

Understanding LLM Performance Degradation in Multi-Instance Processing: The Roles of Instance Count and Context Length

Jingxuan Chen Mohammad Taher Pilehvar Jose Camacho-Collados

School of Computer Science and Informatics, Cardiff University
 {ChenJ192,PilehvarMT,CamachoColladosJ}@cardiff.ac.uk

Abstract

Users often rely on Large Language Models (LLMs) for processing multiple documents or performing analysis over a number of instances. For example, analysing the overall sentiment of a number of movie reviews requires an LLM to process the sentiment of each review individually in order to provide a final aggregated answer. While LLM performance on such individual tasks is generally high, there has been little research on how LLMs perform when dealing with multi-instance inputs. In this paper, we perform a comprehensive evaluation of the multi-instance processing (MIP) ability of LLMs for tasks in which they excel individually. The results show that all LLMs follow a pattern of slight performance degradation for small numbers of instances ($\approx 20-100$), followed by a performance collapse on larger instance counts. Crucially, our analysis shows that while context length is associated with this degradation, the number of instances has a stronger effect on the final results. This finding suggests that when optimising LLM performance for MIP, attention should be paid to both context length and, in particular, instance count.¹

1 Introduction

LLMs have demonstrated remarkable capabilities across a wide range of natural language processing tasks and beyond (Wang et al., 2024). However, these capabilities have been predominantly evaluated in settings where a single instance is provided to the model at a time, which we refer to as single-instance processing (SIP). In contrast, many real-world applications, such as data analytics, document analysis, and large-scale information processing, require multi-instance processing (MIP), where the model generates individual predictions for multiple instances and subsequently aggregates them into a single, cohesive final prediction. The

¹Data and code are available at <https://github.com/jingxuanchen916/multi-instance-processing>.

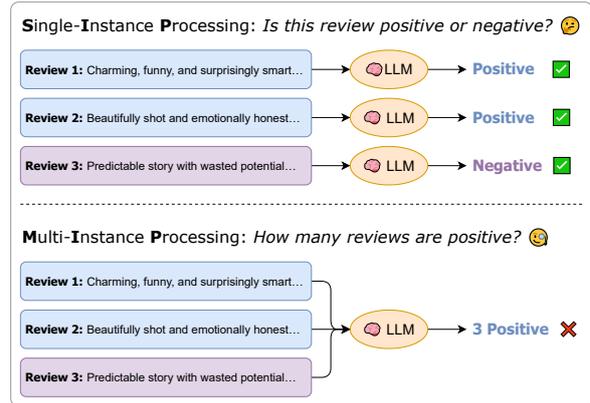


Figure 1: A toy example of SIP and MIP settings for sentiment analysis, where an LLM succeeds under SIP but fails under MIP given the same instances.

ability of LLMs to process multiple instances has been studied in recent literature in different contexts such as data analytics (Chen et al., 2025; Rahman et al., 2025; Sun et al., 2025) and long-context settings (Bertsch et al., 2025; Liu et al., 2025; Qiu et al., 2024; Shaham et al., 2023; Wolfson et al., 2026), which we unify and formalise under the MIP umbrella in this work.

Compared to SIP, MIP poses additional challenges due to its long-context nature and the need to perform repeated reasoning and aggregation over multiple instances, making it substantially more demanding for current LLMs (Bertsch et al., 2025). Ensuring reliable performance in such settings therefore requires a careful understanding of model failure modes, in order to inform the development of effective mitigation strategies for MIP.

As a motivating example, consider a non-expert user who inputs multiple movie reviews and wishes to determine how many of them are positive. Figure 1 illustrates a toy comparison between SIP and MIP: while an LLM can correctly classify the sentiment of each review in isolation, it may fail when required to process and aggregate all instances

within a single prompt. Given that the model is capable of accurately handling individual instances, it is crucial to understand the nature of the errors that arise when multiple instances must be processed jointly. Although alternative solutions, such as agentic designs that process instances separately or require users to manually batch instances or write code, are possible, these approaches are often impractical for non-expert users in real-world settings.

Existing work has extensively examined the challenges posed by long-context inputs, showing that model performance often degrades as context length increases, even when inputs remain within the model’s nominal context window (An et al., 2025; Moon and Lim, 2025). However, in most long-context benchmarks, increasing input length is accompanied by a simultaneous increase in task complexity (Hsieh et al., 2024). This coupling makes it difficult to disentangle whether observed performance degradation arises from longer inputs per se or from increased semantic and reasoning demands. In parallel, multi-instance or batch processing settings have been explored in prior work, but primarily from an efficiency or cost-reduction perspective, and typically with relatively small batch sizes (Cheng et al., 2023; Lin et al., 2024a). As a result, the effect of scaling the number of instances itself on model performance remains underexplored.

Given these gaps, we propose two research questions (RQs) to better understand LLM behaviour in MIP settings:

1. How does LLM performance change in MIP settings, and what failure behaviours emerge as the number of instances increases?
2. What are the primary drivers of performance degradation in MIP? Do instance count and context length have similar effects?

To answer these RQs, we evaluate sixteen LLMs across a diverse set of eight tasks, including calculation, token-level and sentence-level classification, for which they are known to perform well when only individual instances are provided. We analyse a broad range of open-weight and closed-source LLMs, and systematically examine the limits and factors that influence their performance in MIP. Our results show that model performance typically degrades gradually for small instance counts before collapsing at larger scales, and that this degradation is more strongly associated with the number of instances than with context length alone, even when the latter is substantially increased.

2 Related Work

The practical demand for MIP is typically driven by the use of LLM-based agents for complex data analytics, such as data wrangling, exploratory analysis, and multi-file reasoning (Guo et al., 2024; Hong et al., 2025; Nam et al., 2025). In these workflows, data processing is a foundational step and often requires transforming heterogeneous inputs into intermediate representations or executable code (Lin et al., 2024b; Shankar et al., 2024). While current approaches typically adopt modular pipelines to improve reliability through task decomposition (Nam et al., 2025; Shankar et al., 2024), these systems are ultimately bounded by the LLM’s ability to handle high-density information. Understanding the limits of such analytical workflows therefore requires examining LLM performance along two critical dimensions: the capacity to maintain coherence over long contexts and the ability to process multiple inputs efficiently through batch processing.

Long Context. Recent benchmarks have examined LLM performance degradation in long-context settings, spanning retrieval-based evaluations (Hsieh et al., 2024; Levy et al., 2025; Yang et al., 2025b), long-context reasoning over dispersed evidence (Kuratov et al., 2024; Vodrahalli et al., 2024; Zhang et al., 2024), and broader application-oriented suites (Yen et al., 2025). Complementary work also studies specialised regimes, including long procedural generation (Ye et al., 2025), scalable mathematical reasoning (Zhou et al., 2025), narrative understanding (Hamilton et al., 2025), and long-term conversational memory (Wu et al., 2025). Closest to our work, Bertsch et al. (2025) evaluates long-context information aggregation via many in-context instances, but varies task difficulty primarily through context length rather than isolating the effect of increasing the number of instances.

Batch Processing. Recent works study whether LLMs can answer multiple questions within a single prompt, typically motivated by reducing inference cost and balancing capacity limits (Cheng et al., 2023; Ji et al., 2025; Lin et al., 2024a). These studies consistently observe that only a small number of instances can be processed reliably before accuracy degrades. However, they do not investigate scaling behaviour beyond this regime, which arises in many practical settings where non-expert users may directly pass all instances to an LLM.

Moreover, their evaluations primarily focus on independent question answering (Wang et al., 2025) or classification tasks (Gozzi and Di Maio, 2024), rather than controlled settings in which task semantics are fixed and the instance count itself is the primary source of difficulty.

3 Multi-Instance Processing

We study MIP, where an LLM is required to reason over multiple input instances within a single prompt. In contrast to retrieval-augmented generation (RAG), which typically uses only a subset of retrieved inputs to produce the final answer, MIP requires processing all provided instances. As shown in Figure 1, LLMs must iterate over all instances individually in MIP to produce intermediate results, which are then aggregated into a final answer.

An instance is defined as a single data entry x (e.g., a movie review, a user post, a sentence, or a number). Let \mathcal{X} denote the set of all available instances. The model M is provided with a subset $\mathcal{X}' \subseteq \mathcal{X}$, where the number of instances $n = |\mathcal{X}'|$ may vary across inputs. SIP is treated as a special case of MIP where $n = 1$.

3.1 Formulation

Given a task instruction prompt τ and an input instance set $\mathcal{X}' = \{x_1, \dots, x_n\} \subseteq \mathcal{X}$, the model M , parameterised by θ , generates an output

$$o \sim p_\theta(\cdot \mid \tau, \mathcal{X}'),$$

where p_θ denotes the model’s conditional output distribution.

A valid model output typically contains an aggregated prediction y^{agg} together with a natural language explanation r . Let y^{agg^*} denote the corresponding ground-truth aggregated label. An output is considered correct when the aggregated prediction matches the ground truth, i.e., $y^{\text{agg}} = y^{\text{agg}^*}$. Otherwise, it is considered a wrong output and belongs to the set \mathcal{W} . The model may also produce an invalid output, in which case it fails to generate a well-formed prediction. Such outputs belong to the set \mathcal{I} . We define both wrong answers (i.e., valid but incorrect outputs) and invalid outputs as failures:

$$\mathcal{E}_{\text{fail}} = \mathcal{W} \cup \mathcal{I}.$$

We write $o \in \mathcal{E}_{\text{fail}}$ to denote that a model output is a failure.

3.2 Filtering for Controlled Difficulty

An important aspect of our methodology is ensuring that the evaluated tasks are simple enough for LLMs to solve when individual instances are provided (i.e., in SIP settings). Therefore, to control instance-level task difficulty when evaluating MIP, we construct inputs based on SIP outcomes.

Let $\mathcal{X}_{\text{SIP}} \subseteq \mathcal{X}$ denote the subset of instances for which all comparison models can produce the correct prediction under the SIP setting. MIP inputs are then formed by uniformly sampling subsets $\mathcal{X}' \subseteq \mathcal{X}_{\text{SIP}}$ using fixed random seeds. To further ensure reliable evaluation, we retain only models whose average SIP task success rate exceeds 95% and whose per-task SIP success rate exceeds 90%. We also keep only tasks for which agreement among all retained models exceeds 85%, measured as the proportion of instances that all comparison models answer correctly prior to filtering.

Our filtering procedure ensures that failures observed in the MIP setting are not attributable to intrinsic instance difficulty or ambiguity, but instead reflect the model’s ability to reason over and aggregate multiple instances, which is the primary focus of this work.

3.3 Evaluation Metrics

We define an experiment as a specific evaluation configuration represented as a tuple $e = (M, \tau, \mathcal{X}')$, which produces a model output o_e . We define accuracy as a binary metric. Let y^{agg}_e denote the aggregated prediction extracted from the model output o_e . Then:

$$\text{Acc}(e) = \begin{cases} 1, & \text{if } y^{\text{agg}}_e = y^{\text{agg}^*}_e, \\ 0, & \text{otherwise.} \end{cases}$$

Let \mathcal{D} denote the set of all evaluated experiments. The success rate (SR) is defined as the average accuracy across experiments:

$$\text{SR} = \frac{1}{|\mathcal{D}|} \sum_{e \in \mathcal{D}} \text{Acc}(e).$$

The invalid rate (IR) measures the fraction of experiments in which the model produces an invalid output:

$$\text{IR} = \frac{|\{e \in \mathcal{D} \mid o_e \in \mathcal{I}\}|}{|\mathcal{D}|}.$$

4 Experimental Setting

In this section, we describe our general experimental setting.

Name	Task
<i>Arithmetic</i>	Solve arithmetic problems & Sum of answers
<i>Category</i>	Classify news category & Aggregate class counts
<i>Language</i>	Identify language & Aggregate class counts
<i>NER</i>	Count “person” entities & Aggregate total counts
<i>Parity</i>	Detect odd or even number & Aggregate counts
<i>Sentiment</i>	Detect sentiment polarity & Aggregate counts
<i>Word</i>	Count target word “women” & Aggregate total counts
<i>WSD</i>	Identify “apple” word sense & Aggregate counts

Table 1: Overview of the selected tasks and their aggregation logic in the MIP setting. More details can be found in Appendix A.

4.1 Individual Tasks

We consider eight heterogeneous tasks² for our analysis, as summarised in Table 1. Each task is chosen such that it can be solved individually in the SIP setting by standard LLMs. When multiple instances are provided in the MIP setting, the model is tasked with additionally aggregating outputs across all instances (e.g., counting how many movie reviews are classified as positive in sentiment analysis). Detailed task descriptions and examples are provided in Appendix A.

4.2 Models and Prompting

We use OpenRouter³ to evaluate sixteen LLMs, including nine open-weight models (*DeepSeek R1*, *DeepSeek V3*, *gpt-oss-120b*, *gpt-oss-20b*, *Llama 3.3*, *Llama 4 Maverick*, *MiniMax M2.5*, *Qwen3-Instruct* and *Qwen3-Thinking*) and seven closed-source models (*Claude Sonnet 4.6*, *Gemini 2.5 Flash*, *Gemini 3.1 Pro*⁴, *GPT-5*, *GPT-5 Nano*, *Grok 4* and *Grok 4 Fast*).⁵

For prompting, we use a temperature of 0 and a maximum output length of 20K tokens for consistency across models. To allow limited tolerance

²We removed three additional tasks whose SIP performance fell below our requirements.

³<https://openrouter.ai>

⁴We used the preview version, which was the only version available at the time of experimentation (March 2026).

⁵As with the task filtering, we removed two open-weight LLMs, *Llama 4 Scout* and *Mistral NeMo*, whose SIP performance did not meet our criteria.

to formatting errors, we permit up to three retries when a model produces an invalid output belonging to \mathcal{I} . The full set of prompting templates is provided in Appendix B.

4.3 Single-Instance Filtering

As described in Section 3.2, we ensure that each instance can be successfully solved in the SIP setting. To this end, we conduct SIP experiments on 2,500 instances for each task.⁶ We report each LLM’s SIP performance for each task in Appendix C.1. Moreover, Table 7 in Appendix C.2 reports the percentage of instances retained (i.e., agreement) for each task (from 89% to nearly 100%), and the corresponding maximum and minimum SIP success rates across models, all exceeding 93%. Finally, as described in Section 3.2, this filtering retains only instances for which all comparison models agreed on the correct answer, thereby excluding potentially ambiguous instances and annotation errors.

4.4 MIP Sampling

After single-instance filtering, for each task τ we construct MIP inputs by sampling instances from \mathcal{X}_{SIP} using five different random seeds ($s \in \{1, 2, 3, 4, 5\}$). We evaluate ten MIP sample sizes $n \in \mathcal{N} = \{2, 5, 10, 20, 50, 100, 200, 500, 1000, 2000\}$. For each (τ, n, s) , we prompt each model M with the corresponding instance set, retaining only instances for which all models are correct in the SIP setting.

5 RQ1: Performance and Failure Behaviours

Our main goal is to evaluate LLM performance and failure behaviours in MIP settings, particularly as the number of instances increases.

5.1 Performance Analysis

Gradual success rate degradation followed by collapse as the number of instances increases.

Figure 2 reports success rates aggregated across tasks for each model. We observe a consistent performance degradation as the number of instances increases. In particular, all models show noticeable drops above 200 instances and near-collapse beyond 1,000 instances, with success rates falling

⁶The original dataset of *Category* contains fewer than 2,500 instances, and each instance is substantially longer. We therefore use 250 instances instead of 2,500 for this task. Correspondingly, the maximum MIP sample size for *Category* is also ten times smaller than for the other tasks.

below 20% at 2,000 instances. Figure 3 shows success rates aggregated across models for each task and reveals a similar downward trend. With the exception of *Arithmetic*, all tasks achieve success rates above 60% when fewer than 50 instances. Performance then deteriorates steadily as the instance count grows. Complete results by model and task are reported in Appendix D.1.

LLM comparison. Table 2 reports the success rate and invalid rate for all models. Overall, closed-source LLMs do not consistently exhibit superior performance. While frontier proprietary models (*GPT-5*, *Gemini 3.1 Pro* and *Grok 4*) achieve the highest success rates, this advantage comes at substantially higher cost.⁷ Beyond them, *Qwen3-Thinking*, *gpt-oss-120b*, *DeepSeek R1*, *Grok 4 Fast* and *GPT-5 Nano* achieve the highest success rates (above 65%). Notably, only *Llama 4 Maverick* and *Grok 4 Fast* produce no invalid outputs, indicating greater robustness. While all models achieve success rates above 35% on average, Figure 2 indicates that most successful cases occur when fewer than 500 instances are processed.

Robustness to instance order. To examine whether instance order affects LLM performance, we conduct an additional robustness experiment. For all experiments whose instance sets were originally sampled with random seed $s = 1$, we randomly shuffle the instance order twice (using $s = 6$ and $s = 7$) and rerun the evaluation. Figure 4 reports the resulting success rates as the number of instances increases. The degradation patterns remain highly consistent across different orderings of the same instance sets, suggesting that instance order has little effect on overall performance.

5.2 Failure Behaviours

Beyond our default setting, which requires only an aggregated answer, we introduce an additional variant for more fine-grained analysis. In this variant, models are required to produce instance-level predictions $\{y_i\}_{i=1}^n$ before providing the aggregated answer. Even with such instance-level predictions, which provide explicit intermediate reasoning, the relative performance of models remains similar to the aggregated-only setting.

As described in Section 3.1, we consider two broad categories of failures: wrong answers and

⁷Given the high cost of frontier models and practical budget constraints, we restrict subsequent experiments to only lightweight closed-source and open-weight models.

Model	Size / Cost	Success (%)	Invalid (%)
DeepSeek R1	A37B	67.5±2.6	2.9±0.6
DeepSeek V3	A37B	39.0±3.6	2.9±0.6
gpt-oss-120b	117B	68.3±2.8	3.6±1.1
gpt-oss-20b	21B	60.8±2.5	4.9±0.6
Llama 3.3	70B	39.0±3.8	2.9±0.6
Llama 4 Maverick	17B	43.1±1.1	0.0 ±0.0
MiniMax M2.5	A10B	62.3±1.8	16.1±2.7
Qwen3-Instruct	A22B	37.9±3.6	1.3±0.0
Qwen3-Thinking	A22B	69.4±2.4	3.9±1.6
Claude Sonnet 4.6	\$ 4.68	60.3±2.7	0.3±0.6
Gemini 2.5 Flash	\$ 0.40	37.7±3.9	4.2±3.2
Gemini 3.1 Pro	\$ 6.28	80.3±1.4	2.6±0.9
GPT-5	\$ 2.64	81.8 ±2.6	1.8±0.7
GPT-5 Nano	\$ 0.13	66.5±3.8	7.5±0.6
Grok 4	\$ 5.54	70.6±1.7	1.3±0.0
Grok 4 Fast	\$ 0.26	67.0±2.8	0.0 ±0.0

Table 2: Model success rate and invalid rate (mean±std), averaged across all tasks and instance counts. Standard deviation is computed over five random seeds. The top nine models are open-weight LLMs (DeepSeek R1, DeepSeek V3, MiniMax M2.5, Qwen3-Instruct, and Qwen3-Thinking are all mixture-of-experts LLMs, with total parameter counts of 671B, 671B, 230B, 235B, and 235B, respectively), where we report model size. The bottom seven are closed-source LLMs with undisclosed sizes, for which we report the average cost per task across five runs.

invalid outputs. Wrong answers include errors at the individual-instance level, the aggregation level, or both. Invalid outputs include (1) parsing errors, where the model output cannot be reliably parsed into the expected structured format, and (2) overlong input errors, where the input exceeds the model’s allowable context length. To analyse failure behaviours in greater detail, we use the instance-level variant in the following experiments.

Different failure types emerge as the number of instances increases. Figure 5 presents a stacked bar plot showing the contribution of each failure type as the instance count increases. Blue bars correspond to wrong answers, while orange bars correspond to invalid outputs. Wrong answers can occur even with as few as two instances. When the instance count exceeds 200, parsing errors increase substantially and reach nearly 30% at 2,000 instances. Overlong input errors emerge mainly beyond 200 instances, primarily due to the *Language* task, as non-English inputs typically require more prompt tokens (Petrov et al., 2023).

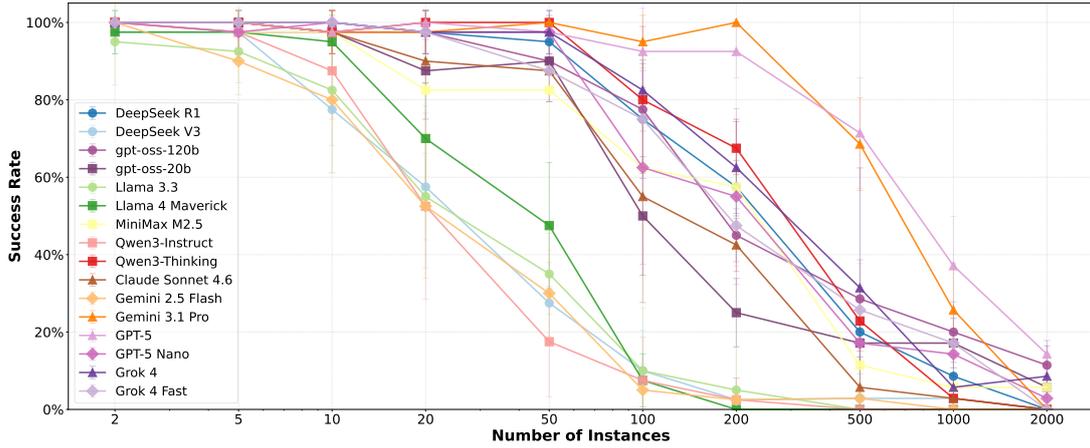


Figure 2: Model success rates (averaged across all tasks) as a function of the number of instances. Error bars indicate standard deviation across five random seeds. LLMs from the same company share the same colour family, while markers denote categories: ● (open-weight, $\geq 37\text{B}$ active parameters), ■ (open-weight, $\leq 22\text{B}$ active parameters), ▲ (frontier closed-source), and ◆ (lightweight closed-source).

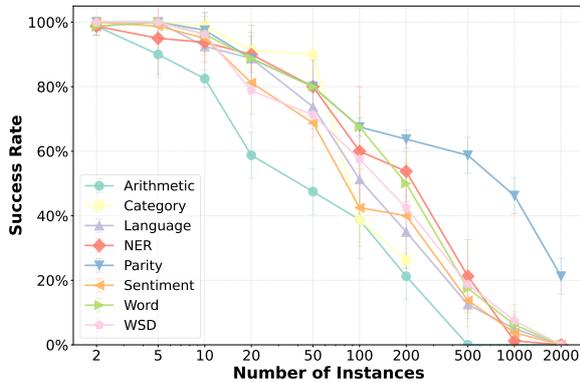


Figure 3: Task success rates (averaged across all LLMs) as a function of the number of instances. Error bars indicate standard deviation across five random seeds.

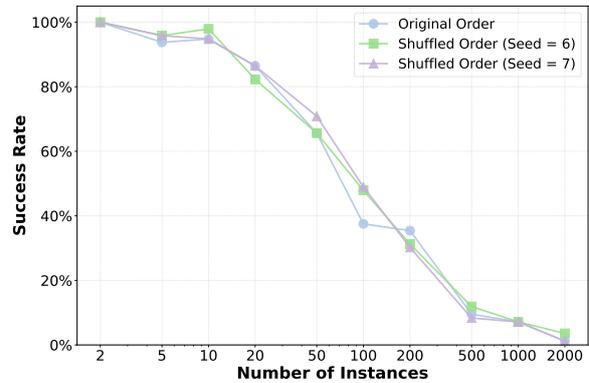


Figure 4: Success rates as a function of the number of instances for the original instance order and two shuffled variants constructed from the same instance sets.

Mistakes in individual instances and aggregation. To further characterise wrong answers, we categorise them into four error types: key mistakes (keys do not span from 1 to n), individual mistakes (at least one instance-level prediction is incorrect), aggregation mistakes (all instance-level predictions are correct but the aggregated answer is incorrect), and combined mistakes (both instance-level and aggregation errors occur). As shown in Figure 5, when the instance count exceeds 100, combined mistakes increase markedly, accounting for approximately 25% to 45% of failures. Moreover, when aggregation mistakes and combined mistakes are considered together, aggregation remains challenging for LLMs regardless of instance-level correctness or instance count. Complete results by model and task are reported in Appendix D.2.

Model differences in making mistakes. Table 3 compares models in terms of how individual-instance mistakes are distributed across experiments. Some models tend to concentrate many individual mistakes within a small number of failed experiments, while others exhibit more frequent but sparser errors. Focusing on individual mistakes, we observe that each failed experiment of *Grok 4 Fast* typically contains many incorrect instance-level predictions. In contrast, *Gemini 2.5 Flash* more often produces failed experiments with only a small number of individual mistakes, and *MiniMax M2.5* shows low values in both metrics.

Self-awareness of limitations. Ideally, an LLM should recognise its own capability boundaries. A desirable behaviour is for the model to explicitly acknowledge in its reasoning r when it cannot han-

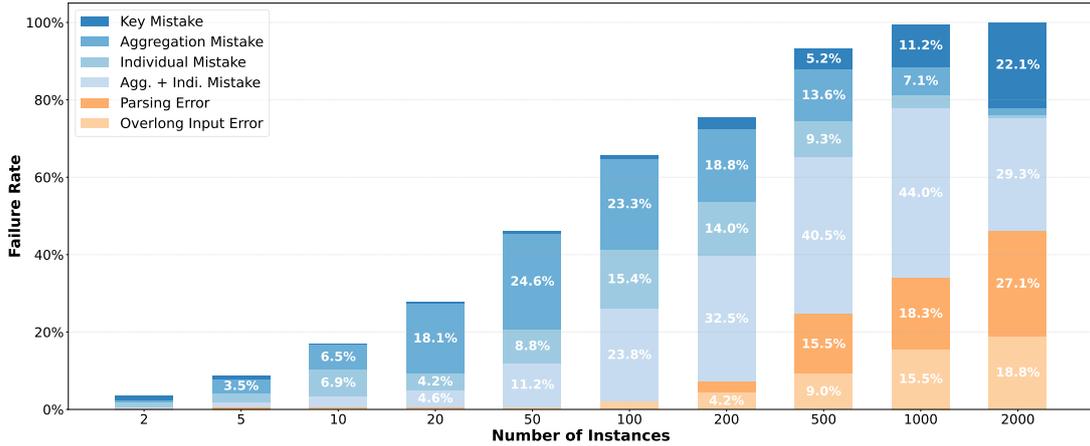


Figure 5: Breakdown of failure types. Key mistakes, aggregation mistakes, individual mistakes, and combined mistakes (Agg.+Indi.) are categorised as wrong answers (blue), while parsing errors and overlong input errors are categorised as invalid outputs (orange).

Model	Wrong Exp.	Wrong Indv. Per Exp.
DeepSeek R1	21.0±3.7	49.5±9.9
DeepSeek V3	29.6±2.7	34.8±13.1
gpt-oss-120b	16.2±2.9	37.5±17.4
gpt-oss-20b	15.8±2.4	11.8±5.0
Llama 3.3	35.6±2.2	53.8±26.7
Llama 4 Maverick	33.6±4.2	38.4±4.6
MiniMax M2.5	8.2±1.5	5.8±5.5
Qwen3-Instruct	29.4±2.1	20.8±3.1
Qwen3-Thinking	16.0±2.9	15.4±4.0
Gemini 2.5 Flash	13.4±1.3	5.5±3.2
GPT-5 Nano	11.2±3.1	10.3±15.3
Grok 4 Fast	22.4±1.7	117.3±8.0

Table 3: Percentage of experiments with at least one incorrect instance-level prediction (Wrong Exp.) and the average number of incorrect instance-level predictions per failed experiment (mean±std). Standard deviation is computed over five random seeds.

de a large number of instances (e.g., by suggesting batch-wise processing or explicitly stating that the instance count exceeds its capacity). Such behaviour can be reflected by the model producing predictions for only the first few instances and omitting the remaining ones, which corresponds to a key mistake (missing key). After inspecting the model outputs, we found that only 171 out of 4,620 experiments exhibit this omission, almost exclusively at instance counts of 500 or more, as expected. However, a manual analysis shows that only 27 out of these 171 experiments explicitly suggest batch-wise processing, and most models do not warn users about such limitations. *GPT-5 Nano* demon-

strates this behaviour most frequently, explicitly indicating difficulty in 19 out of 28 cases. In contrast, *DeepSeek V3*, *gpt-oss-120b*, and *gpt-oss-20b* do so in fewer than 20% of its omission cases, while no such behaviour is observed for the remaining models. Notably, although *Qwen3-Instruct* and *Gemini 2.5 Flash* do not acknowledge limitations, they also produce almost no cases of this omission.

5.3 Discussion

Overall, our findings answer RQ1 by showing that as the instance count increases, all LLMs experience degraded performance, characterised by lower success rates and more frequent failures. In particular, when the instance count reaches 1,000 or higher, no LLM achieves a success rate above 40%.

Among the evaluated LLMs, *GPT-5*, *Gemini 3.1 Pro*, *Grok 4*, *gpt-oss-120b* and *Qwen3-Thinking* exhibit the strongest overall performance. Focusing on the lightweight and open-weight models, although *Gemini 2.5 Flash* achieves relatively low success rates, its failures typically involve a relatively small number of incorrect instances. At the same time, invalid outputs become increasingly common as the instance count grows. Notably, *GPT-5 Nano* is the only model that consistently identifies when a task exceeds its capacity.

6 RQ2: Context Length vs Number of Instances

In the previous section, we analysed LLM behaviour in MIP settings as the number of instances increases. The results revealed consistent performance degradation across all models, albeit at dif-

ferent rates and to different extents. A natural question is whether this degradation is driven by increased context length, as has been observed in prior work across a range of tasks and settings (see Section 2). In this section, we analyse the effects of context length and instance count jointly, and examine their respective contributions.

6.1 Context Length Augmentation

To study the impact of context length, we design a setting in which the length of each individual instance is artificially increased without altering its original content.⁸ For each sampled instance x , we construct a perturbed instance $x' = x + \epsilon$, where ϵ denotes injected noise. Following Hsieh et al. (2024), we define ϵ as the string “- IRRELEVANT CONTEXT START -” followed by seven repetitions of the sentence “The grass is green. The sky is blue. The sun is yellow. Here we go. There and back again.”, plus “- IRRELEVANT CONTEXT END -”. This choice is motivated by prior work showing that even irrelevant context can degrade model performance (Shi et al., 2023; Yang et al., 2025a). After noise injection, the average length of each instance (measured in the number of prompt tokens) more than doubles, increasing from approximately 136 tokens in the default setting to 326 tokens in the artificially augmented setting.

Results. Figure 6 compares average performance across all tasks and models between the default and artificially augmented settings, as the number of instances increases. When the number of instances is held constant, the success rates of the two settings remain broadly similar, despite the average context length being more than twice as large in the augmented setting. This indicates that artificially increasing the length of individual instances does not substantially impact performance, suggesting that context length alone is unlikely to fully explain the performance degradation observed in MIP settings. Complete results by model and task are reported in Appendix D.2.2.

Robustness to noise position. To further examine whether the position of injected noise affects model performance, we vary where the noise is inserted for the same instance set sampled with

⁸We exclude *Parity* from this experiment because each instance contains only a single number. For consistency with context length constraints, we also cap the maximum sample size at $n = 1000$ in this setting ($n = 100$ for *Category*), as discussed below.

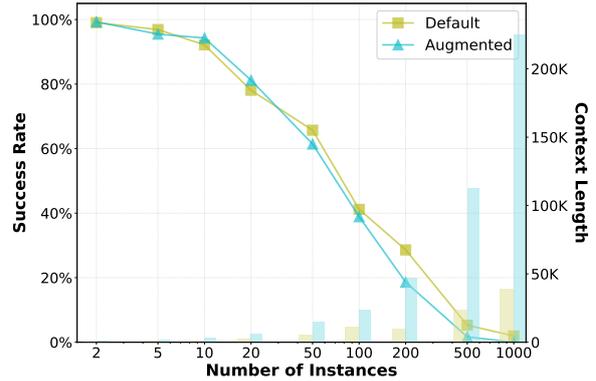


Figure 6: Success rate (lines) and total prompt token length (bars) in the artificial length setting as the number of instances increases. Error bars indicate standard deviation across five random seeds.

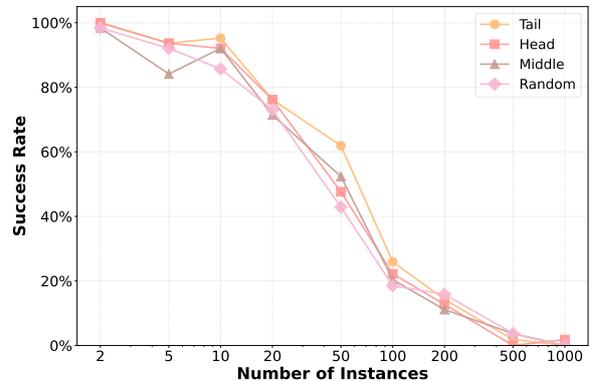


Figure 7: Success rates as a function of the number of instances for different injected-noise positions constructed from the same instance sets.

$s = 1$. In addition to the tail setting used above, we also insert the noise at the head, in the middle, and at random positions within each instance. Figure 7 shows that the performance remains broadly similar across these variants, indicating that the position of injected noise does not lead to substantial performance differences.

6.2 Correlation Analysis

Motivated by the above findings, we conduct a correlation analysis to further examine the relationships of instance count and context length to overall model performance. Recall that in our default setting, instance samples are constructed using five random seeds. While this random sampling introduces some variation in individual instance lengths, the average prompt length tends to be similar across samples, especially when the number of instances is high. To obtain a wider range of context lengths for correlation analysis, we augment the default

Number of Instances	Success rate and Context length	
	Correlation	P-Value
2	0.043	0.498
5	0.044	0.491
10	0.044	0.490
20	-0.042	0.510
50	-0.033	0.602
100	-0.014	0.823
200	0.056	0.420
500	0.117	0.102
1000	-0.053	0.475

Table 4: Correlation between success rate and prompt token length when the number of instances is fixed.

samples with two additional variants drawn from \mathcal{X}_{SIP} : a *long* set, consisting of the longest instances in each dataset, and a *short* set, consisting of the shortest instances. This design enables a clearer comparison across different context lengths by introducing greater variation in total input length.

Correlation with instance count and total context length. As an initial analysis, we compute the Spearman correlation between success rate and each factor independently: the number of instances and the total context length. Both factors exhibit strong negative correlations, indicating that performance decreases as either quantity increases. However, the number of instances shows a notably stronger relationship, with Spearman correlations of -0.61 , compared to that of -0.37 for the context length. In both cases, the corresponding p -values are below 0.001 , indicating that the correlations are highly unlikely to arise by chance.

Correlation conditioned on the number of instances. The two factors are inherently related, as total context length grows with the number of instances. To better understand their effects, we additionally compute correlations while holding the number of instances fixed and examining variation only in context length. As shown in Table 4, the resulting correlations between context length and success rate are substantially weaker, with values ranging between -0.15 and 0.15 , and p -values consistently above 0.1 . These results indicate that context length alone has limited explanatory power in this setting and suggest that the performance degradation observed in Section 6.1 is more strongly associated with the number of instances than total input length.

6.3 Discussion

Based on the above results, we find that the number of instances plays a stronger role than context length in determining model success rates. When the number of instances is held fixed, the effect of context length appears to be comparatively limited according to our correlation analysis. LLMs have been shown to struggle when required to perform many repeated operations (Son et al., 2024; Fu et al., 2024), which is precisely what MIP entails as the instance count increases. This behaviour contrasts with RAG settings, where models primarily need to identify relevant contexts rather than process all inputs exhaustively. In MIP settings, LLMs must process each instance individually and aggregate the resulting outputs. While prior work has shown that LLMs can handle increasingly long contexts (Liu et al., 2025), our findings suggest that improving reliability in MIP settings may also require training strategies that explicitly target multi-instance reasoning and aggregation.

7 Conclusion

In this paper, we presented a comprehensive evaluation of LLMs in MIP settings, that is, tasks that require aggregation of information from multiple instances to produce a final answer. The results show that LLMs are generally able to solve tasks involving a small number of instances, but begin to make mistakes as the number of instances increases. While the errors are initially small, this has important implications for user trust, since models are able to consistently solve the task when only a single instance is given. Moreover, as the number of instances further increases, models’ performance starts collapsing, in most cases without any warning to the user that this may happen.

Crucially, our experiments on context length highlight the importance of reasoning at the instance-count level, rather than focusing solely on context length as is commonly done. This has implications for how context should be processed in batches, for example when developing data science agents. Instead of relying only on context length for batching decisions in MIP settings, we should also consider the number of instances as a relevant factor. More generally, our results suggest that models may benefit from training strategies that better support multi-instance reasoning and aggregation, especially in settings where accuracy and user trustworthiness are paramount.

Limitations

This work focuses on diagnosing failure modes of LLMs in multi-instance processing (MIP) settings, rather than proposing or validating concrete solutions. While our controlled evaluation reveals performance degradation as the number of instances increases, we do not evaluate mitigation strategies such as task decomposition, external tool use, verification, or agentic designs. As a result, this paper should be interpreted primarily as an empirical characterisation of model behaviour rather than as a prescription for improving MIP reliability.

Our experiments emphasise exact aggregation tasks (e.g., counting, summation, exact class frequencies). Although these tasks are common in analytics-style applications, they may overemphasise brittleness in settings where approximate or semantic aggregation would suffice. The extent to which our findings generalise to softer aggregation objectives (e.g., majority voting, summarisation, or trend identification) therefore remains an open question. In addition, our experiments rely on a fixed prompt template across tasks, which may not fully capture the variability of prompt formulations encountered in real-world applications.

Despite our efforts to examine the roles of instance count and context length, these two factors are not entirely independent in practice. While correlation and controlled noise-injection analyses suggest that instance count has a stronger relationship with degradation, more fine-grained causal analyses (e.g., controlled computational complexity or attention-level diagnostics) are left for future work. Moreover, we did not stress-test LLMs with extremely large context inputs, which might have given different results and perhaps a stronger context length effect.

Finally, our study does not include model-internal interpretability analyses such as attention patterns or hidden-state dynamics, nor does it evaluate training-time interventions. In addition, our experiments are English-centric and limited to a set of selected LLMs, due to time and budget constraints. Future work could explore cross-lingual MIP behaviour and architectural or training modifications explicitly designed for multi-instance reasoning.

Ethical Considerations

This work presents an empirical evaluation of LLMs in multi-instance processing settings and does not involve the collection of new data, in-

teraction with human subjects, or the deployment of models in real-world decision-making systems. All datasets used in our experiments are publicly available and open-source, and were accessed and processed in accordance with their original licences and intended research use. Also, we do not claim that any particular model or provider is inherently unsafe or unsuitable. Rather, our results reflect general limitations of current LLMs under specific experimental conditions. We hope that this work contributes to more responsible deployment of LLM-based systems by encouraging practitioners to consider instance-level reliability, aggregation strategies, and potential failure modes when designing real-world applications.

Acknowledgments

Jose Camacho-Collados was supported by a UKRI Future Leaders Fellowship.

References

- Chenxin An, Jun Zhang, Ming Zhong, Lei Li, Shansan Gong, Yao Luo, Jingjing Xu, and Lingpeng Kong. 2025. Why does the effective context length of llms fall short? In *The Thirteenth International Conference on Learning Representations*.
- Francesco Barbieri, Jose Camacho-Collados, Luis Espinosa Anke, and Leonardo Neves. 2020. [TweetEval: Unified benchmark and comparative evaluation for tweet classification](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1644–1650, Online. Association for Computational Linguistics.
- Amanda Bertsch, Adithya Pratapa, Teruko Mitamura, Graham Neubig, and Matthew R Gormley. 2025. Oolong: Evaluating long context reasoning and aggregation capabilities. *arXiv preprint arXiv:2511.02817*.
- Ke Chen, Peiran Wang, Yaoning Yu, Xianyang Zhan, and Haohan Wang. 2025. Large language model-based data science agent: A survey. *arXiv preprint arXiv:2508.02744*.
- Zhoujun Cheng, Jungo Kasai, and Tao Yu. 2023. [Batch prompting: Efficient inference with large language model APIs](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: Industry Track*, pages 792–810, Singapore. Association for Computational Linguistics.
- Alexis Conneau and Douwe Kiela. 2018. [SentEval: An evaluation toolkit for universal sentence representations](#). In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, Miyazaki, Japan. European Language Resources Association (ELRA).

- Tairan Fu, Raquel Ferrando, Javier Conde, Carlos Ariaga, and Pedro Reviriego. 2024. Why do large language models (llms) struggle to count letters? *arXiv preprint arXiv:2412.18626*.
- Manuel Gozzi and Federico Di Maio. 2024. Comparative analysis of prompt strategies for large language models: Single-task vs. multitask prompts. *Electronics*, 13(23):4712.
- Derek Greene and Pádraig Cunningham. 2006. Practical solutions to the problem of diagonal dominance in kernel document clustering. In *Proceedings of the 23rd international conference on Machine learning*, pages 377–384.
- Siyuan Guo, Cheng Deng, Ying Wen, Hechang Chen, Yi Chang, and Jun Wang. 2024. Ds-agent: automated data science by empowering large language models with case-based reasoning. In *Proceedings of the 41st International Conference on Machine Learning*, pages 16813–16848.
- Sil Hamilton, Rebecca MM Hicke, Matthew Wilkens, and David Mimno. 2025. Too long, didn’t model: Decomposing llm long-context understanding with novels. *arXiv preprint arXiv:2505.14925*.
- Sirui Hong, Yizhang Lin, Bang Liu, Bangbang Liu, Binhao Wu, Ceyao Zhang, Danyang Li, Jiaqi Chen, Jiayi Zhang, Jinlin Wang, Li Zhang, Lingyao Zhang, Min Yang, Mingchen Zhuge, Taicheng Guo, Tuo Zhou, Wei Tao, Robert Tang, Xiangtao Lu, and 9 others. 2025. [Data interpreter: An LLM agent for data science](#). In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 19796–19821, Vienna, Austria. Association for Computational Linguistics.
- Cheng-Ping Hsieh, Simeng Sun, Samuel Kriman, Shantanu Acharya, Dima Rekish, Fei Jia, and Boris Ginsburg. 2024. Ruler: What’s the real context size of your long-context language models? In *First Conference on Language Modeling*.
- Zhaoxuan Ji, Xinlu Wang, Zhaojing Luo, Zhongle Xie, and Meihui Zhang. 2025. Optimized batch prompting for cost-effective llms. *Proceedings of the VLDB Endowment*, 18(7):2172–2184.
- Yury Kuratov, Aydar Bulatov, Petr Anokhin, Ivan Rodkin, Dmitry Sorokin, Artyom Sorokin, and Mikhail Burtsev. 2024. Babilong: Testing the limits of llms with long context reasoning-in-a-haystack. *Advances in Neural Information Processing Systems*, 37:106519–106554.
- Shahar Levy, Nir Mazon, Lihi Shalmon, Michael Hasid, and Gabriel Stanovsky. 2025. [More documents, same length: Isolating the challenge of multiple documents in RAG](#). In *Findings of the Association for Computational Linguistics: EMNLP 2025*, pages 19539–19547, Suzhou, China. Association for Computational Linguistics.
- Jianzhe Lin, Maurice Diesendruck, Liang Du, and Robin Abraham. 2024a. Batchprompt: Accomplish more with less. In *The Twelfth International Conference on Learning Representations*.
- Yiming Lin, Madelon Hulsebos, Ruiying Ma, Shreya Shankar, Sepanta Zeigham, Aditya G Parameswaran, and Eugene Wu. 2024b. Towards accurate and efficient document analytics with large language models. *arXiv preprint arXiv:2405.04674*.
- Jiaheng Liu, Dawei Zhu, Zhiqi Bai, Yancheng He, Huanxuan Liao, Haoran Que, Zekun Wang, Chenchen Zhang, Ge Zhang, Jiebin Zhang, and 1 others. 2025. A comprehensive survey on long context language modeling. *arXiv preprint arXiv:2503.17407*.
- Daniel Loureiro, Kiamehr Rezaee, Mohammad Taher Pilehvar, and Jose Camacho-Collados. 2021. [Analysis and evaluation of language models for word sense disambiguation](#). *Computational Linguistics*, 47(2):387–443.
- Hyeonseok Moon and Heuseok Lim. 2025. Needlechain: Measuring intact long-context reasoning capability of large language models. *arXiv preprint arXiv:2507.22411*.
- Jaehyun Nam, Jinsung Yoon, Jiefeng Chen, and Tomas Pfister. 2025. Ds-star: Data science agent via iterative planning and verification. *arXiv preprint arXiv:2509.21825*.
- Aleksandar Petrov, Emanuele La Malfa, Philip Torr, and Adel Bibi. 2023. Language model tokenizers introduce unfairness between languages. *Advances in neural information processing systems*, 36:36963–36990.
- Zexuan Qiu, Jingjing Li, Shijue Huang, Xiaoqi Jiao, Wanjun Zhong, and Irwin King. 2024. Clongeval: A chinese benchmark for evaluating long-context large language models. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 3985–4004.
- Afshin Rahimi, Yuan Li, and Trevor Cohn. 2019. [Massively multilingual transfer for NER](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 151–164, Florence, Italy. Association for Computational Linguistics.
- Mizanur Rahman, Amran Bhuiyan, Mohammed Saidul Islam, Md Tahmid Rahman Laskar, Ridwan Mahbub, Ahmed Masry, Shafiq Joty, and Enamul Hoque. 2025. Llm-based data science agents: A survey of capabilities, challenges, and future directions. *arXiv preprint arXiv:2510.04023*.
- David Saxton, Edward Grefenstette, Felix Hill, and Pushmeet Kohli. 2019. Analysing mathematical reasoning abilities of neural models. In *International Conference on Learning Representations*.

- Uri Shaham, Maor Ivgi, Avia Efrat, Jonathan Berant, and Omer Levy. 2023. Zeroscrolls: A zero-shot benchmark for long text understanding. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 7977–7989.
- Shreya Shankar, Tristan Chambers, Tarak Shah, Aditya G Parameswaran, and Eugene Wu. 2024. Docetl: Agentic query rewriting and evaluation for complex document processing. *arXiv preprint arXiv:2410.12189*.
- Freda Shi, Xinyun Chen, Kanishka Misra, Nathan Scales, David Dohan, Ed H Chi, Nathanael Schärli, and Denny Zhou. 2023. Large language models can be easily distracted by irrelevant context. In *International Conference on Machine Learning*, pages 31210–31227. PMLR.
- Natalia Silveira, Timothy Dozat, Marie-Catherine de Marneffe, Samuel Bowman, Miriam Connor, John Bauer, and Chris Manning. 2014. A gold standard dependency corpus for English. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC’14)*, pages 2897–2904, Reykjavik, Iceland. European Language Resources Association (ELRA).
- Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and Christopher Potts. 2013. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 1631–1642, Seattle, Washington, USA. Association for Computational Linguistics.
- Guijin Son, SangWon Baek, Sangdae Nam, Ilgyun Jeong, and Seungone Kim. 2024. Multi-task inference: Can large language models follow multiple instructions at once? In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5606–5627, Bangkok, Thailand. Association for Computational Linguistics.
- Maojun Sun, Ruijian Han, Binyan Jiang, Houduo Qi, Defeng Sun, Yancheng Yuan, and Jian Huang. 2025. A survey on large language model-based agents for statistics and data science. *The American Statistician*, pages 1–14.
- Martin Thoma. 2018. The wili benchmark dataset for written language identification. *arXiv preprint arXiv:1801.07779*.
- Kiran Vodrahalli, Santiago Ontanon, Nilesch Tripuraneni, Kelvin Xu, Sanil Jain, Rakesh Shivanna, Jeffrey Hui, Nishanth Dikkala, Mehran Kazemi, Bahare Fatemi, and 1 others. 2024. Michelangelo: Long context evaluations beyond haystacks via latent structure queries. *arXiv preprint arXiv:2409.12640*.
- Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, and 1 others. 2024. A survey on large language model based autonomous agents. *Frontiers of Computer Science*, 18(6):186345.
- Zhengxiang Wang, Jordan Kodner, and Owen Rambow. 2025. Exploring limitations of LLM capabilities with multi-problem evaluation. In *The Sixth Workshop on Insights from Negative Results in NLP*, pages 121–140, Albuquerque, New Mexico. Association for Computational Linguistics.
- Tomer Wolfson, Harsh Trivedi, Mor Geva, Yoav Goldberg, Dan Roth, Tushar Khot, Ashish Sabharwal, and Reut Tsarfaty. 2026. Monaco: More natural and complex questions for reasoning across dozens of documents. *Transactions of the Association for Computational Linguistics*, 14:23–46.
- Di Wu, Hongwei Wang, Wenhao Yu, Yuwei Zhang, Kai-Wei Chang, and Dong Yu. 2025. Longmemeval: Benchmarking chat assistants on long-term interactive memory. In *The Thirteenth International Conference on Learning Representations*.
- Minglai Yang, Ethan Huang, Liang Zhang, Mihai Surdeanu, William Yang Wang, and Liangming Pan. 2025a. How is LLM reasoning distracted by irrelevant context? an analysis using a controlled benchmark. In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 13329–13347, Suzhou, China. Association for Computational Linguistics.
- Yijun Yang, Zeyu Huang, Wenhao Zhu, Zihan Qiu, Fei Yuan, Jeff Z Pan, and Ivan Titov. 2025b. A controllable examination for long-context language models. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.
- Xi Ye, Fangcong Yin, Yinghui He, Joie Zhang, Howard Yen, Tianyu Gao, Greg Durrett, and Danqi Chen. 2025. Longproc: Benchmarking long-context language models on long procedural generation. *arXiv preprint arXiv:2501.05414*.
- Howard Yen, Tianyu Gao, Minmin Hou, Ke Ding, Daniel Fleischer, Peter Izsak, Moshe Wasserblat, and Danqi Chen. 2025. Helmet: How to evaluate long-context models effectively and thoroughly. In *The Thirteenth International Conference on Learning Representations*.
- Xinrong Zhang, Yingfa Chen, Shengding Hu, Zihang Xu, Junhao Chen, Moo Khai Hao, Xu Han, Zhen Leng Thai, Shuo Wang, Zhiyuan Liu, and 1 others. 2024. \$infty\$Bench: Extending long context evaluation beyond 100k tokens. *arXiv preprint arXiv:2402.13718*.
- Yang Zhou, Hongyi Liu, Zhuoming Chen, Yuandong Tian, and Beidi Chen. 2025. Gsm-infinite: How do your llms behave over infinitely increasing context length and reasoning complexity? *arXiv preprint arXiv:2502.05252*.

A Task and Data Source

Here we introduce the details for each task and their data source.

A.1 Arithmetic

The task is about **arithmetic** calculation, which uses questions from Mathematics Dataset (Saxton et al., 2019), including only easy addition or subtraction questions from its training split. The average word count for this task is 4.7.

Example Input:

1. What is the difference between -2 and 251860?
2. $-9,259,432 + 1$
3. What is 1,141.09 less than 1?

SIP Ground Truth Output:

1. 251,862
2. -9,259,431
3. -1,140.09

MIP Question: Solve all the provided arithmetic questions and calculate the sum of all answers.

MIP Ground Truth Output: -9,008,709.09

A.2 Category

The task is about news **category** classification, where a news article can belong to “tech” or “business” or “entertainment” or “politics” or “sport”. The dataset is BBC News (Greene and Cunningham, 2006), where we use its training split. The average word count for this task is 371.4.

Example Input:

1. german business confidence slides german business confidence fell in february knocking hopes of a speedy recovery in europe s largest economy...
2. bbc poll indicates economic gloom citizens in a majority of nations surveyed in a bbc world service poll believe the world economy is worsening...
3. lifestyle governs mobile choice faster better or funkier hardware alone is not going to help phone firms sell more handsets research suggests...

SIP Ground Truth Output:

1. business
2. business
3. tech

MIP Question: Count how many of the provided news articles belong to the “tech” category.

MIP Ground Truth Output: 1

A.3 Language

The task is about **language** identification, where a paragraph can belong to “English” or “Chinese” or “Persian” or “Spanish”. The dataset is WiLI-2018 (Thoma, 2018). The average word count for this task is 55.8.

Example Input:

1. Nordahl Road is a station served by North County Transit District’s SPRINTER light rail line...
2. En Navidad de 1974, poco después de que interpretó la canción en francés película Papillon (Toi qui Regarde la Mer)...
3. A talk by Takis Fotopoulos about the Internationalization of the Capitalist Market Economy and the project of Inclusive Democracy...

SIP Ground Truth Output:

1. English
2. Spanish
3. English

MIP Question: Count how many paragraphs are in English.

MIP Ground Truth Output: 2

A.4 NER

The task is about **named entity recognition**, which uses data from WikiANN (Rahimi et al., 2019). The average word count for this task is 16.0.

Example Input:

1. we love everything about the fence .
2. i want to hook up with that girl paige in the brown leather jacket .
3. in addition , there is a reduction of 22,101 mmbtu which is the difference between the scada values (best available) that anita showed on the february 29th storage sheet and the " official " february 29th values that gary wilson received from mips .

SIP Ground Truth Output:

1. 0
2. 1
3. 2

MIP Question: Count occurrences of the entity 'PERSON' in all sentences.

MIP Ground Truth Output: 3

A.5 Parity

The task is about **parity** classification (i.e., identify if a number is “odd” or “even”), where we use synthetic data generated by ourselves. The average word count for this task is 1 since only a single number is provided.

Example Input:

1. 18010
2. 10160
3. 89449

SIP Ground Truth Output:

1. even
2. even
3. odd

MIP Question: Count how many of the provided numbers are odd.

MIP Ground Truth Output: 1

A.6 Sentiment

The task is about **sentiment** analysis, where a movie review can belong to “positive” or “negative”. The dataset is Sentiment Treebank (Socher et al., 2013), where we only use the “most” positive and negative reviews to avoid ambiguity. The average word count for this task is 18.7.

Example Input:

1. High Crimes is a cinematic misdemeanor , a routine crime thriller remarkable only for its lack of logic and misuse of two fine actors , Morgan Freeman and Ashley Judd .
2. One of the worst movies of the year .
3. A mix of gritty realism , crisp storytelling and radiant compassion that effortlessly draws you in .

SIP Ground Truth Output:

1. negative
2. negative
3. positive

MIP Question: Count how many of the provided movie reviews are positive.

MIP Ground Truth Output: 1

A.7 Word

The task is about tweets **word** occurrence (i.e., count a target word's occurrences in given tweets). The dataset is TweetEval (Barbieri et al., 2020), where we use its stance detection subset. The average word count for this task is 17.3.

Example Input:

1. IF FEMINISTS WERE HONEST “I want a worldwide matriarchal dictatorship with all men enslaved to women” #GamerGate #SemST
2. What the fuck do women even do? I mean seriously they're just useless other than sex. #womensrights #Feminist #SemST
3. DEAR FEMINISTS Start asking for accountability from man-haters instead of shielding them for convenient concealment. #SemST

SIP Ground Truth Output:

1. 1
2. 2
3. 0

MIP Question: Count occurrences of the word “women” in all tweets.

MIP Ground Truth Output: 3

A.8 WSD

The task is about word sense disambiguation, where a target word “apple” is required to be distinguished as meaning either “company” or “fruit” based on its context. The dataset is CoarseWSD-20 (Loureiro et al., 2021), where we use its “apple” subset. The average word count for this task is 31.4.

Example Input:

1. both seasons are available for download from apple 's itunes store .
2. in klayman ii , the plaintiffs sued the same government defendants and in addition , facebook , yahoo! , google , microsoft , youtube , aol , paltalk , skype , sprint , at&t , apple again alleging the bulk metadata collection violates the first , fourth and fifth amendment and constitutes divulgence of communication records in violation of section 2702 of stored communications act .
3. description alongside dried pears the filling also contains raisin , walnut and other dried fruit such as apple or figs .

SIP Ground Truth Output:

1. company
2. company
3. fruit

MIP Question: Count how many paragraphs contain the word "apple" referring to the company (Apple Inc.), not the fruit.

MIP Ground Truth Output: 2

A.9 Excluded Tasks

Beyond the tasks mentioned above, we have three additional tasks that have been filtered out due to unsatisfactory SIP performance:

- **Bigram Shift** detection: from SentEval (Conneau and Kiela, 2018), which checks whether a bigram in a sentence has been shifted, with binary outcomes (i.e., shifted or not).
- **Subject Number** identification: from SentEval (Conneau and Kiela, 2018), which checks whether the subject of a sentence is “plural” or “singular”.

- **Voice** classification: from Universal Dependencies (Silveira et al., 2014), which checks whether a sentence is in the “active” or “passive” voice. Since no ground-truth labels are available, we use rule-based approaches to annotate the dataset.

B Prompt Template

Here we present the prompt templates we use in our experiments. We use the example inputs in Appendix A for illustration.

B.1 Default Setting

B.1.1 Arithmetic

Task: Solve all the provided arithmetic questions and calculate the sum of all answers.

Instructions:

- Calculate the answer for each arithmetic operation
- Sum all individual answers
- Provide exact values without unnecessary trailing zeros (e.g., "5" not "5.0")
- Prefix negatives with '-' (e.g., "-42")
- For decimal results, keep only necessary decimal places (e.g., "3.14" not "3.140")
- If you only receive one question, the sum is just its answer

Response format:

Return a JSON object with:

- "reasoning": briefly explain your calculation process
- "answer": sum of all answers as a string

Example:

```
{"reasoning": "your approach here",  
  "answer": "42"}
```

```
=== Here are the arithmetic questions  
-> ===
```

```
Question 1: What is the difference  
-> between -2 and 251860?
```

```
Question 2: -9259432 + 1
```

```
Question 3: What is 1141.09 less than 1?  
...
```

```
=== End of arithmetic questions ===
```

B.1.2 Category

Task: Count how many of the provided
→ news articles belong to the 'tech'
→ category.

Background:

- Each news article belongs to one of 5
→ categories:
- * business
- * entertainment
- * politics
- * sport
- * tech
- You need to classify each article
→ based on its content and context

Instructions:

- Read each news article carefully
- Identify which category each article
→ belongs to based on the content
- Count how many articles belong to the
→ 'tech' category
- Do not count articles from other
→ categories (business, entertainment,
→ politics, sport)
- If you only receive one article,
→ return 1 if it's tech, else return 0

Response format:

Return a JSON object with:

- "reasoning": briefly explain your
→ approach to classifying news
→ categories and how you counted
- "answer": integer count of tech news
→ articles

Example:

```
{"reasoning": "your approach here",  
  "answer": 42}
```

=== Here are the news articles ===

Article 1: german business confidence
→ slides german business confidence
→ fell in february knocking hopes of a
→ speedy recovery in europe s largest
→ economy...

Article 2: bbc poll indicates economic
→ gloom citizens in a majority of
→ nations surveyed in a bbc world
→ service poll believe the world
→ economy is worsening...

Article 3: lifestyle governs mobile
→ choice faster better or funkier
→ hardware alone is not going to help
→ phone firms sell more handsets
→ research suggests...

...

=== End of news articles ===

B.1.3 Language

Task: Count how many of the provided
→ paragraphs are written in English.

Background:

- The paragraphs are written in one of
→ four languages:
- English (label 0)
 - Chinese (label 1)
 - Persian (label 2)
 - Spanish (label 3)

Instructions:

- Read each paragraph carefully
- Identify the language of each
→ paragraph
- Count how many paragraphs are written
→ in English
- Do not count paragraphs in other
→ languages
- If you only receive one paragraph,
→ return 1 if it's English, else
→ return 0

Response format:

Return a JSON object with:

- "reasoning": briefly explain your
→ approach to identifying English
→ paragraphs
- "answer": integer count of English
→ paragraphs

Example:

```
{"reasoning": "your approach here",  
  "answer": 5}
```

=== Here are the paragraphs ===

Paragraph 1: Nordahl Road is a station
→ served by North County Transit
→ District's SPRINTER light rail
→ line...

Paragraph 2: En Navidad de 1974, poco
→ después de que interpretó la canción
→ en francés película Papillon (Toi
→ qui Regarde la Mer)...

Paragraph 3: A talk by Takis Fotopoulos

- ↪ about the Internationalization of
- ↪ the Capitalist Market Economy and
- ↪ the project of Inclusive
- ↪ Democracy...

...

=== End of paragraphs ===

B.1.4 NER

Task: Count how many times the entity

- ↪ "PERSON" appears across all provided
- ↪ sentences.

Background:

- An entity may consist of multiple
 - ↪ words that form a contiguous
 - ↪ fragment in the text
- You need to first identify entities in
 - ↪ the sentence (named entity
 - ↪ recognition), then count them
- Two entities may appear consecutively
 - ↪ without punctuation or words between
 - ↪ them
- No entity overlaps occur (each word
 - ↪ belongs to at most one entity)
- Some words do not belong to any entity

Entity Definition:

- PERSON: names of people, real or
 - ↪ fictional, but not nominals

Instructions:

- Identify all PERSON entities across
 - ↪ all sentences
- Count the total number of PERSON
 - ↪ entity mentions (not unique
 - ↪ entities, but total occurrences)
- Each distinct mention counts as one
 - ↪ occurrence, even if it refers to the
 - ↪ same person

Response format:

Return a JSON object with:

- "reasoning": briefly explain how you
 - ↪ identified the entities and counted
 - ↪ them
- "answer": integer count of PERSON
 - ↪ entities

Example:

```
{"reasoning": "your approach here",  
  "answer": 5}
```

=== Here are the sentences ===

Sentence 1: we love everything about the
↪ fence .

Sentence 2: i want to hook up with that
↪ girl paige in the brown leather
↪ jacket .

Sentence 3: in addition , there is a
↪ reduction of 22,101 mmbtu which is
↪ the difference between the scada
↪ values (best available) that anita
↪ showed on the february 29th storage
↪ sheet and the " official " february
↪ 29th values that gary wilson
↪ received from mips .

...

=== End of sentences ===

B.1.5 Parity

Task: Count how many of the provided

- ↪ numbers are odd.

Background:

- An odd number is an integer that is
 - ↪ not evenly divisible by 2
- Odd numbers end in 1, 3, 5, 7, or 9
- An even number is an integer that is
 - ↪ evenly divisible by 2
- Even numbers end in 0, 2, 4, 6, or 8

Instructions:

- Check each number to determine if it
 - ↪ is odd or even
- Count how many numbers are odd
- Do not count even numbers
- If you only receive one number, return
 - ↪ 1 if it's odd, else return 0

Response format:

Return a JSON object with:

- "reasoning": briefly explain your
 - ↪ approach to identifying odd numbers
- "answer": integer count of odd numbers

Example:

```
{"reasoning": "your approach here",  
  "answer": 3}
```

=== Here are the numbers ===

Number 1: 18010

Number 2: 10160

Number 3: 89449

...

=== End of numbers ===

B.1.6 Sentiment

Task: Count how many of the provided
→ movie reviews are positive.

Instructions:

- Each review has a sentiment: positive
→ or negative
- Count only the reviews with positive
→ sentiment
- Return the total count of positive
→ reviews
- If you only receive one review, return
→ 1 if it's positive, else return 0

Response format:

Return a JSON object with:

- "reasoning": briefly explain how you
→ solved this task
- "answer": integer count of positive
→ reviews

Example:

```
{"reasoning": "your approach here",  
  "answer": 42}
```

=== Here are the movie reviews ===

Review 1: High Crimes is a cinematic
→ misdemeanor , a routine crime
→ thriller remarkable only for its
→ lack of logic and misuse of two fine
→ actors , Morgan Freeman and Ashley
→ Judd .

Review 2: One of the worst movies of the
→ year .

Review 3: A mix of gritty realism ,
→ crisp storytelling and radiant
→ compassion that effortlessly draws
→ you in .

...

=== End of movie reviews ===

B.1.7 Word

Task: Count how many times the word
→ "women" appears in the provided
→ tweets.

Instructions:

- Search is case-insensitive (e.g.,
→ "women", "Women", "WOMEN" all count)

- Count occurrences that include the
→ substring "women" (e.g., "women",
→ "womens", "women's", "womenfolk")
- Do not count forms that lack the
→ substring "women" (e.g., "woman",
→ "womankind")
- Count all occurrences across all
→ tweets, not just unique tweets
- If you only receive one tweet, just
→ return the occurrence for that tweet

Response format:

Return a JSON object with:

- "reasoning": briefly explain how you
→ counted the matches
- "answer": integer total count

Example:

```
{"reasoning": "your approach here",  
  "answer": 42}
```

=== Here are the tweets ===

Tweet 1: IF FEMINISTS WERE HONEST "I
→ want a worldwide matriarchal
→ dictatorship with all men enslaved
→ to women" \#GamerGate \#SemST

Tweet 2: What the fuck do women even do?
→ I mean seriously they're just
→ useless other than sex.

Tweet 3: DEAR FEMINISTS Start asking
→ for accountability from man-haters
→ instead of shielding them for
→ convenient concealment. \#SemST

...

=== End of tweets ===

B.1.8 WSD

Task: Count how many paragraphs contain
→ the word "apple" referring to the
→ company (Apple Inc.), not the fruit.

Background:

- Each paragraph contains exactly one
→ occurrence of the word "apple"
→ (case-insensitive, as a complete
→ word, not as part of another word)
- This "apple" can mean either:
 - * The company: Apple Inc., the
→ technology company
 - * The fruit: the edible fruit that grows
→ on apple trees

- You need to determine the meaning of
 - ↳ "apple" in each paragraph based on
 - ↳ context

Instructions:

- Read each paragraph carefully
- Identify whether "apple" refers to the
 - ↳ company or the fruit based on
 - ↳ contextual clues
- Count how many paragraphs where
 - ↳ "apple" means the company (Apple
 - ↳ Inc.)
- Do not count paragraphs where "apple"
 - ↳ means the fruit
- If you only receive one paragraph,
 - ↳ return 1 if it means the company,
 - ↳ else return 0

Response format:

Return a JSON object with:

- "reasoning": briefly explain your
 - ↳ approach to disambiguating the word
 - ↳ sense and how you counted
- "answer": integer count of paragraphs
 - ↳ where "apple" means the company

Example:

```
{"reasoning": "your approach here",
  "answer": 42}
```

=== Here are the paragraphs ===

Paragraph 1: both seasons are available
 ↳ for download from apple 's itunes
 ↳ store .

Paragraph 2: in klayman ii , the
 ↳ plaintiffs sued the same government
 ↳ defendants and in addition ,
 ↳ facebook , yahoo! , google ,
 ↳ microsoft , youtube , aol , paltalk
 ↳ , skype , sprint , at&t , apple
 ↳ again alleging the bulk metadata
 ↳ collection violates the first ,
 ↳ fourth and fifth amendment and
 ↳ constitutes divulgence of
 ↳ communication records in violation
 ↳ of section 2702 of stored
 ↳ communications act .

Paragraph 3: description alongside dried
 ↳ pears the filling also contains
 ↳ raisin , walnut and other dried
 ↳ fruit such as apple or figs .

...

=== End of paragraphs ===

B.2 Instance-Level Setting

B.2.1 Arithmetic

Task: Solve each of the provided

- ↳ arithmetic questions and calculate
- ↳ the sum of all answers.

Instructions:

- Calculate the answer for each
 - ↳ arithmetic question
- Use 1-based indexing for question
 - ↳ numbers ("1", "2", "3", ...)
- Provide exact values without
 - ↳ unnecessary trailing zeros (e.g.,
 - ↳ "5" not "5.0")
- Prefix negatives with '-' (e.g.,
 - ↳ "-42")
- For decimal results, keep only
 - ↳ necessary decimal places (e.g.,
 - ↳ "3.14" not "3.140")
- Sum all individual answers

Response format:

Return a JSON object with:

- One key per question: "1", "2", "3",
 - ↳ ... , mapping to the answer as a
 - ↳ string
- "sum": sum of all answers as a string
- "reasoning": brief explanation of your
 - ↳ approach

Example for 3 questions:

```
{"1": "8", "2": "-30", "3": "5.6",
  "sum": "-16.4", "reasoning": "your
  approach here"}
```

=== Here are the arithmetic questions

↳ ===

Question 1: What is the difference

- ↳ between -2 and 251860?

Question 2: $-9259432 + 1$

Question 3: What is 1141.09 less than 1?

...

=== End of arithmetic questions ===

B.2.2 Category

Task: For each news article, classify it

- ↳ into one of 5 categories, then
- ↳ provide total counts for each
- ↳ category.

Background:

- Each news article belongs to one of 5
 - ↪ categories:
- * business (label 0)
- * entertainment (label 1)
- * politics (label 2)
- * sport (label 3)
- * tech (label 4)
- You need to classify each article
 - ↪ based on its content and context

Instructions:

- Read each news article carefully and
 - ↪ classify it into the most
 - ↪ appropriate category
- For each article, assign:
 - * 0 if it's business news
 - * 1 if it's entertainment news
 - * 2 if it's politics news
 - * 3 if it's sport news
 - * 4 if it's tech news
- Count the total number of articles for
 - ↪ each category
- Provide classification for each
 - ↪ article along with summary counts

Response format:

- Return a JSON object with:
- One key per article: "1", "2", "3",
 - ↪ . . . , mapping to the category label
 - ↪ (0-4)
 - "business": integer count of articles
 - ↪ classified as business
 - "entertainment": integer count of
 - ↪ articles classified as entertainment
 - "politics": integer count of articles
 - ↪ classified as politics
 - "sport": integer count of articles
 - ↪ classified as sport
 - "tech": integer count of articles
 - ↪ classified as tech
 - "reasoning": brief explanation of your
 - ↪ approach to news classification

Example for 3 articles:

```
{"1": 4, "2": 0, "3": 3, "business": 1,  
↪ "entertainment": 0, "politics": 0,  
↪ "sport": 1, "tech": 1, "reasoning":  
↪ "your approach here"}
```

=== Here are the news articles ===

Article 1: german business confidence
↪ slides german business confidence
↪ fell in february knocking hopes of a
↪ speedy recovery in europe s largest
↪ economy...

Article 2: bbc poll indicates economic
↪ gloom citizens in a majority of
↪ nations surveyed in a bbc world
↪ service poll believe the world
↪ economy is worsening...

Article 3: lifestyle governs mobile
↪ choice faster better or funkier
↪ hardware alone is not going to help
↪ phone firms sell more handsets
↪ research suggests...

...

=== End of news articles ===

B.2.3 Language

Task: For each paragraph, identify which
↪ language it is written in, then
↪ provide summary counts for all
↪ categories.

Background:

- The paragraphs are written in one of
- ↪ four languages:
- English (label 0)
 - Chinese (label 1)
 - Persian (label 2)
 - Spanish (label 3)

Instructions:

- Read each paragraph carefully and
 - ↪ identify its language
- Classify each paragraph using the
 - ↪ labels:
- * 0 = English
- * 1 = Chinese
- * 2 = Persian
- * 3 = Spanish
- Provide the classification for each
 - ↪ individual paragraph
- Also provide summary counts for all
 - ↪ four language categories

Response format:

- Return a JSON object with:
- One key per paragraph: "1", "2", "3",
 - ↪ . . . , mapping to the classification
 - ↪ (0, 1, 2, or 3)
 - "english": integer count of paragraphs
 - ↪ classified as English

- "chinese": integer count of paragraphs
 - ↳ classified as Chinese
- "persian": integer count of paragraphs
 - ↳ classified as Persian
- "spanish": integer count of paragraphs
 - ↳ classified as Spanish
- "reasoning": brief explanation of your approach

Example for 5 paragraphs:

```
{
  "1": 0, "2": 1, "3": 2, "4": 3, "5": 0,
  "english": 2, "chinese": 1,
  "persian": 1, "spanish": 1,
  "reasoning": "your approach here"
}
```

=== Here are the paragraphs ===

Paragraph 1: Nordahl Road is a station
 ↳ served by North County Transit
 ↳ District's SPRINTER light rail
 ↳ line...

Paragraph 2: En Navidad de 1974, poco
 ↳ después de que interpretó la canción
 ↳ en francés película Papillon (Toi
 ↳ qui Regarde la Mer)...

Paragraph 3: A talk by Takis Fotopoulos
 ↳ about the Internationalization of
 ↳ the Capitalist Market Economy and
 ↳ the project of Inclusive
 ↳ Democracy...

...

=== End of paragraphs ===

B.2.4 NER

Task: Count how many times the entity
 ↳ "PERSON" appears in each sentence
 ↳ and provide the total count.

Background:

- An entity may consist of multiple
 - ↳ words that form a contiguous
 - ↳ fragment in the text
- You need to first identify entities in
 - ↳ each sentence (named entity
 - ↳ recognition), then count them
- Two entities may appear consecutively
 - ↳ without punctuation or words between
 - ↳ them
- No entity overlaps occur (each word
 - ↳ belongs to at most one entity)
- Some words do not belong to any entity

Entity Definition:

- PERSON: names of people, real or
 - ↳ fictional, but not nominals

Instructions:

- Identify all PERSON entities in each
 - ↳ sentence
- Count the number of PERSON entity
 - ↳ mentions in each sentence separately
- Provide the count for each sentence
 - ↳ along with the total count across
 - ↳ all sentences
- Each distinct mention counts as one
 - ↳ occurrence, even if it refers to the
 - ↳ same person

Response format:

Return a JSON object with:

- One key per sentence: "1", "2", "3",
 - ↳ . . . , mapping to the integer count
 - ↳ of PERSON mentions in that sentence
- "total": total count of PERSON
 - ↳ entities across all sentences
- "reasoning": brief explanation of how
 - ↳ you identified and counted the
 - ↳ PERSON entities

Example:

```
{
  "1": 0, "2": 2, "3": 1, "total": 3,
  "reasoning": "your approach here"
}
```

=== Here are the sentences ===

Sentence 1: we love everything about the
 ↳ fence .

Sentence 2: i want to hook up with that
 ↳ girl paige in the brown leather
 ↳ jacket .

Sentence 3: in addition , there is a
 ↳ reduction of 22,101 mmbtu which is
 ↳ the difference between the scada
 ↳ values (best available) that anita
 ↳ showed on the february 29th storage
 ↳ sheet and the " official " february
 ↳ 29th values that gary wilson
 ↳ received from mips .

...

=== End of sentences ===

B.2.5 Parity

Task: For each number, identify whether
 ↳ it is odd or even, then provide
 ↳ summary counts for both categories.

Background:

- An odd number is an integer that is
 - ↪ not evenly divisible by 2
- Odd numbers end in 1, 3, 5, 7, or 9
 - ↪ (label 1)
- An even number is an integer that is
 - ↪ evenly divisible by 2
- Even numbers end in 0, 2, 4, 6, or 8
 - ↪ (label 0)

Instructions:

- Check each number to determine if it
 - ↪ is odd or even
- Classify each number using the labels:
 - * 1 = odd
 - * 0 = even
- Provide the classification for each
 - ↪ individual number
- Also provide summary counts for both
 - ↪ odd and even categories

Response format:

Return a JSON object with:

- One key per number: "1", "2", "3",
 - ↪ . . . , mapping to the classification
 - ↪ (0 or 1)
- "odd": integer count of numbers
 - ↪ classified as odd
- "even": integer count of numbers
 - ↪ classified as even
- "reasoning": brief explanation of your
 - ↪ approach

Example for 5 numbers:

```
{"1": 0, "2": 1, "3": 0, "4": 1, "5": 1,  
  ↪ "odd": 3, "even": 2, "reasoning":  
  ↪ "your approach here"}
```

=== Here are the numbers ===

Number 1: 18010

Number 2: 10160

Number 3: 89449

...

=== End of numbers ===

B.2.6 Sentiment

Task: For each movie review, classify

- ↪ whether it is positive or negative,
- ↪ then provide summary counts.

Instructions:

- Classify each review as either:
- 0 = negative sentiment

- 1 = positive sentiment
- Provide the classification for each
 - ↪ individual review
- Also provide summary counts for
 - ↪ negative and positive reviews

Response format:

Return a JSON object with:

- One key per review: "1", "2", "3",
 - ↪ . . . , mapping to the classification
 - ↪ (0 or 1)
- "negative": integer count of reviews
 - ↪ classified as negative
- "positive": integer count of reviews
 - ↪ classified as positive
- "reasoning": brief explanation of your
 - ↪ approach to classification

Example for 3 reviews:

```
{"1": 1, "2": 0, "3": 1, "negative": 1,  
  ↪ "positive": 2, "reasoning": "your  
  ↪ approach here"}
```

=== Here are the movie reviews ===

Review 1: High Crimes is a cinematic

↪ misdemeanor , a routine crime
↪ thriller remarkable only for its
↪ lack of logic and misuse of two fine
↪ actors , Morgan Freeman and Ashley
↪ Judd .

Review 2: One of the worst movies of the
↪ year .

Review 3: A mix of gritty realism ,
↪ crisp storytelling and radiant
↪ compassion that effortlessly draws
↪ you in .

...

=== End of movie reviews ===

B.2.7 Word

Task: For each tweet, count how many

- ↪ times the word "women" appears, then
- ↪ provide the total count.

Instructions:

- Search is case-insensitive (e.g.,
 - ↪ "women", "Women", "WOMEN" all count)
- Count occurrences that include the
 - ↪ substring "women" (e.g., "women",
 - ↪ "womens", "women's", "womenfolk")
- Do not count forms that lack the
 - ↪ substring "women" (e.g., "woman",
 - ↪ "womankind")

- Count occurrences in each tweet
 - ↳ separately
- Provide the per-tweet counts plus the
 - ↳ overall total

Response format:

Return a JSON object with:

- One key per tweet: "1", "2", "3",
 - ↳ . . . , mapping to the count of
 - ↳ "women" in that tweet (integer)
- "total": integer representing the
 - ↳ total count of "women" across all
 - ↳ tweets
- "reasoning": brief explanation of your
 - ↳ counting approach

Example for 3 tweets:

```
{ "1": 2, "2": 0, "3": 1, "total": 3,
  "reasoning": "your approach here" }
```

=== Here are the tweets ===

Tweet 1: IF FEMINISTS WERE HONEST "I
 ↳ want a worldwide matriarchal
 ↳ dictatorship with all men enslaved
 ↳ to women" \#GamerGate \#SemST

Tweet 2: What the fuck do women even do?
 ↳ I mean seriously they're just
 ↳ useless other than sex.
 ↳ \#womensrights \#Feminist \#SemST

Tweet 3: DEAR FEMINISTS Start asking
 ↳ for accountability from man-haters
 ↳ instead of shielding them for
 ↳ convenient concealment. \#SemST

...

=== End of tweets ===

B.2.8 WSD

Task: For each paragraph, identify
 ↳ whether the word "apple" refers to
 ↳ the company or the fruit, then
 ↳ provide total counts.

Background:

- Each paragraph contains exactly one
 - ↳ occurrence of the word "apple"
 - ↳ (case-insensitive, as a complete
 - ↳ word, not as part of another word)
- This "apple" can mean either:
 - * The company (label 0): Apple Inc., the
 - ↳ technology company
 - * The fruit (label 1): the edible fruit
 - ↳ that grows on apple trees

- You need to determine the meaning of
 - ↳ "apple" in each paragraph based on
 - ↳ context

Instructions:

- Read each paragraph carefully and
 - ↳ classify the meaning of "apple"
- For each paragraph, assign:
 - * 0 if "apple" means the company (Apple
 - ↳ Inc.)
 - * 1 if "apple" means the fruit
- Count the total number of paragraphs
 - ↳ for each category
- Provide classification for each
 - ↳ paragraph along with summary counts

Response format:

Return a JSON object with:

- One key per paragraph: "1", "2", "3",
 - ↳ . . . , mapping to either 0 (company)
 - ↳ or 1 (fruit)
- "company": integer count of paragraphs
 - ↳ where "apple" means the company
- "fruit": integer count of paragraphs
 - ↳ where "apple" means the fruit
- "reasoning": brief explanation of your
 - ↳ approach to word sense
 - ↳ disambiguation

Example for 3 paragraphs:

```
{ "1": 0, "2": 1, "3": 0, "company": 2,
  "fruit": 1, "reasoning": "your
  approach here" }
```

=== Here are the paragraphs ===

Paragraph 1: both seasons are available
 ↳ for download from apple 's itunes
 ↳ store .

Paragraph 2: in klayman ii , the
 ↳ plaintiffs sued the same government
 ↳ defendants and in addition ,
 ↳ facebook , yahoo! , google ,
 ↳ microsoft , youtube , aol , paltalk
 ↳ , skype , sprint , at&t , apple
 ↳ again alleging the bulk metadata
 ↳ collection violates the first ,
 ↳ fourth and fifth amendment and
 ↳ constitutes divulgence of
 ↳ communication records in violation
 ↳ of section 2702 of stored
 ↳ communications act .

Task	Agreement (%)
BShift	69.3
SubjNum	82.5
Voice	78.3

Table 5: Task filtering across tasks that have been removed.

Paragraph 3: description alongside dried
↪ pears the filling also contains
↪ raisin , walnut and other dried
↪ fruit such as apple or figs .
...
=== End of paragraphs ===

C Single Instance Filtering Results

We note that not all models in our evaluation are included in this SIP filtering analysis, as the filtering procedure is defined based on a fixed subset of comparison models; however, the excluded models generally exhibit stronger performance in the MIP setting, and are likely to perform similarly under the SIP settings (e.g., for this filtering, we used the smaller versions of the closed-source LLMs considered).

C.1 Individual Results

Table 5 reports the agreements for the tasks that are removed, while Table 6 shows the average SIP success rate for each model and task. Although *BShift* achieves a relatively high SIP success rate for the selected models, its agreement across models is low.

C.2 Task Filtering Results

Table 7 shows the retained instance percentage for each task, along with the corresponding maximum and minimum SIP success rates across models.

D Experimental Results

Here we present the success rate and failure breakdown for each model and task.

D.1 Success Rate

D.1.1 Task Success Rate for Models

Figure 8, Figure 9, Figure 10 and Figure 11 show task success rate for each model.

D.1.2 Model Success Rate for Tasks

Figure 12 and Figure 13 show model success rate for each task.

Model	Arithmetic	Language	NER	News	Parity	BShift	SubjNum	Sentiment	Tweets	Voice	WSD	Average
DeepSeek V3	98.3	99.2	95.8	97.6	100.0	92.7	88.3	99.3	99.6	89.5	99.4	96.3
gpt-oss-120b	97.2	99.5	96.2	97.6	100.0	92.6	89.4	98.5	100.0	91.4	99.3	96.5
gpt-oss-20b	97.7	99.5	96.2	97.2	100.0	90.4	89.6	97.8	100.0	90.6	99.1	96.2
Llama 3.3	93.7	99.5	93.6	98.4	100.0	90.9	87.4	99.7	99.5	86.9	99.2	95.3
Llama 4 Maverick	96.7	99.3	94.3	99.2	100.0	90.9	88.9	99.4	99.9	91.7	99.4	96.3
Llama 4 Scout	93.4	97.9	93.9	97.2	100.0	86.9	87.3	99.0	99.5	88.7	99.1	94.8
Mistral NeMo	72.0	96.2	78.1	96.8	75.3	62.9	75.2	83.8	60.3	72.7	98.4	79.2
Qwen3-Instruct	99.5	99.4	97.3	99.2	100.0	91.3	88.9	99.5	99.7	92.8	99.3	97.0
Gemini 2.5 Flash	96.2	99.6	96.5	98.0	100.0	92.8	88.9	99.0	99.4	92.0	99.4	96.5
GPT-5 Nano	98.0	99.5	96.5	97.2	100.0	93.0	89.0	99.2	100.0	91.3	99.2	96.6
Grok 4 Fast	97.9	99.5	97.2	98.8	100.0	95.7	88.6	99.5	100.0	90.2	99.4	97.0

Table 6: SIP success rate (%) across tasks. Open-weight models are listed first, followed by closed-source models. Underlined rows and columns indicate the tasks and LLMs that are excluded because they do not satisfy the filtering criteria described in Section 3.2.

Task	Agreement (%)	Max (%)	Min (%)
Arithmetic	89.0	99.5	93.7
Category	94.8	99.2	97.2
Language	98.8	99.6	99.2
NER	87.6	97.3	93.6
Parity	100.0	100.0	100.0
Sentiment	96.4	99.7	97.8
Word	98.3	100.0	99.4
WSD	98.6	99.4	99.1

Table 7: Task filtering results showing agreement (i.e., the percentage of instances for which all LLMs produce correct SIP predictions), as well as the maximum and minimum SIP success rates across all LLMs (the actual agreement and minimum rate for *Parity* is 99.96%).

D.2 Failure Breakdown

D.2.1 Failure Breakdown for Models

Figure 14, Figure 15 and Figure 16 show failure breakdown for models, averaged across all tasks.

D.2.2 Failure Breakdown for Tasks

Figure 17 and Figure 18 show failure breakdown for tasks, averaged across all models.

D.3 Context Length

D.3.1 Default and Augmented Success Rate

Figure 19 compares average performance across all tasks and models between the default and artificially augmented settings, as the context length increases.

D.3.2 Default and Augmented Success Rate for Models

Figure 20, Figure 21 and Figure 22 show task success rate for each model.

D.3.3 Default and Augmented Success Rate for Tasks

Figure 23 and Figure 24 show model success rate for each task.

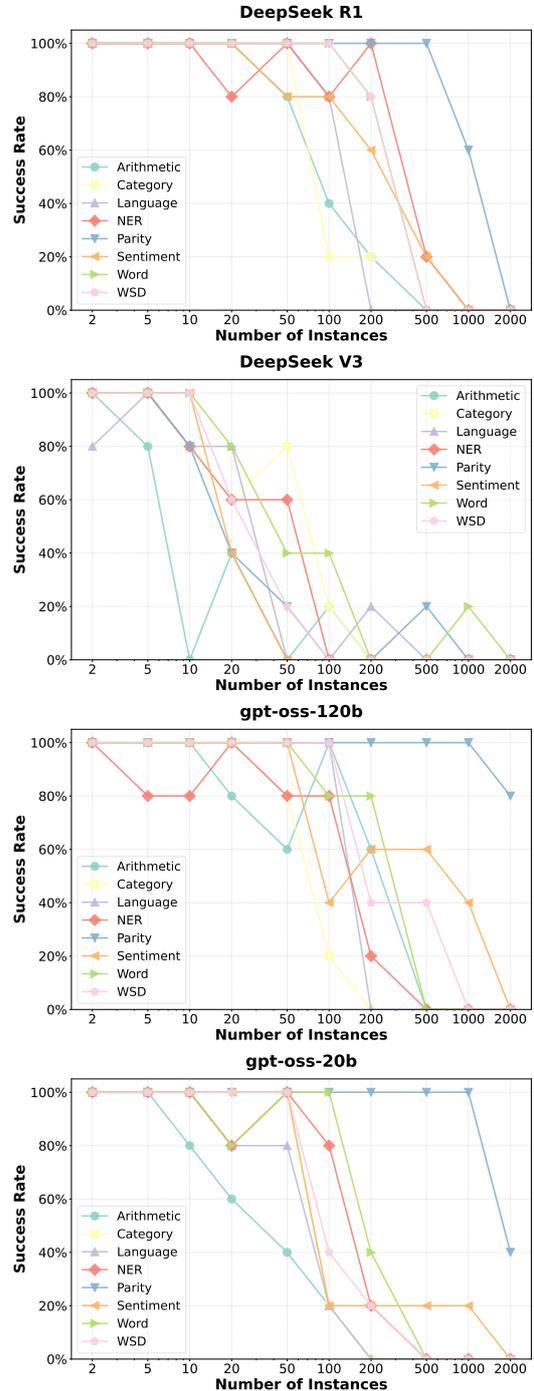


Figure 8: Success rate of models.

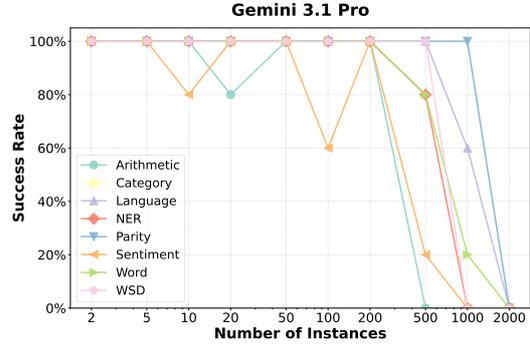
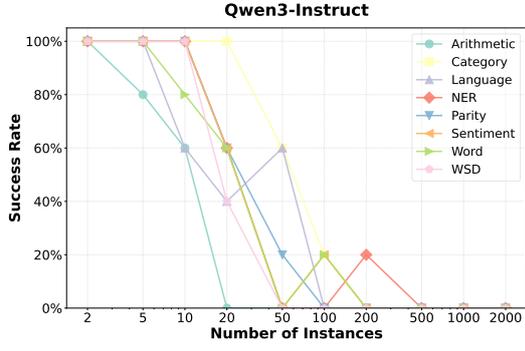
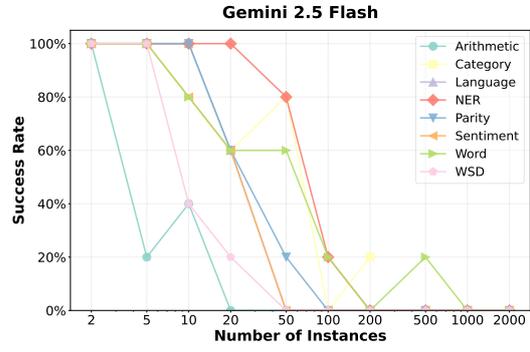
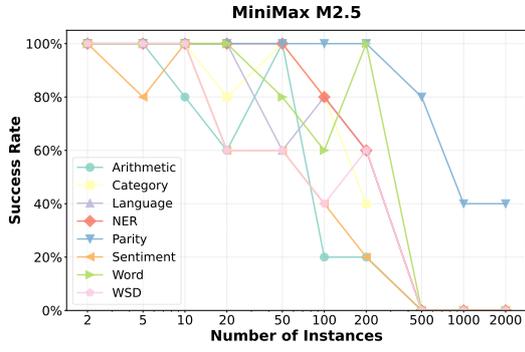
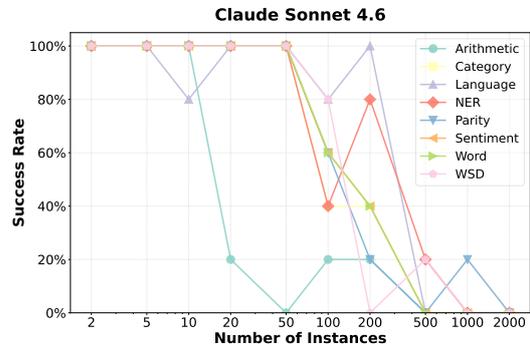
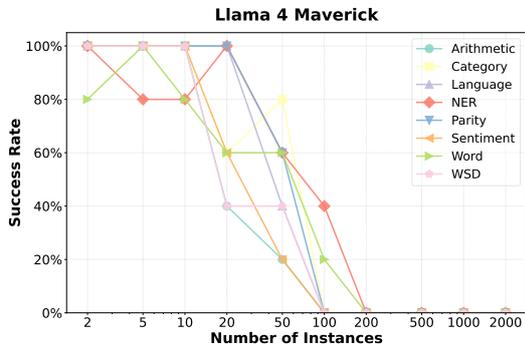
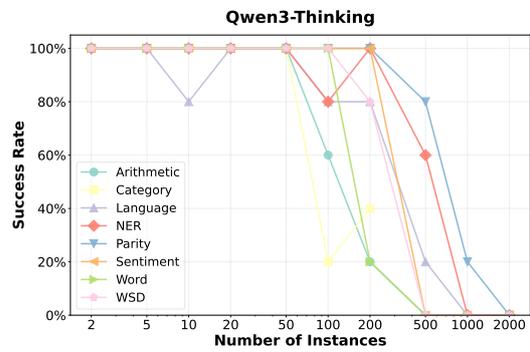
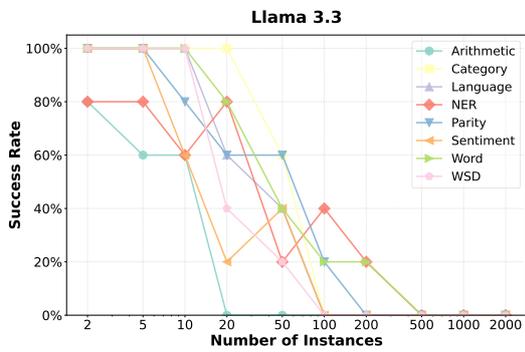


Figure 9: Success rate of models.

Figure 10: Success rate of models.

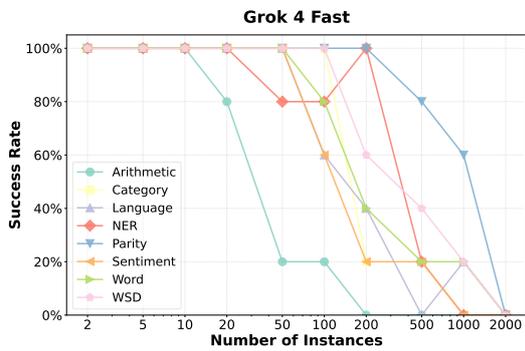
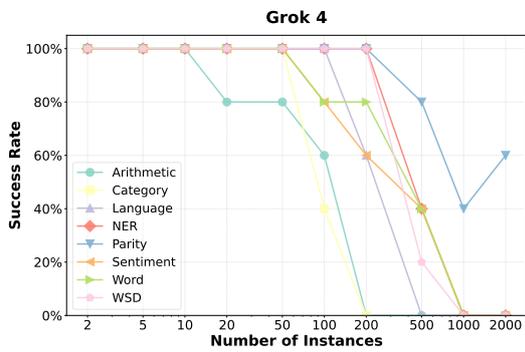
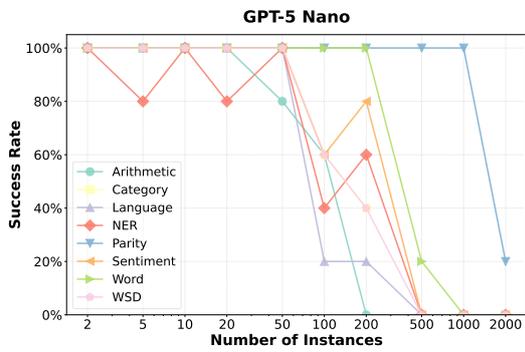
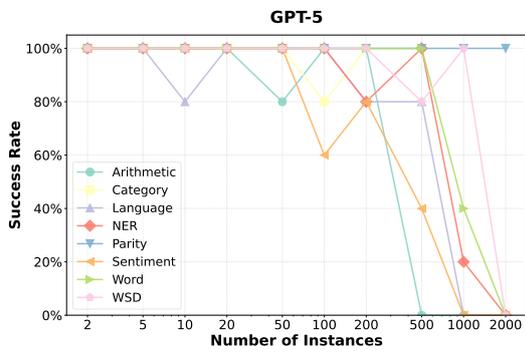


Figure 11: Success rate of models.

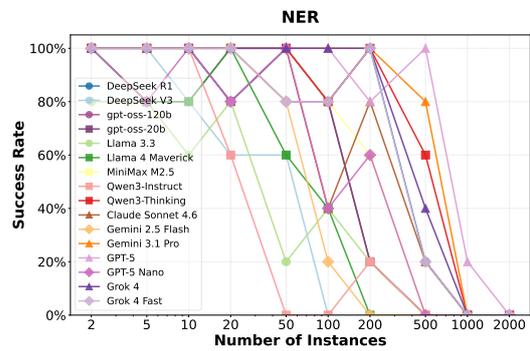
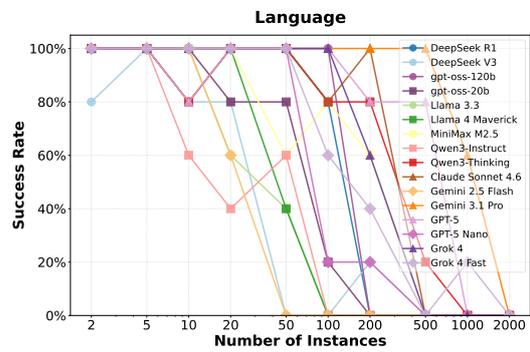
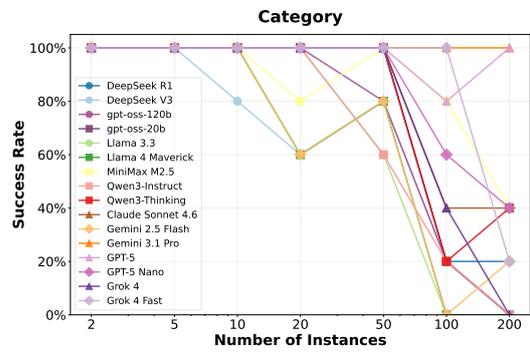
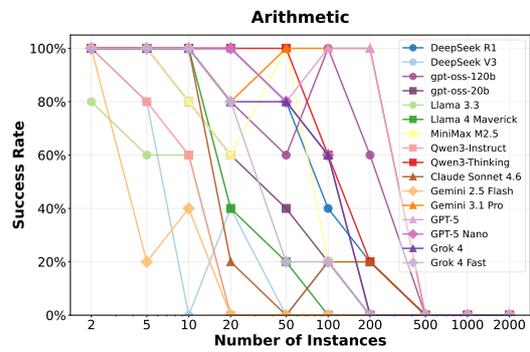


Figure 12: Model success rate for tasks.

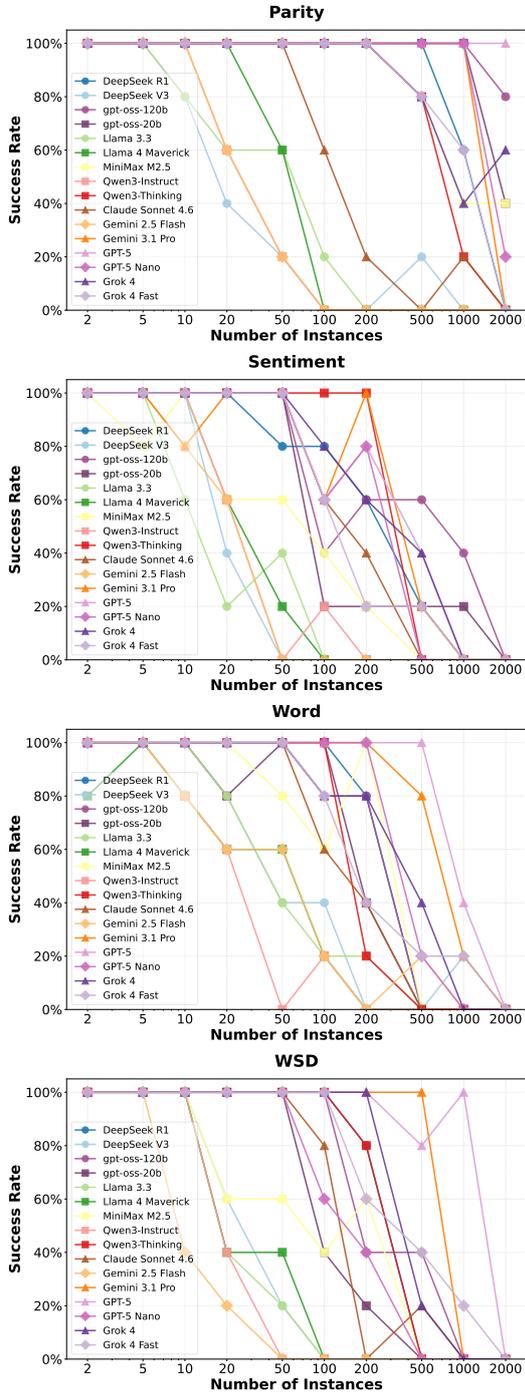


Figure 13: Model success rate for tasks.

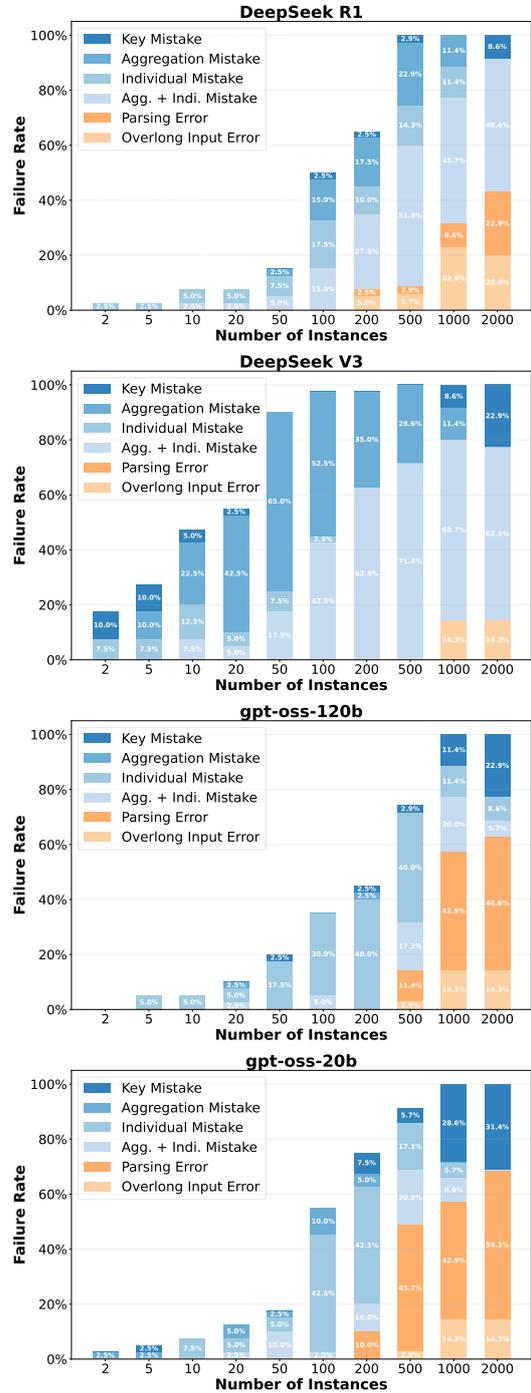


Figure 14: Failure breakdown for models.

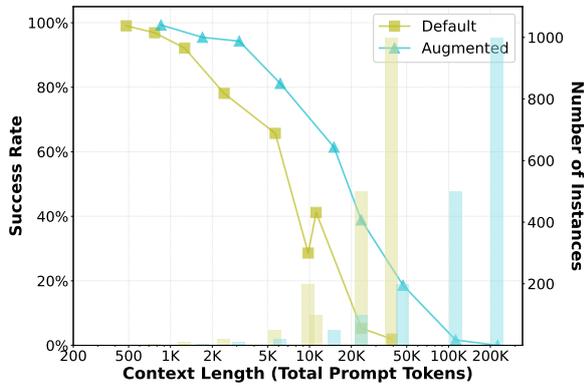


Figure 19: Success rate (lines) and the number of instances (bars) in the artificial length setting as total prompt token length increases. Error bars indicate standard deviation across five random seeds.

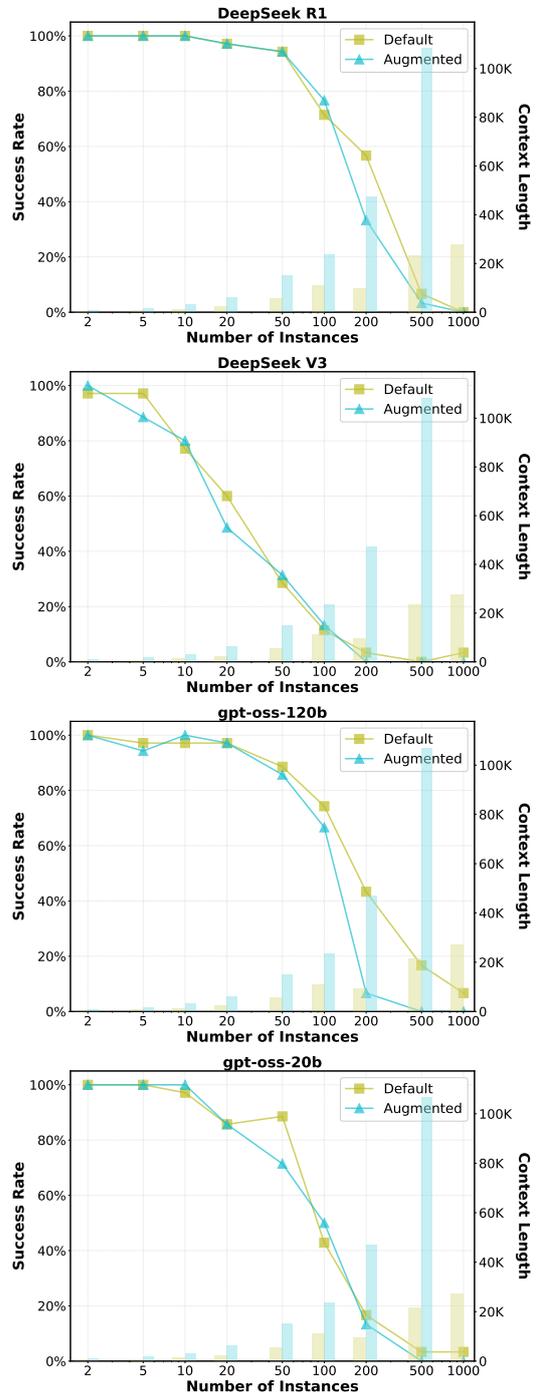


Figure 20: Success rate of models.

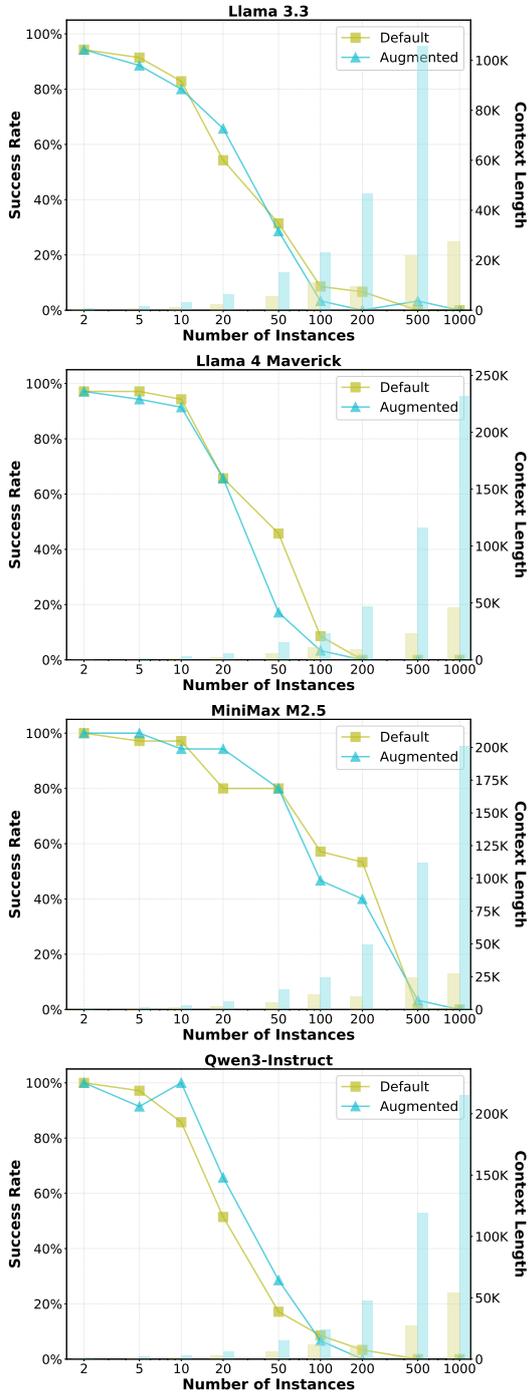


Figure 21: Success rate of models.

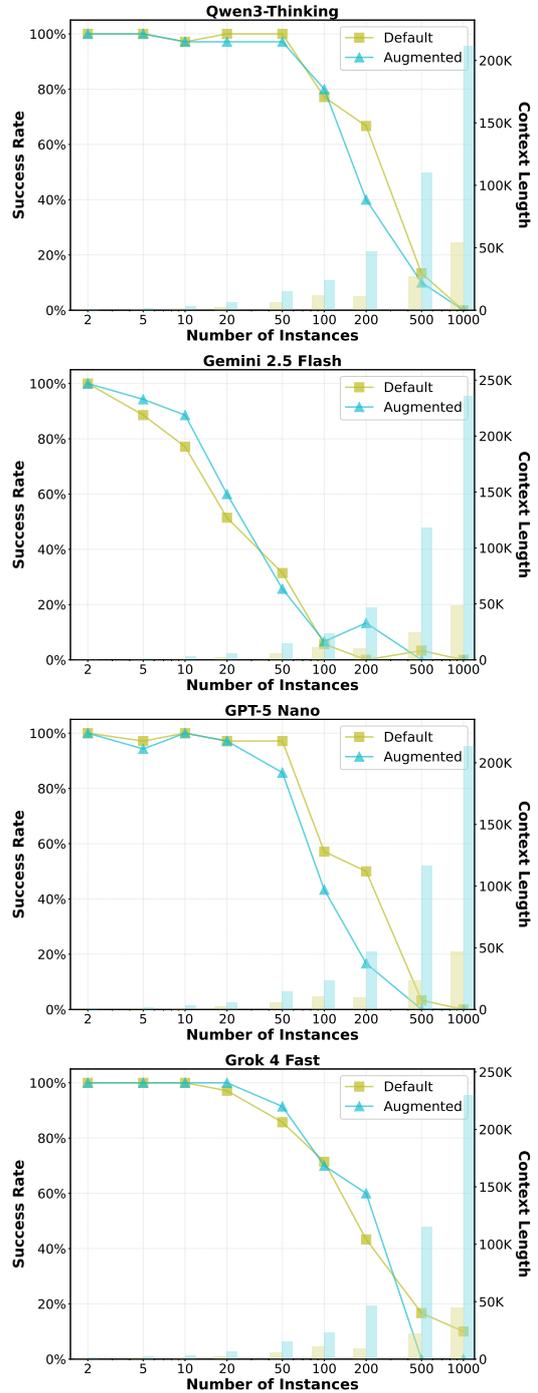


Figure 22: Success rate of models.

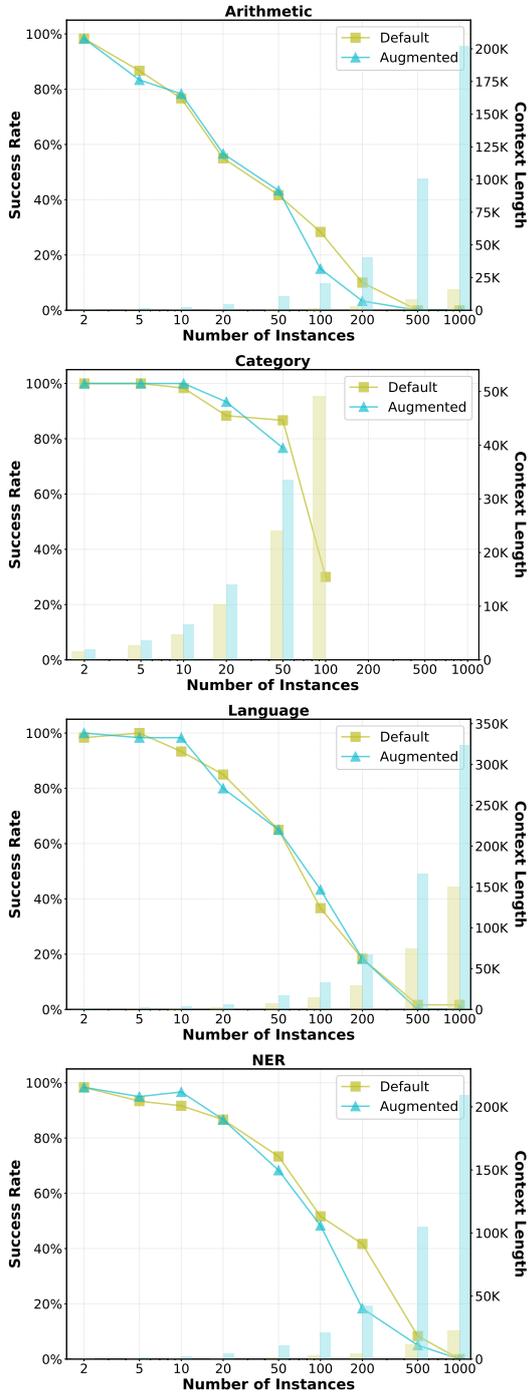


Figure 23: Success rate for tasks.

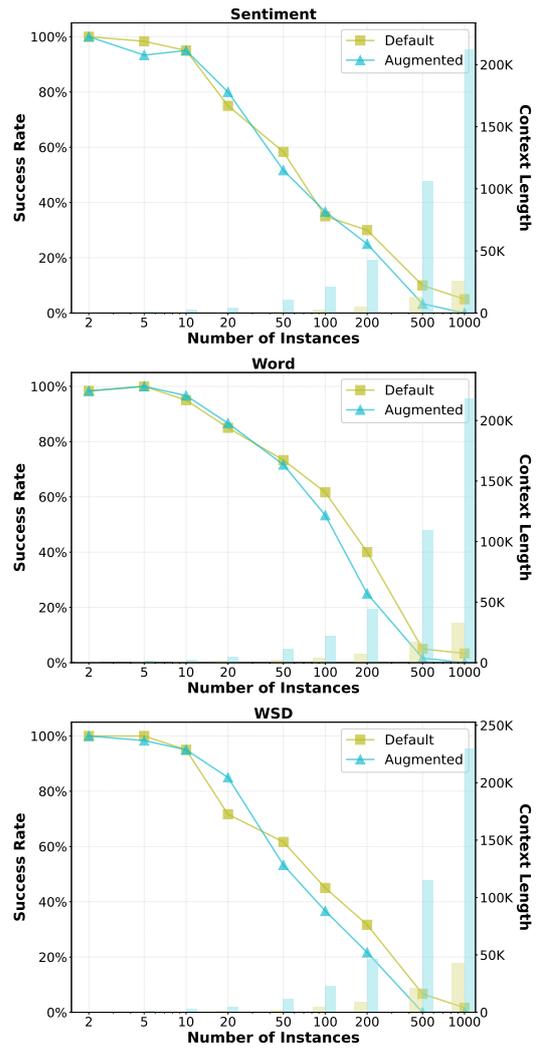


Figure 24: Success rate for tasks.