Beyond the	Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang,	919
imitation game: Quantifying and extrapolating the	Zhifang Guo, and Zhihao Fan. 2024. Qwen2 techni-	920
capabilities of language models. <i>arXiv preprint</i>	cal report . <i>Preprint</i> , <i>arXiv:2407.10671</i> .	921
<i>arXiv:2206.04615</i> .		
Mirac Suzgun, Nathan Scales, Nathanael Schärli, Se-	Dian Yu, Kai Sun, Claire Cardie, and Dong Yu. 2020.	922
bastian Gehrmann, Yi Tay, Hyung Won Chung,	Dialogue-based relation extraction . In <i>Proceedings</i>	923
Aakanksha Chowdhery, Quoc V Le, Ed H Chi, Denny	<i>of the 58th Annual Meeting of the Association for</i>	924
Zhou, , and Jason Wei. 2022. Challenging big-bench	<i>Computational Linguistics</i> , pages 4927–4940, Online.	925
tasks and whether chain-of-thought can solve them.	Association for Computational Linguistics.	926
<i>arXiv preprint arXiv:2210.09261</i> .		
Gemini Team. 2024. Gemini 1.5: Unlocking multi-	Fengji Zhang, Bei Chen, Yue Zhang, Jacky Keung, Jin	927
modal understanding across millions of tokens of	Liu, Daoguang Zan, Yi Mao, Jian-Guang Lou, and	928
context . <i>Preprint</i> , <i>arXiv:2403.05530</i> .	Weizhu Chen. 2023. RepoCoder: Repository-level	929
	code completion through iterative retrieval and gen-	930
Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-	eration . In <i>Proceedings of the 2023 Conference on</i>	931
Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan	<i>Empirical Methods in Natural Language Processing</i> ,	932
Schalkwyk, Andrew M Dai, Anja Hauth, Katie	pages 2471–2484, Singapore. Association for Com-	933
	putational Linguistics.	934

Xinrong Zhang, Yingfa Chen, Shengding Hu, Zihang Xu, Junhao Chen, Moo Khai Hao, Xu Han, Zhen Leng Thai, Shuo Wang, Zhiyuan Liu, and Maosong Sun. 2024. [∞bench: Extending long context evaluation beyond 100k tokens](#). *Preprint*, arXiv:2402.13718.

A Additional Related Work

In-context learning (ICL) enables models to quickly recognize and perform tasks during inference by conditioning on a set of provided demonstrations (Brown et al., 2020). Many previous works have sought to understand the mechanisms behind in-context learning (ICL). Xie et al. (2022) suggests that models implicitly perform Bayesian inference during inference, retrieving relevant skills learned during pretraining. Additionally, Lin and Lee (2024) introduces the concept of a dual operating mode in ICL: task learning and task retrieval. With sufficient demonstrations, models can adapt to unseen tasks learned during pretraining, thereby enhancing performance as the number of demonstrations increases. To explore how many-shot ICL operates, Bertsch et al. (2024) modified the attention patterns by restricting attention among individual examples. Their findings suggest that performance improvements primarily arise from retrieving similar examples rather than comprehending the task. However, their experiment is limited to classification tasks. It may also be biased when comparing full attention and block attention, as block attention allows access to more demonstrations. Our work tries to design better experiments to investigate during many-shot ICL what skill each task mainly requires from LCLMs.

B Datasets

BANKING77 (Casanueva et al., 2020) is an intent classification task in the banking domain. It has over 10k customer service queries labeled with 77 intents.

GoEmotions (Demszky et al., 2020) contains 58 Reddit comments labeled for 27 emotion categories or Neutral.

DialogRE (Yu et al., 2020) is a relation extraction dataset that is built based on transcripts of an American TV show Friends. It comprises 10,168 relation triples for 1,788 dialogues and 36 total relations types. We only focus on relation classification for this dataset.

TREC (Li and Roth, 2002; Hovy et al., 2001) is a question classification dataset with six coarse and 50 fine class labels. It contains 5,500 questions in the training set and 500 in the test set.

CLINC150 (Larson et al., 2019) is an intent classification dataset with 150 intents from 10 domains.

MATH (Hendrycks et al., 2021) is a dataset of 12,5000 challenging completion mathematics prob-

lems. Each problem has a full step-by-step solution. We use four subdomains from the dataset: algebra, geometry, counting and probability, and number theory.

GSM8K (Hendrycks et al., 2021) consists of 8.5K high quality grade school math problems created by human problem writers. These problems take between 2 and 8 steps to solve, and solutions primarily involve performing a sequence of elementary calculations using basic arithmetic operations (+ - / *) to reach the final answer.

BBH (Srivastava et al., 2022) is a subset of 23 challenging BIG-Bench tasks (Suzgun et al., 2022), which include task categories such as mathematics, commonsense reasoning, and question answering. We use four subtasks from BBH-Hard: geometric shape, salient translation error detection, word sorting, and dyck languages.

ARC (Clark et al., 2018) is a dataset of 7,787 genuine grade-school level, multiple-choice science questions. The dataset is partitioned into a Challenge Set and Easy Set, where the former contains only questions answered incorrectly by both a retrieval-based algorithm and a word co-occurrence algorithm.

GPQA (Rein et al., 2023) is a dataset of 448 multiple-choice questions with detailed explanations written by domain experts in biology, physics, and chemistry.

XLSUM (Hasan et al., 2021) is a summarization dataset that focuses on news articles from BBC. In this work, we focus only on English news articles.

FLORES-200 (NLLB Team, 2022) is a translation benchmark that contains many low-resource languages. We follow Agarwal et al. (2024) and choose the translation task from Tamil to English. Additionally, we also test models on Chinese and Spanish.

C Models

Llama-3.1 8B and 70B (Dubey et al., 2024): We use both the 8B and 70B Llama 3.1 Instruction models. These multilingual models are trained on a 128k context window using position interpolation. The models are further fine-tuned with synthetic long-text Supervised Fine-Tuning (SFT) data and also undergo Direct Preference Optimization (DPO) (Rafailov et al., 2024).

GLM-4-9B-Chat (GLM et al., 2024): This is a 9-billion-parameter multilingual model, also trained on a 128k context window with position

interpolation. It is further fine-tuned with labeled long-text SFT data and undergoes a DPO stage.

Mistral Family (Mistral AI, 2024): We use both 12-billion-parameter and 123-billion-parameter multilingual models, trained on a 128k context window.

Qwen2 7B and 72B (Yang et al., 2024): These two models are trained with a context size of 32k tokens, and their context window is extended to 128k by YARN (Peng et al., 2023), a dynamic position interpolation technique.

Phi-3 (Abdin et al., 2024): We use the mini (3.8B), small (7B), and medium (14B) versions of Phi-3 models. They are trained with the context size of 4k tokens on high quality data, and LongRope (Ding et al., 2024) extends their context size to 128k.

Jamba-1.5-Mini (Team et al., 2024): It’s a hybrid SSM-Transformer model with 12B of active parameters and 52B of total parameters with a context size of 256k tokens.

Gemini-1.5-Pro (Team, 2024): It is a commercial model introduced by Google and has a context size of 2 million tokens.

D Quantization vs. Regular

We compare the 4-bit quantized version and unquantized version of both Llama-3.1 8B and Llama-3.1-70B. In both Figure 4 and Figure 5, we can observe that the quantized version experiences a little performance drop but exhibits the same trend as the unquantized version with the increasing context length.

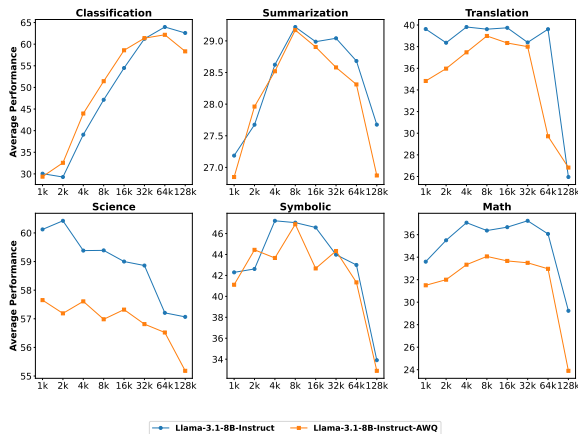


Figure 4: Comparison between Llama-3.1-8B and 4-bit quantized Llama-3.1-8B. There are some performance gaps between two models on translation, science, and math tasks, but with the increasing context size, the performance trend is the same for both models.

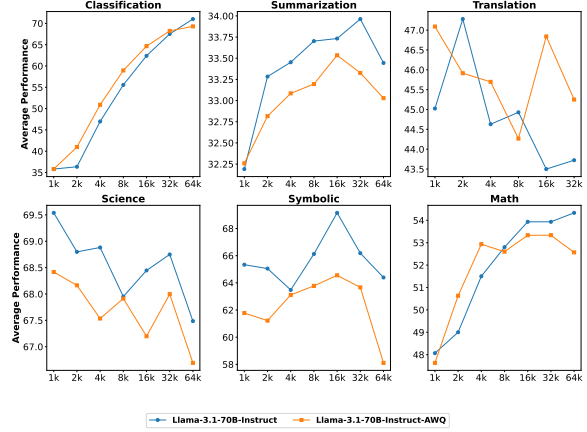


Figure 5: Comparison between Llama-3.1-70B and 4-bit quantized Llama-3.1-70B. Similar to the smaller model, the performance trends hold for both models except the translation tasks. In our benchmark, we exclude all the translation tasks because of the inconsistent multilingual ability of LCLMs.

E More analysis on Section 4 many-shot ICL

Subjective tasks do not benefit from more examples. The GoEmotions task, though being a classification problem, exhibits a fluctuating performance trend across all models with increasing shots in Figure 9. We attribute this inconsistency to the subjective nature of the task, where nuanced emotional categories may lead to low annotator agreement (Demszky et al., 2020). This variance in the annotated labels may results in a weaker correlation between context length and performance. This finding highlights a limitation in using ICL tasks with ambiguous ground truths to evaluate LCLMs, as their performance does not improve with more demonstrations.

Models’ performance fluctuates on translation tasks. As shown in Figure 10, the performance curves for all models across different languages differ. For the low-resource language, models show larger performance gap than those in the high-resource language, e.g., Spanish. In Chinese, models become spikier than in other languages across different context sizes. In Figure 1a, translation tasks show a very flat curve, with no significant improvement as the number of demonstrations increases. This result contrasts with Agarwal et al. (2024), where the Gemini-1.5 Pro model demonstrated consistent performance improvements in Kurdish and Tamil translation tasks as the context size increased. We think the performance incon-

sistency is caused by the mismatched multilingual capability of models and different model sizes.

Math tasks benefit from additional demonstrations, particularly for stronger models. In math reasoning tasks, only the Llama-3.1 and Qwen2 model families show significant performance improvements with additional demonstrations. Notably, Qwen2 performance plateaus at 16k length, while Llama-3.1 continues to improve until 64k. The models with larger parameter sizes tend to exhibit more consistent performance gains, supporting findings from Agarwal et al. (2024) who have demonstrated that Gemini 1.5 Pro improves on math tasks with more examples.

F Sample Learning Ratio with Replacement

To ensure the performance downgrade is not caused by the absence of certain labels in the experiment from Section 5, we replace similar examples with distant examples with the same labels. The new sample learning ratio formula is $\frac{score_{original}}{score_{replace}}$. We use Llama-3.1 and Qwen2 models and conduct this experiment from 1k to 32k with BM25 and from 1k to 32k with SBERT (Reimers and Gurevych, 2019) retrievers.

BM25: The trend in Figure 2 matches the results of Figure 6. All the classification tasks downgrade performance more when similar examples are replaced. However, the degree of downgrade is less significant than removing similar examples.

SBERT: For SentenceTransformer, we use all-MiniLM-L6-v2 as the base model. The trends observed from Figure 2 and Figure 6 still hold in Figure 7. That is, all the classification tasks still have a higher ratio and the non-classification tasks have a ratio close to 1.

G A deeper look into all-sample learning task

In this section, we investigate which ASL tasks tend to benefit more from additional demonstrations and whether models use all the demonstrations to understand the task during ICL.

To that end, we propose another metric, Global Context Index, to measure the global context understanding skill required by a task i.e., skills required to learn from all samples. Specifically, for each ASL task, we create two variants of demonstrations, both starting with the same demonstrations used in the 1k context-length experiment. For

context lengths l ranging from 2k to 64k, the *unique variant* keeps adding unique demonstrations to the prompt, whereas the *duplicate variant* repeats the same demonstrations as in the 1k context. At each length l , the performance of the unique variant is denoted as $\text{Perf}_{\text{unique}}^{(l)}$, and the performance of the duplicate variant as $\text{Perf}_{\text{duplicate}}^{(l)}$. We then average the percentage difference between $\text{Perf}_{\text{unique}}^{(l)}$ and $\text{Perf}_{\text{duplicate}}^{(l)}$ for $l = 2k$ to $l = 64k$ across 6 context lengths as:

$$\text{Global Context Index} = \frac{1}{6} \sum_{l=2k}^{64k} \left(1 - \frac{\text{Perf}_{\text{duplicate}}^{(l)}}{\text{Perf}_{\text{unique}}^{(l)}} \right)$$

If duplicating examples results in worse performance on an ASL task than adding unique examples, the global context index will be positive and suggests that the model benefits more from providing unique demonstrations. This means that performance improvements come from learning from diverse examples rather than simply picking up on formatting patterns or relying on spurious correlations between in-domain tokens and predictions. We use Llama-3.1-70B for the preliminary analysis because it is best at using additional demonstrations out of all models we have tested so far, e.g., it shows a high positive correlation between context lengths and performance in Figure 1b.

In Figure 8, tasks such as the math problems, summarization, Dyck languages, translation error detection from BBH, and GPQA with explanations all have worse performance with duplicated demonstrations. This means that *they necessitate a greater degree of global context understanding rather than relying on the retrieval of similar examples*. These tasks are often complex reasoning challenges, for which models may lack pretraining skills to solve perfectly, underscoring the need for additional demonstrations or deeper task comprehension.

ARC-Easy, ARC-Challenge, GPQA, and BBH word sorting tasks are indifferent to duplicating examples. This indicates that these tasks do not benefit from additional demonstrations. Most of these tasks assess the intrinsic abilities of the models reasoning with their parametric knowledge, thus a few demonstrations suffice. Adding more demonstrations may introduce distractions rather than improve performance. Interestingly, GPQA with “chain-of-thoughts” benefit from additional examples. We suspect that without these solution steps,

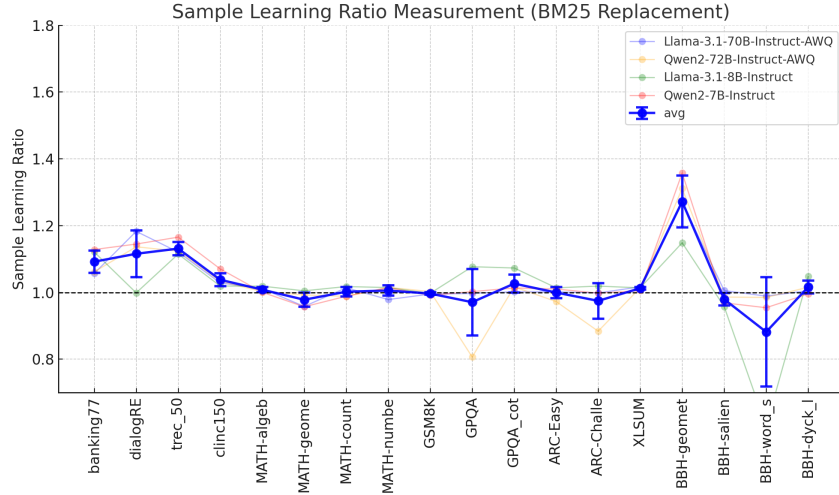


Figure 6: Sample Learning Ratio under the replacement setting with BM25 on all tasks except XLSUM from 1k to 32k tokens. The ratio of 1 indicates models are not doing retrieval during ICL because similar demonstrations don’t help models perform better. Similar to Figure 2, classification is the only category of tasks that has a higher ratio, which means classification tasks largely require model retrieval skills during ICL. The rest of the tasks are close to 1, and the models’ performance on these tasks does not rely on retrieving similar examples.

GPQA is too challenging for the model to understand even after seeing many demonstrations with answers only.

9 Strong ASL Tasks. Based on the analysis of the global context index, we remove ARC-Challenge and BBH word sorting tasks to create a new category of tasks called **strong ASL Tasks**. Evaluation results on strong ASL Tasks are summarized in Table 3.

H Task Performance

In this section, we present the models’ performance on individual tasks and group them by the task categories: classification (Figure 9), translation (Figure 10), summarization (Figure 9), and reasoning (Figures 11, 12, and 13).

I Error Analysis

In this section, we include examples from MATH, BBH, and XLSUM tasks, highlighting the performance of Llama-3.1-70B and Qwen2-72B at both 16k and 128k context lengths. We analyze the strengths and shortcomings of each model’s outputs, focusing on accuracy, consistency, and the impact of extended context lengths on their reasoning processes. MATH (Table 5, 6, and 7), BBH-word_sorting (Table 11, 12, and 13), BBH-dyck_languages (Table 8, 9, and 10), and XLSUM (Table 14).

Qwen2 and GLM-4 show relatively robust capabilities on both tasks. The Qwen2-72B model consistently maintains performance across both SSL and ASL tasks, demonstrating its adaptability for longer contexts. Trained on data with up to 32k tokens, Qwen2 models employ modified RoPE frequency and training-free positional interpolation methods to handle longer contexts. However, the Qwen2 family models drop their performance from 16k to 32k in the ASL tasks but maintain their performance after 32k. This raises the question of whether the training-free length extension methods enable models to use additional demonstrations or merely maintain their performance in the short context length and ignore additional examples during many-shot ICL. Meanwhile, GLM-4-chat also shows a relatively robust performance at a longer context size and is the only model to experience a performance increase from 64k to 128k on SSL tasks. GLM-4’s training methodology closely mirrors that of Llama 3.1 models, with adjustments to the RoPE base and continuous training on long-context data. The difference is, during SFT, GLM-4-9B follows LongAlign (Bai et al., 2024), which determines the length distribution of the long-context SFT data carefully. GLM-4-9B also goes through the RLHF stage with both short and long data.

Does LCLM’s poor performance beyond 16k limit their potential in the world application?

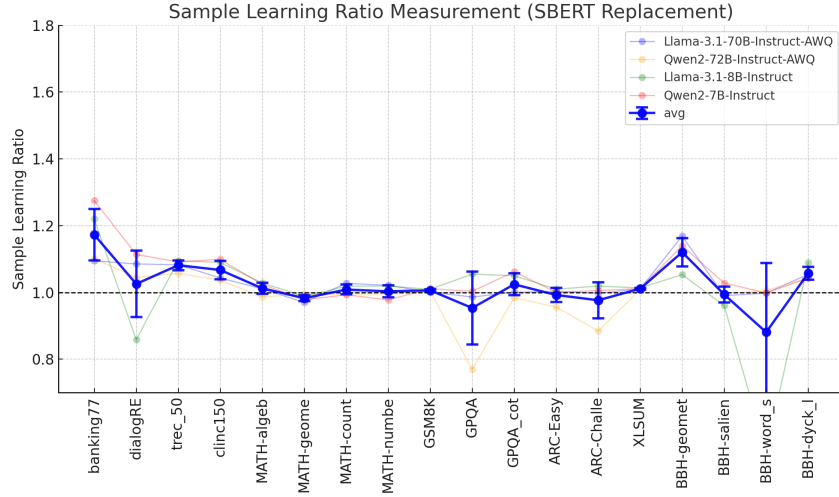


Figure 7: Sample Learning Ratio under the replacement setting with SBERT on selective tasks from 1k to 32k tokens. A ratio of 1 signifies that models do not perform retrieval during in-context learning (ICL), as similar demonstrations do not enhance their performance. As shown in Figure 2, classification tasks are the only category with a higher retrieval load ratio, indicating a strong dependence on retrieval during ICL. In contrast, other tasks exhibit ratios close to 1, suggesting minimal reliance on retrieval, with models’ performance largely unaffected by retrieval-based demonstrations.

In our ASL task evaluation, we distinguish between more practical tasks, such as XLSUM and MATH, and less real-world-relevant tasks, such as symbolic reasoning and ARC. To assess whether poor ASL performance beyond 16K tokens translates to real-world limitations, we analyzed the correlation between these tasks for three models that exhibit significant degradation after 16K tokens. In Figure 14, for Phi-3-Medium and Mistral-Large, we observed strong correlations across all tasks, which likely stems from their uniformly poor performance across the board. In contrast, Llama-3.1-70B exhibits a distinct pattern: while word sorting (a less practical task) does not correlate well with other tasks, all other ASL tasks show strong mutual correlations. Based on these findings, we believe that LCLMs’ poor performance on ASL tasks beyond 16K tokens is a fair indicator of their limitations in real-world applications, where long-context understanding is often essential.

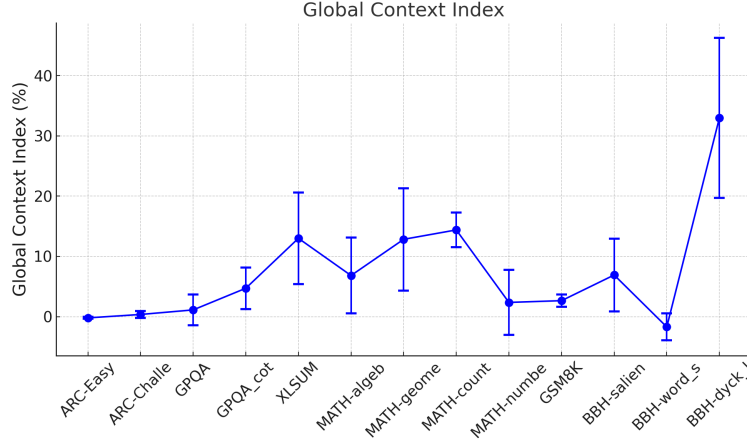


Figure 8: Global context index is the average % difference between adding duplicated vs. unique examples from 2k to 64k context for non-retrieval tasks. 0% means duplicating does not harm the model’s performance. Easy tasks such as ARC and word sorting do not benefit from additional information. When a task is too difficult, e.g., GPQA, the model cannot effectively learn all demonstrations unless explanations are provided.

ASL Tasks	1k	2k	4k	8k	16k	32k	64k	128k	AVG.	AVG.L.
GLM-4-9b-Chat	36.79	36.23	38.30	39.30	37.60	37.94	36.53	35.45	37.27	36.64
Mistral-Nemo-Instruct	33.94	34.88	34.92	34.72	28.22	28.64	26.28	23.23	30.60	26.05
Mistral-Large-Instruct-AWQ	57.09	56.30	56.21	56.12	56.43	53.33	42.98	13.10	48.94	36.47
Llama-3.1-8B-Instruct-AWQ	31.31	32.79	33.02	34.50	34.25	35.22	33.71	27.88	32.84	32.27
Llama-3.1-70B-Instruct-AWQ	45.53	47.60	48.39	49.08	49.64	49.83	47.74	13.88	43.99	37.23
Qwen2-7B-Instruct-AWQ	37.75	39.47	43.86	44.55	42.83	35.17	33.00	32.70	38.67	33.62
Qwen2-72B-Instruct-AWQ	47.38	49.03	50.32	50.69	50.78	48.56	48.18	48.68	49.20	48.47
Phi-3-Mini-Instruct	29.86	29.20	26.61	26.95	27.65	26.34	25.54	23.08	26.90	24.98
Phi-3-Medium-Instruct	37.74	37.15	31.49	32.02	33.04	33.19	33.06	24.56	32.78	30.27
Phi-3-Small-Instruct	38.40	38.40	38.35	31.69	34.04	34.59	33.74	32.46	35.21	33.60
Jamba-1.5-Mini	27.86	29.04	28.93	28.86	27.86	24.92	23.12	22.42	26.63	23.48
Gemini-1.5-Pro	58.26	60.88	61.30	65.20	65.05	65.12	62.38	63.61	62.73	63.70

Table 3: Model performance on strong ASL tasks. AVG. is the average model performance of all context lengths. AVG.L. is the average model performance of 32k, 64k and 128k. **Red** indicates performance improvement compared to 1k. **Blue** indicates performance downgrade compared to 1k. A darker color means higher improvement or downgrade. **BOLD** number means the largest number of a column.

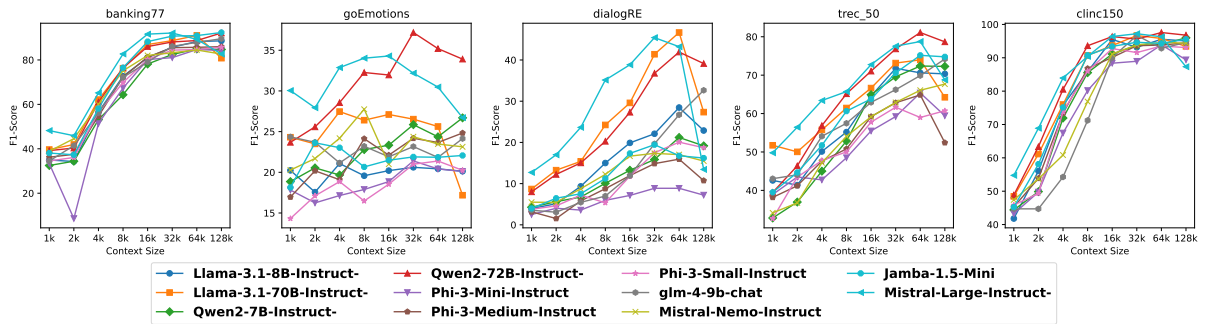


Figure 9: Models’ performance on all classification tasks. All tasks except GoEmotions show a very consistent gain with increasing context size. We excluded GoEmotions from our benchmark because of the data’s strong subjectivity.

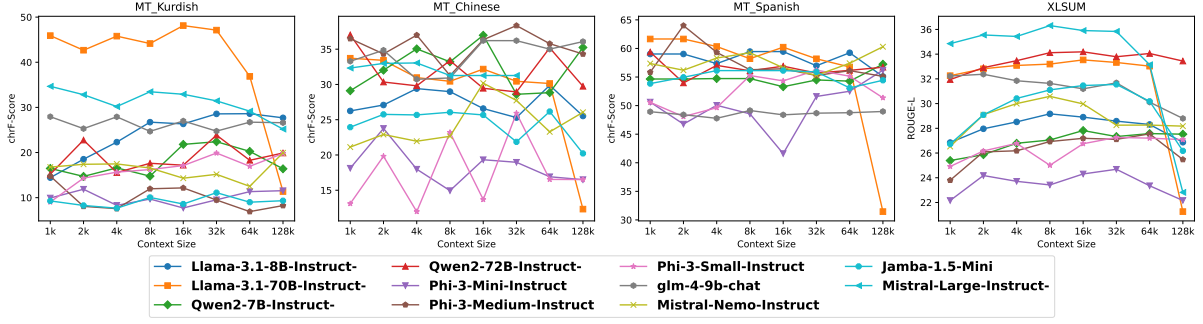


Figure 10: Models' performance on all translation tasks and the summarization task. For translation tasks, we do not observe a clear pattern among different languages and models, which can be caused by LCLMs' different multilingual abilities. We can see a slightly positive trend for the summarization task.

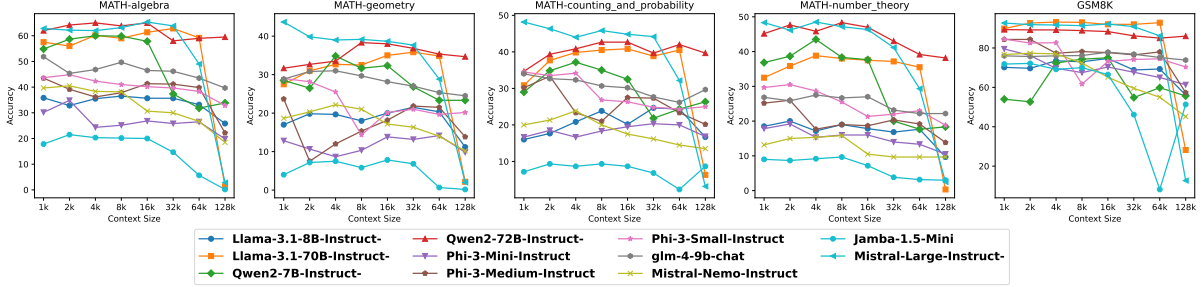


Figure 11: Models' performance on all math tasks. Overall, the larger and stronger models benefit more from the increasing context window size on math tasks.

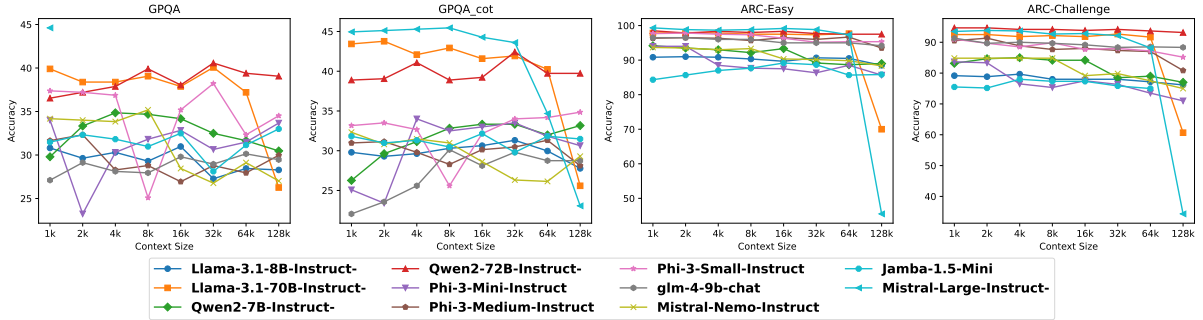


Figure 12: Models' performance on all science tasks. For the ARC task, the performance of all models stays the same across all context sizes. For GPQA, we can see larger and more robust LCLMs keep or increase their performance with the increasing context size.

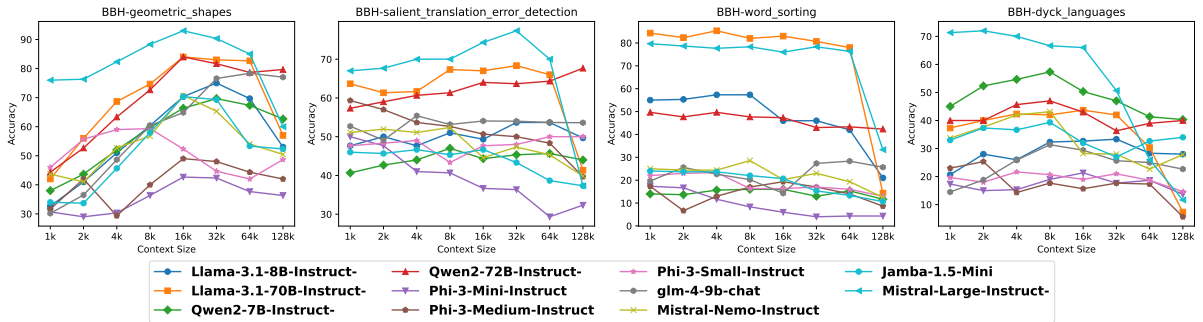
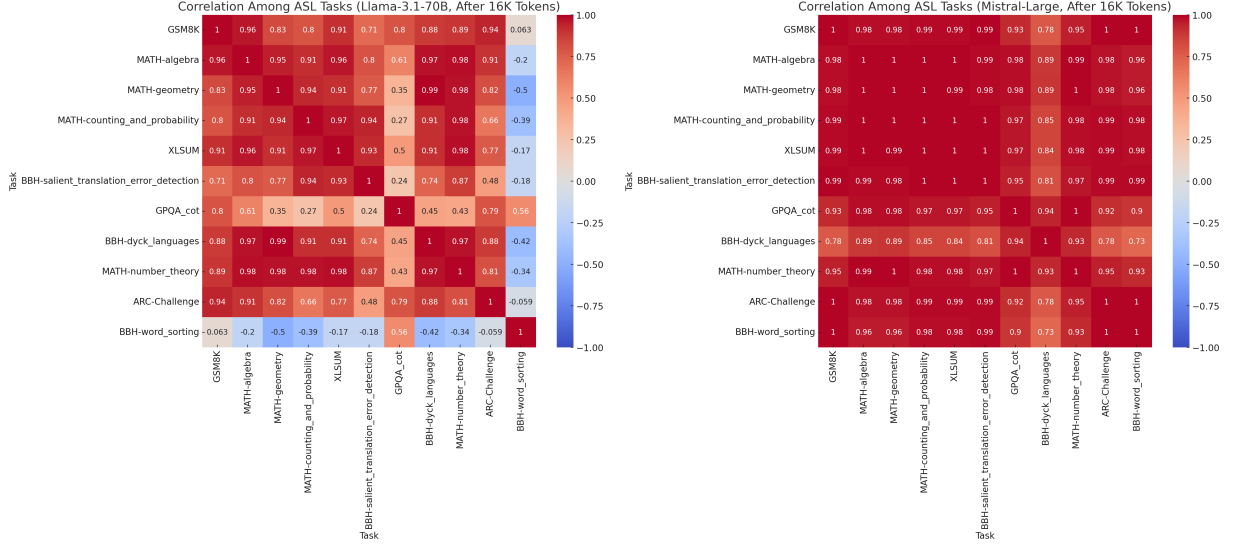
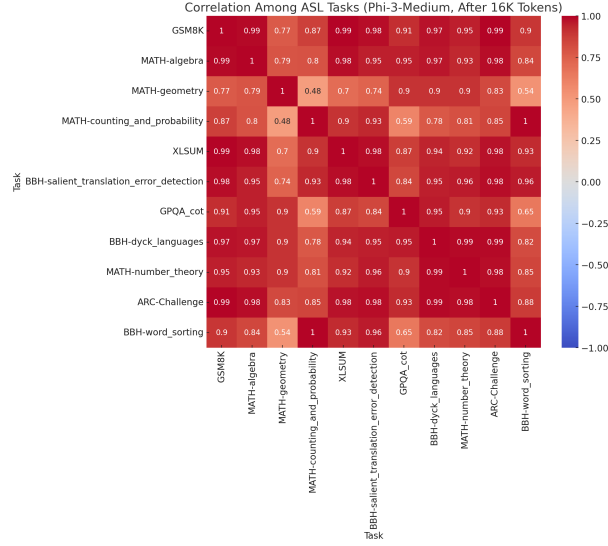


Figure 13: Models' performance on all symbolic tasks. For the geometric shape and translation error detection tasks, we can all model benefit from the increasing context length. We suspect the word sorting task may too easy for the models, so the lines are flat. For the dyck language task, the models experience performance gain up 16k context length but start downgrading afterward.



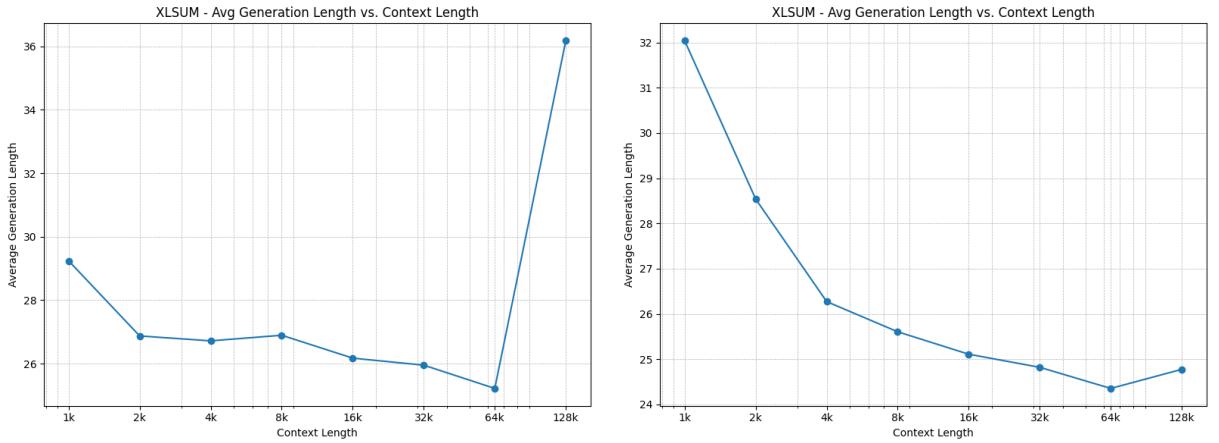
(a) Llama-3.1-70B-Instruct's correlation

(b) Mistral-Large's correlation



(c) Phi-3-Medium's correlation

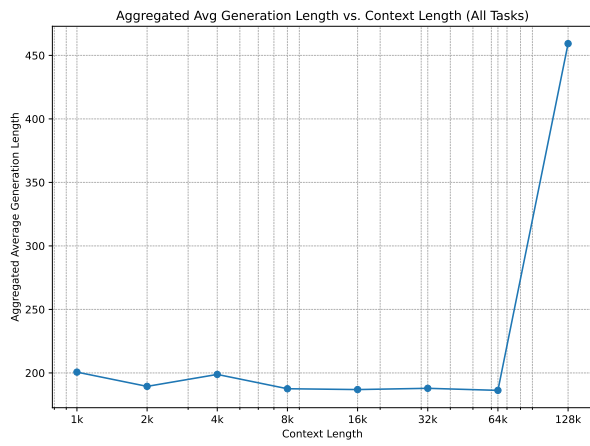
Figure 14: Correlation table among ASL tasks.



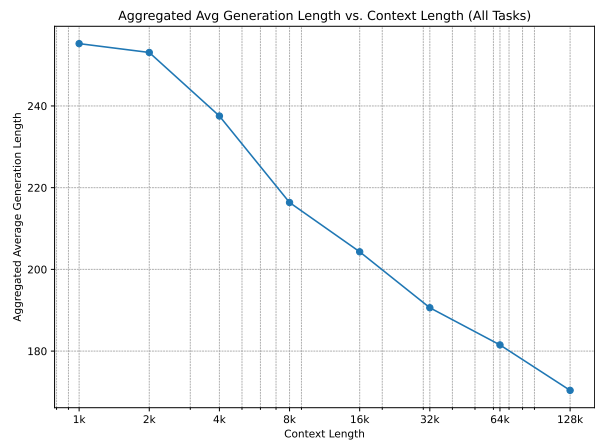
(a) Llama-3.1-70B-Instruct's generation length on XLSUM

(b) Qwen2-72B-Instruct's generation length on XLSUM

Figure 15: XLSUM generation length



(a) Llama-3.1-70B-Instruct's generation length on MATH



(b) Qwen2-72B-Instruct's generation length on MATH

Figure 16: MATH generation length

Task	Example
banking77	Query: I have multiple of the same transaction Intent: transaction_charged_twice
goEmotions	Comment: It is bad tho. Category: disapproval
dialogRE	Dialogue: Speaker 1: I'm divorced! I'm only 26 and I'm divorced! Speaker 2: Shut up! Speaker 3: You must stop! Speaker 1: That only took me an hour. The list of 1 relations are (Speaker 1,26) The respective relations between each entity pair are: per:age
trec_50	Question: Describe the Finnish music personality Salonen 's appearance . Type: DESC:desc
clinc150	Query: i need a good joke about office parties Intent: tell_joke
MATH-algebra	Problem: How many cubic feet are in three cubic yards? Solution: Cubing both sides of 1 yard = 3 feet we find that 1 cubic yard equals 27 cubic feet. Therefore, 3 cubic yards are equal to $27 \cdot 3 = \boxed{81}$ cubic feet.
MATH-geometry	<p>Problem: In quadrilateral $ABCD$, $BC = 8$, $CD = 12$, $AD = 10$, and $m\angle A = m\angle B = 60^\circ$. Given that $AB = p + \sqrt{q}$, where p and q are positive integers, find $p + q$.</p> <p>Solution: <code>[asy]draw((0,0)--(20.87,0)--(15.87,8.66)--(5,8.66)--cycle);</code> <code>draw((5,8.66)--(5,0)); draw((15.87,8.66)--(15.87,0)); draw((5,8.66)--</code> <code>(16.87,6.928)); label("A",(0,0),SW); label("B",(20.87,0),SE);</code> <code>label("E",(15.87,8.66),NE); label("D",(5,8.66),NW); label("P",(5,0),S);</code> <code>label("Q",(15.87,0),S); label("C",(16.87,7),E); label("12",(10.935,7.794),S);</code> <code>label("10",(2.5,4.5),W); label("10",(18.37,4.5),E); [asy]</code> Draw line segment DE such that line DE is concurrent with line BC. Then, $ABED$ is an isosceles trapezoid so $AD = BE = 10$, and $BC = 8$ and $EC = 2$. We are given that $DC = 12$. Since $\angle CED = 120^\circ$, using the Law of Cosines on $\triangle CED$ gives $12^2 = DE^2 + 4 - 2(2)(DE)(\cos 120^\circ),$ which simplifies to $144 - 4 = DE^2 + 2DE$. Adding 1 to both sides yields $(DE + 1)^2 = 141$, so $DE = \sqrt{141} - 1$. In the $30-60-90$ triangles $\triangle DAP$ and $\triangle EBQ$, we have $AP = BQ = 5$. Since $PQ = DE$, it follows that $AB = AP + PQ + BQ = 5 + (\sqrt{141} - 1) + 5 = 9 + \sqrt{141},$ so $(p, q) = (9, 141)$ and $p + q = \boxed{150}$. </p>
MATH-counting	Problem: Sarah wants to order a pizza with 2 different toppings. She has 9 different toppings to choose from. How many different pizzas could she order? Solution: There are 9 choices of toppings, and we need to choose 2 distinct toppings. This is represented by the number of 2-element subsets of a 9-element set. We use the binomial coefficient $\binom{9}{2} = \boxed{36}$ to compute this.
MATH-number_theory	Problem: Express eleven in base 2. Solution: We have $11 = 1 \cdot 2^3 + 0 \cdot 2^2 + 1 \cdot 2^1 + 1 \cdot 2^0$, so $11 = \boxed{1011_2}$.

BBH-geometric_shapes	<p>Input: This SVG path element <path d="M 53.64,29.71 L 61.55,33.55 M 61.55,33.55 L 65.49,42.75 M 65.49,42.75 L 60.75,49.85 M 60.75,49.85 L 54.92,52.75 L 48.80,50.52 M 48.80,50.52 L 44.97,43.03 M 44.97,43.03 L 46.15,33.55 M 46.15,33.55 L 53.64,29.71"/> draws a</p> <p>Options: (A) circle, (B) heptagon, (C) hexagon, (D) kite, (E) line, (F) octagon, (G) pentagon, (H) rectangle, (I) sector, (J) triangle</p> <p>Target: (F)</p>
BBH-salient_translation_error_detection	<p>Input: The following translations from German to English contain a particular error. That error will be one of the following types: Named Entities (an entity is changed), Numerical Values (values or units are changed), Modifiers or Adjectives (modifiers are changed), Negation or Antonyms (negations or opposites are altered), Facts (trivial factual errors), Dropped Content (significant content is removed). Please identify that error.</p> <p>Source: Die unvollständige Liste der Baudenkmale in Barsinghausen enthält Baudenkmale der Barsinghausener Kernstadt sowie der Ortsteile Bantorf, Barrigsen, Eckerde, Egestorf, Göxe, Großgoltern, Groß Munzel, Hohenbostel, Holtensen, Kirchdorf, Landringhausen, Langreder, Nordgoltern, Ostermunzel, Stemmen, Wichtringhausen und Winninghausen.</p> <p>Translation: The complete list of architectural monuments in Barsinghausen contains architectural monuments of the Barsinghausen core town as well as the districts bantorf, Barrigsen, Eckerde, Egestorf, Göxe, Großgoltern, Groß Munzel, Hohenbostel, Holtensen, Kirchdorf, Landringhausen, Langreder, Nordgoltern, Ostermunzel, Stemmen, Wichtringhausen and Winninghausen.</p> <p>The translation contains an error pertaining to</p> <p>Options: (A) Modifiers or Adjectives, (B) Numerical Values, (C) Negation or Antonyms, (D) Named Entities, (E) Dropped Content, (F) Facts</p> <p>Target: (C)</p>
BBH-word_sorting	<p>Input: Sort the following words alphabetically: List: thrill splutter panicking scorch same dot prod obstetric malton onus drumhead delmarva barn embezzle it&t damp guru subsist entirety greene</p> <p>Target: barn damp delmarva dot drumhead embezzle entirety greene guru it&t malton obstetric onus panicking prod same scorch splutter subsist thrill</p>
BBH-dyck_languages	<p>Input: Complete the rest of the sequence, making sure that the parentheses are closed properly. Input: ([[[]] < [< [] >] > ></p> <p>Target:])</p>
GPQA	<p>Question: Determine which set of states mentioned below are only entangled states:</p> <p>(a) $\frac{1}{\sqrt{30}}(00\rangle + 2i 01\rangle - 3 10\rangle - 4i 11\rangle)$</p> <p>(b) $\frac{1}{5}(00\rangle + 2i 01\rangle - 2 10\rangle - 4i 11\rangle)$</p> <p>(c) $\frac{1}{2}(00\rangle + 01\rangle + 10\rangle - 11\rangle)$</p> <p>(d) $\frac{1}{2}(00\rangle + 01\rangle - 10\rangle - 11\rangle)$</p> <p>Options: A. a,b B. b,d C. c,d D. a,c</p> <p>Answer: D</p>

GPQA_cot	<p>Question: Determine which set of states mentioned below are only entangled states:</p> <p>(a) $\frac{1}{\sqrt{30}}(00\rangle + 2i 01\rangle - 3 10\rangle - 4i 11\rangle)$</p> <p>(b) $\frac{1}{5}(00\rangle + 2i 01\rangle - 2 10\rangle - 4i 11\rangle)$</p> <p>(c) $\frac{1}{2}(00\rangle + 01\rangle + 10\rangle - 11\rangle)$</p> <p>(d) $\frac{1}{2}(00\rangle + 01\rangle - 10\rangle - 11\rangle)$</p> <p>Options: A. a,b B. b,d C. c,d D. a,c</p> <p>Answer: D</p> <p>Explanation: For a state $a 00\rangle + b 01\rangle + c 10\rangle + d 11\rangle$, separability requires $a \cdot d = b \cdot c$. This condition is not met for option a,c; hence both states are entangled.</p>
ARC-Challenge	<p>Question: One important difference between living things and nonliving things is that only living things have</p> <p>1. compounds 2. elements 3. molecules 4. cells</p> <p>Answer: 4</p>
ARC-Easy	<p>Question: Infections may be caused by</p> <p>1. mutations 2. microorganisms 3. toxic substances 4. climate changes</p> <p>Answer: 2</p>
MT_Kurdish	<p>English: Every corner store is filled with a confusing array of pre-paid phone cards that can be used from pay phones or ordinary telephones.</p> <p>Kurdish: ...</p>
MT_Chinese	<p>English: Every corner store is filled with a confusing array of pre-paid phone cards that can be used from pay phones or ordinary telephones.</p> <p>Chinese: ...</p>
MT_Spanish	<p>English: Every corner store is filled with a confusing array of pre-paid phone cards that can be used from pay phones or ordinary telephones.</p> <p>Spanish: En las tiendas de cada esquina hay una enorme y confusa variedad de tarjetas de telefonía prepagas que pueden utilizarse en teléfonos públicos o comunes.</p>
GSM8K	<p>Question: If one Burmese python can eat one 50-cm alligator per week, how many Burmese pythons would it take to eat fifteen 50-centimeter alligators in three weeks?</p> <p>Answer: 15 alligators eaten over 3 weeks gives a rate of $15/3 = 5$ alligators per week. Thus, at 1 alligator per week per python, it would take 5 pythons to eat 15 alligators.</p> <div style="border: 1px solid black; padding: 2px; display: inline-block;">5</div>

Table 4: Input and output example from each task

Task: Math-Algebra

Problem:

A triangle has vertices at coordinates (11, 1), (2, 3) and (3, 7). What is the number of units in the length of the longest side of the triangle?

Difficulty: Level 3

Solution:

We must find the distance between each pair of points by using the distance formula. The distance between (11, 1) and (2, 3) is $\sqrt{(11-2)^2 + (1-3)^2} = \sqrt{81+4} = \sqrt{85}$. The distance between (2, 3) and (3, 7) is $\sqrt{(2-3)^2 + (3-7)^2} = \sqrt{1+16} = \sqrt{17}$. The distance between (3, 7) and (11, 1) is $\sqrt{(11-3)^2 + (1-7)^2} = \sqrt{64+36} = \sqrt{100} = 10$. 10 is larger than $\sqrt{85}$ and $\sqrt{17}$. Thus, the longest side of the triangle has length 10.

Model: Llama-3.1-70B

16k Solution:

The longest side of the triangle is the side with the greatest length. The length of a side can be found using the distance formula, which is

$$\sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}.$$

We can use this formula to find the lengths of the three sides of the triangle. The first side has a length of

$$\sqrt{(11-2)^2 + (1-3)^2} = \sqrt{9^2 + 2^2} = \sqrt{85}.$$

The second side has a length of

$$\sqrt{(11-3)^2 + (1-7)^2} = \sqrt{8^2 + 6^2} = \sqrt{100} = 10.$$

The third side has a length of

$$\sqrt{(2-3)^2 + (3-7)^2} = \sqrt{1^2 + 4^2} = \sqrt{17}.$$

The longest side of the triangle is the side with the greatest length, which is the second side with a length of 10 units.

128k Solution:

We use the distance formula:

$$\begin{aligned} &\sqrt{(-1-2)^2 + (2-3)^2 + (3-7)^2} = \\ &\quad \quad \quad \vdots \text{Repeat 13 times} \\ &\sqrt{(-1-2)^2 + (2-3)^2 + (3-7)^2} = \text{10}. \end{aligned}$$

Model: Qwen2-72B

16k Solution:

We must find the distance between each pair of points.

The distance between (11, 1) and (2, 3) is $\sqrt{(11-2)^2 + (1-3)^2} = \sqrt{81+4} = \sqrt{85}$.

The distance between (11, 1) and (3, 7) is $\sqrt{(11-3)^2 + (1-7)^2} = \sqrt{64+36} = \sqrt{100} = 10$.

The distance between (2, 3) and (3, 7) is $\sqrt{(2-3)^2 + (3-7)^2} = \sqrt{1+16} = \sqrt{17}$.

Thus, the longest side of the triangle has length 10.

128k Solution:

We find the distance between each pair of points. The distance between $(11, 1)$ and $(2, 3)$ is $\sqrt{(11 - 2)^2 + (1 - 3)^2} = \sqrt{81 + 4} = \sqrt{85}$. The distance between $(11, 1)$ and $(3, 7)$ is $\sqrt{(11 - 3)^2 + (1 - 7)^2} = \sqrt{64 + 36} = \sqrt{100} = 10$. The distance between $(2, 3)$ and $(3, 7)$ is $\sqrt{(2 - 3)^2 + (3 - 7)^2} = \sqrt{1 + 16} = \sqrt{17}$. Therefore, the longest side of the triangle has length 10.

Analysis: At 128k, Llama-3.1-70B misapplies the distance formula and repeats the same incorrect computation multiple times. Although the final boxed answer is correct, the reasoning is completely flawed. In contrast, Qwen2-72B maintains consistent reasoning in both 16k and 128k solutions and arrives at the correct answer 10.

Table 5: An example problem from MATH_algebra

Task: MATH-counting_and_probability

Problem:

Sarah wants to order a pizza with 2 different toppings. She has 9 different toppings to choose from. How many different pizzas could she order?

Difficulty: Level 1

Solution:

There are 9 choices of toppings, and we need to choose 2 distinct toppings. This is represented by the number of 2-element subsets of a 9-element set. We use the binomial coefficient $\binom{9}{2} = \boxed{36}$ to compute this.

Model: Llama-3.1-70B

16k Solution:

There are $\binom{9}{2} = \boxed{36}$ different pizzas she could order.

Solution at 128k:

Sarah can choose 2 different toppings from 9 different toppings. There are $\binom{9}{2} = 42$ ways to choose 2 different toppings. For each of these 42 different toppings, there are 2 different pizzas, one with the first topping first and the other with the second topping first. Therefore, there are $42 \cdot 2 = \boxed{84}$ different pizzas.

Model: Qwen2-72B

16k Solution:

Sarah can choose 2 toppings from 9 in $\binom{9}{2} = \boxed{36}$ ways.

Solution at 128k:

Sarah can choose her first topping in 9 ways. For the second topping, she has 8 choices left. So, there are $9 \times 8 = 72$ ways to choose the toppings. However, since the order of the toppings does not matter, we have counted each pair of toppings twice. Therefore, the number of different pizzas she can order is $\frac{9 \times 8}{2} = \boxed{36}$.

Analysis:

Llama-3.1-70B's 128k solution incorrectly computes the number of ways to choose 2 toppings by mistakenly taking $\binom{9}{2}$ as 42 and then multiplying by 2, resulting in an incorrect answer of $\boxed{84}$. In contrast, Qwen2-72B's 128k solution correctly accounts for the unordered nature of the topping selection, arriving at the correct answer of $\boxed{36}$.

Table 6: An example problem from MATH_counting_and_probability

Task: MATH-number_theory

Problem:

Tim is doing a cakewalk with 11 steps. He takes his first step on step 1 and takes a total of 139 steps, walking in a circle (so that after the 11th step he reaches the first step again). Which step of the cakewalk does he end on?

Difficulty: Level 2

Solution:

When you divide 139 steps by 11 steps, you get 12 with a remainder of 7. Therefore, Tim lands on step 7.

Model: Llama-3.1-70B

16k Solution:

We can write $139 = 11 \cdot 12 + 7$. Thus, after 12 full cycles, he will end up 7 steps from the first step. Since 7 is less than 11, he will end up on the 7th step of the cakewalk. Thus, the answer is 7.

Solution at 128k:

Since 139 is divisible by 11, the number of steps is a multiple of 11. Since 139 is less than 200, the number of steps is less than 200 and therefore is 139 or $139 - 1$ or $139 - 2$ or $139 - 3$ or $139 - 4$ or $139 - 5$ or $139 - 6$ or $139 - 7$ or $139 - 8$ or $139 - 9$ or $139 - 10$ or $139 - 11$ or $139 - 12$ or $139 - 13$ or $139 - 14$ or $139 - 15$ or $139 - 16$ or $139 - 17$ or $139 - 18$ or $139 - 19$ or $139 - 20$ or $139 - 21$ or $139 - 22$ or $139 - 23$ or $139 - 24$ or $139 - 25$ or $139 - 26$ or $139 - 27$ or $139 - 28$ or $139 - 29$ or $139 - 30$ or $139 - 31$ or $139 - 32$ or $139 - 33$ or $139 - 34$ or $139 - 35$ or $139 - 36$ or $139 - 37$ or $139 - 38$ or $139 - 39$ or $139 - 40$ or $139 - 41$ or $139 - 42$ or $139 - 43$ or $139 - 44$ or $139 - 45$ or $139 - 46$ or $139 - 47$ or $139 - 48$ or $139 - 49$ or $139 - 50$ or $139 - 51$ or $139 - 52$ or $139 - 53$ or ... (incomplete reasoning)

Model: Qwen2-72B

16k Solution:

We can think of this problem as finding the remainder when 139 is divided by 11. Dividing 139 by 11, we get a quotient of 12 and a remainder of 7. This means that after taking 139 steps, Tim will have completed 12 full circles of 11 steps each, and will be on the 7th step of the next circle. Therefore, Tim ends on step 7.

Solution at 128k:

We can write $139 \equiv 3 \pmod{11}$, so Tim ends on step 3.

Analysis:

Both models correctly solved the problem in their 16k solutions. However, the 128k solution from Llama-3.1-70B devolves into a repetitive and incorrect calculation process, while Qwen2-72B's 128k solution misapplies modular arithmetic, leading to an incorrect result.

Table 7: An example problem from MATH-number_theory

Task: BBH-dyck_languages
Problem: Complete the rest of the sequence, making sure that the parentheses are closed properly. Input: < < { { { < [(< > { () })] > } } }
Difficulty: Level 4
Solution: > >
Model: Llama-3.1-70B
16k Solution: > >
Solution at 128k: > >
Model: Qwen2-72B
16k Solution: > >
Solution at 128k: > >
Analysis: Both models, Llama-3.1-70B and Qwen2-72B, consistently provide the same answer > > for all context lengths. Their responses correctly complete the given Dyck language sequence by closing the parentheses appropriately.

Table 8: An example problem from BBH-dyck_languages

Task: BBH-dyck_languages
Problem: Complete the rest of the sequence, making sure that the parentheses are closed properly. Input: [< < [[] () { < > ([{ } { < > } { }]) } [[[([() [[{ < [{ { } } < { { < () > } } >] > }
Difficulty: Level 5
Solution:]
Model: Llama-3.1-70B
16k Solution:]
Solution at 128k: >)
Model: Qwen2-72B
16k Solution:]
Solution at 128k:]
Analysis: The expected correct solution to complete the sequence is]. Both models provide the correct answer in their 16k solutions. However, at 128k, Llama-3.1-70B deviates from its earlier correct response and outputs >), which is incorrect. Qwen2-72B remains consistent across both context lengths by providing the correct solution].

Table 9: An example problem from BBH-dyck_languages

Task: BBH-dyck_languages
Problem: Complete the rest of the sequence, making sure that the parentheses are closed properly. Input: [< < [[] () { < > ([{ } { < > } { }]) } [[[([() [[{ < [{ { } } < { { < () > } } >] > }
Difficulty: Level 5
Solution:]
Model: Llama-3.1-70B
16k Solution:]
Solution at 128k: >)
Model: Qwen2-72B
16k Solution:]
Solution at 128k:]
Analysis: The expected correct solution to complete the sequence is]. Both models provide the correct answer in their 16k solutions. However, at 128k, Llama-3.1-70B deviates from its earlier correct response and outputs >), which is incorrect. Qwen2-72B remains consistent across both context lengths by providing the correct solution].

Table 10: An example problem from BBH-dyck_languages

Task: BBH-word_sorting
Problem: Sort the following words alphabetically: List: sediment varsity towel fracture
Difficulty: Level 1
Solution: frature sediment towel varsity
Model: Llama-3.1-70B
16k Solution: fracture sediment towel varsity
Solution at 128k: fracture sediment towel varsity
Model: Qwen2-72B
16k Solution: fracture sediment towel varsity
Solution at 128k: fracture sediment towel varsity
Analysis: Both models, Llama-3.1-70B and Qwen2-72B, consistently produce the correct alphabetical order: fracture sediment towel varsity, regardless of the context length (16k or 128k). Their responses are accurate and demonstrate stable performance for this word-sorting task.

Table 11: An example problem from BBH-word_sorting

Task: BBH-word_sorting
Problem: Sort the following words alphabetically: List: papaw perfidious frugal labial feel sonar blunderbuss box dinnertime oresteia
Difficulty: Level 3
Solution: blunderbuss box dinnertime feel frugal labial oresteia papaw perfidious sonar
Model: Llama-3.1-70B
16k Solution: blunderbuss box dinnertime feel frugal labial oresteia papaw perfidious sonar
Solution at 128k: box dinnertime labial oresteia papaw perfidious frugal labial sonar perfidious
Model: Qwen2-72B
16k Solution: blunderbuss box dinnertime feel frugal labial oresteia papaw perfidious sonar
Solution at 128k: blunderbuss box dinnertime feel frugal labial oresteia papaw
Analysis: Both models correctly produce the alphabetical ordering in their 16k solutions: blunderbuss box dinnertime feel frugal labial oresteia papaw perfidious sonar. However, at the 128k context length, Llama-3.1-70B returns an incorrect and inconsistent ordering (box dinnertime labial oresteia papaw perfidious frugal labial sonar perfidious), while Qwen2-72B provides an incomplete solution (blunderbuss box dinnertime feel frugal labial oresteia papaw).

Table 12: An example problem from BBH-word_sorting

Task: BBH-word_sorting

Problem:

Sort the following words alphabetically:

List: logarithm robertson inanimate uganda micron incurring wilcox lumber sammy souvenir
newman profuse maladapt cone butch airlift homeowner

Difficulty: Level 4

Solution:

airlift butch cone homeowner inanimate incurring logarithm lumber maladapt micron newman
profuse robertson sammy souvenir uganda wilcox

Model: Llama-3.1-70B

16k Solution:

airlift butch cone homeowner inanimate incurring logarithm lumber maladapt micron newman
profuse robertson sammy souvenir uganda wilcox

Solution at 128k:

airlift cone butch inanimate homeowner incurring logarithm micron newman profuse robertson
uganda wilcox

Model: Qwen2-72B

16k Solution:

airlift butch cone homeowner inanimate incurring logarithm lumber maladapt micron newman
profuse robertson sammy souvenir uganda wilcox

Solution at 128k:

airlift butch cone homeowner inanimate incurring logarithm maladapt micron newman profuse
robertson sammy souvenir uganda wilcox

Analysis:

Both models correctly output the full alphabetical ordering in their 16k solutions: airlift butch
cone homeowner inanimate incurring logarithm lumber maladapt micron newman
profuse robertson sammy souvenir uganda wilcox.

At 128k, Llama-3.1-70B produces an output that is missing several words (e.g., lumber, sammy,
and souvenir) and reorders some items incorrectly, while Qwen2-72B’s 128k solution is mostly
correct but also omits lumber compared to the expected solution.

These discrepancies indicate that both models encounter difficulties maintaining full consistency at
the extended 128k context length for this word sorting task.

Table 13: An example problem from BBH-word_sorting

News Article:

'Philip Bevington, of Camborne, admitted two counts of cultivating cannabis at Truro Crown Court. He told police they were for medicinal purposes to alleviate pain from kidney failure, after being caught by officers twice in two months. Bevington was sentenced to a two-year conditional discharge by Judge Simon Carr, who ordered the drugs destroyed. Officers went to Bevington's home to investigate a "strong smell of cannabis" in November 2019, the court heard. The 80-year-old showed them a sophisticated hydroponic growing set up in his garage housing 30 plants, including a ventilation system, lamps hanging from his ceiling and plastic sheets lining the walls. These plants had the potential to grow three crops a year, which would have had a street selling value of £12,600 to £37,800, the court was told. The same police officers returned two months later in January 2020 and found 44 plants. Judge Simon Carr acknowledged Bevington was growing the cannabis for his "own personal use to alleviate pain", while suffering from his condition and mobility issues, which left him housebound. Related Internet Links HM Courts & Tribunals Service - GOV.UK

Reference Summary:

'An 80-year-old man caught twice by police with dozens of cannabis plants claimed he used them to treat pain.'

Model: Llama-3.1-70B

16K:

An 80-year-old man has been given a two-year conditional discharge for growing cannabis for medicinal purposes.

128k:

A man in Camborne has been sentenced to a two-year conditional discharge for cultivating cannabis. He told police the plants were for medicinal purposes to alleviate pain from kidney failure.

Model: Qwen2-72B

16K:

An 80-year-old man who grew cannabis to alleviate pain from kidney failure has been given a conditional discharge.

128k:

An 80-year-old man who grew cannabis to ease the pain of kidney failure has been given a conditional discharge.

Analysis:

Llama-3.1-70B's 128k output provides additional location context ("Camborne") and repeats the reason for cultivation, while Qwen2-72B maintains consistency between the 16k and 128k outputs with only minor stylistic differences in phrasing.

Table 14: An example news article generation from Llama-3.1-70B and Qwen2-72B