

Can LLMs subtract numbers?

Mayank Jobanputra¹, Nils Philipp Walter², Maitrey Mehta³, Blerta Veseli¹,
Evan Parker Kelly Chapple¹, Yifan Wang¹, Sneha Chetani¹,
Ellie Pavlick⁴, Antonio Vergari⁵, Vera Demberg¹

¹Saarland University ²CISPA Helmholtz Center for Information Security

³University of Utah ⁴Brown University ⁵University of Edinburgh

mayank@lst.uni-saarland.de

Abstract

We present a systematic study of subtraction in large language models (LLMs). While prior benchmarks emphasize addition and multiplication, subtraction has received comparatively little attention despite being structurally distinct as a non-commutative operation. We evaluate eight pretrained LLMs spanning four families on addition and subtraction problems. Our experiments reveal that subtraction accuracy lags behind addition by a wide margin. We find that the errors for $(a - b)$ are concentrated in cases where $(a < b)$. In such cases, LLMs frequently produce the correct magnitude but omit the negative sign. Probing analyses show that LLMs internally encode whether results should be negative, yet this information is often not reflected in generated outputs. We further test well-known techniques such as few-shot learning and instruction-tuning to see if they can improve the LLMs’ performance. Our results suggest that while few-shot prompting yields modest gains, the instruction-tuned models achieve near-perfect accuracies in generating the negative sign. Together, these findings provide a clearer characterization of the limitations and recoverability of LLMs’ arithmetic capabilities in subtraction.

1 Introduction

Large language models (LLMs) show strong performance on arithmetic tasks, particularly when aided by chain-of-thought (CoT) prompting and its variants (Wei et al., 2022; Kojima et al., 2022; Imani et al., 2023; Li et al., 2025). However, prior work and widely used benchmarks (Balunović et al., 2025; Yang et al., 2025b; Hendrycks et al., 2021) concentrate on addition and multiplication, leaving *subtraction* comparatively underexplored. This gap raises the question of how well LLMs handle subtraction, especially in a zero-shot setting.

Subtraction differs from addition in an important structural way as it is a non-commutative operation ($a - b \neq b - a$). This makes the order

of operands decisive for subtraction. Moreover, subtraction accentuates better positional representations to facilitate accurate borrowing. Borrowing requires robust tracking of digit positions across potentially long sequences — a setting where transformers trained from scratch are known to struggle (McLeish et al., 2024). These factors pose additional demands unique to subtraction, making it unclear whether LLMs can sustain their strong performance on addition (Zhou et al., 2024; Nikankin et al., 2025) when applied to subtraction.

In this work, we systematically evaluate subtraction across LLM families and sizes, benchmarking against addition to establish a clear performance baseline. We evaluate eight pretrained LLMs of different sizes from four model families, Gemma-2 (Riviere et al., 2024), Qwen3 (Yang et al., 2025a), OLMo-2 (OLMo et al., 2024), and Llama-3 (Grattafiori et al., 2024).

We find a surprising disparity between the additive and subtractive capabilities of LLMs. Our findings reveal that pretrained LLMs struggle to generate correct answers when subtraction results are negative. Our main findings are:

a) Under matched complexity, subtraction performance consistently lags addition in state-of-the-art open-source LLMs. In some cases, LLMs reach perfect accuracy on addition while only about half that level on subtraction. (§ 3.1) b) Subtraction $a - b$ is disproportionately error-prone when $(a < b)$, with a dominant failure mode in which LLMs produce the correct magnitude but omit the negative sign. (§ 3.2, 3.3) c) Our probing results suggest that LLMs internally encode when the result should be negative, but they often fail to surface this at decoding time. (§ 3.4) d) Instruction tuning improves subtraction performance to levels comparable with addition in a zero-shot setting. However, adding a few-shot examples yields modest and inconsistent improvements across pretrained LLMs. (§ 4) This candidates subtraction as an inherently harder task

than addition, and can help shed light on the different regimes (i.e. pretrained or instruction-tuned) in which LLMs learn and reason.

2 Experimental Setup

Following prior work (Zhou et al., 2024; Nikankin et al., 2025), we primarily focus on *single-token* numbers for our analysis. However, we also analyze *multi-token* numbers in Appendix A that show similar results.

2.1 Data generation

To evaluate LLMs on the integer subtraction task, we generate synthetic datasets using controlled settings described below.

Sampling operands a and b . We sample operands a and b uniformly from the range of numeric values that are represented as single tokens by each LLM’s tokenizer (see *Tokenizer Range* in Table 1). Consequently, each LLM is evaluated only on problems involving numbers within its own tokenizer range. For every sampled pair (a, b) , we generate arithmetic problems and categorize them based on the relationship between operands (e.g., $a > b$, $a < b$, $a = b$). The final dataset (see #Used Samples in Table 1) is balanced, containing an equal number of problems with $a > b$ and $a < b$. For our work, we do not make use of operands where $a = b$.

Prompt Variants. To ensure the generalizability of our results and minimize the risk of spurious correlations, we use five different prompt formats, ranging from minimal equation style input to verbose template style input. We provide the exact prompt templates in the Appendix D.

Zero-shot vs. n-shot. *Zero-shot prompts* contained only the query equation formatted under one of the five variants. *N-shot prompts* include solved examples (up to ten) before the query, allowing us to probe in-context learning abilities of LLMs.

We use the same procedure to generate the data for both addition (+) and subtraction (−) operators. This procedure yields a balanced dataset, and systematically varied across operators, operand order, and prompt phrasings. We provide full dataset statistics in Table 1.

2.2 LLM Inference

For pretrained LLMs, we use greedy sampling with a temperature of 0 and sample up to 20 new tokens. For instruction-tuned LLMs, we use the sampling strategy and parameters suggested by the LLM creators. Additionally, we use the default system prompt for each instruction-tuned LLM and sample up to 500 new tokens.

We use a fixed seed across all experiments to ensure reproducibility. We extract the final numerical answer from the LLMs’ generated text using a robust parsing mechanism. We run all experiments on 4x H100 GPUs using vLLM (Kwon et al., 2023) without any quantization.

3 Results

In this section, we analyze LLM performance on subtraction, using addition as a reference point. We examine how accuracy varies across LLMs, operand order, and analyze the errors that emerge.

3.1 Can LLMs subtract numbers?

We start our exploration by comparing subtraction accuracy with that of addition across prompt variants in a zero-shot setting. Figure 1 shows that while most LLMs achieve near-perfect accuracy on addition, subtraction is substantially harder. For example, Qwen3-8B reaches almost 100% on addition but only around 57% on subtraction. Similarly, OLMo-2-32B scores above 99% on addition, but drops to roughly 56% on subtraction. In contrast, smaller LLMs such as Gemma-2-9B and Llama-3-8B remain poor at both tasks. In general, the performance of the subtraction lags the addition by 30 – 50 points in several LLMs, showing that the LLMs struggle to compute the correct answer for the subtraction. This leads us to explore further if there are any patterns in the incorrect answers.

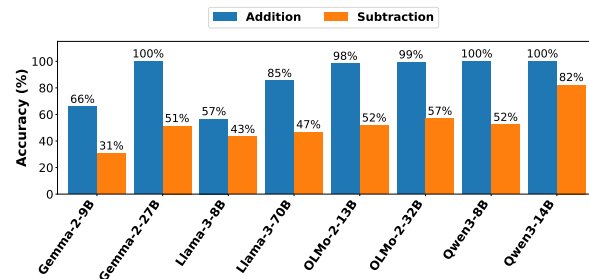


Figure 1: Zero-shot performance of LLMs on addition and subtraction problems (averaged across prompt variants). Subtraction is consistently harder, even for LLMs that perform well on addition.

LLM Family	Tokenizer Range	Total Samples	#Samples $a > b$	#Samples $a < b$	Prompt Variants	#Samples Used
Qwen3	[0, 9]	100	45	45	5	450
Gemma-2	[0, 9]	100	45	45	5	450
Llama-3	[0, 999]	1,000,000	499,500	499,500	5	10,000
OLMo-2	[0, 999]	1,000,000	499,500	499,500	5	10,000

Table 1: Dataset statistics for subtraction experiments across different LLM families. *Tokenizer Range* indicates the continuous range of numeric values that are single tokens for each LLM. *Total Samples* reflects the total number of possible subtraction problems. For LLMs where *Total Samples* exceeds 2000, we uniformly sample 2000 problems from the full set. We then apply our prompt variants to each sampled problem and obtain all the samples for inference (i.e. *#Samples Used*).

3.2 Subtraction with $a > b$ versus $a < b$

We divide our input data into two subsets, $a > b$ and $a < b$, to see if there is any concrete error pattern exists there. Figure 2 suggests a strong asymmetry in performance here as nearly all LLMs succeed in performing $a - b$ when $a > b$, but accuracy collapses for $a < b$. For instance, Qwen3-8B, Gemma-2-27B, and Llama-3.1-70B achieve near-perfect scores when the answer is positive, but below 5% when it is negative.

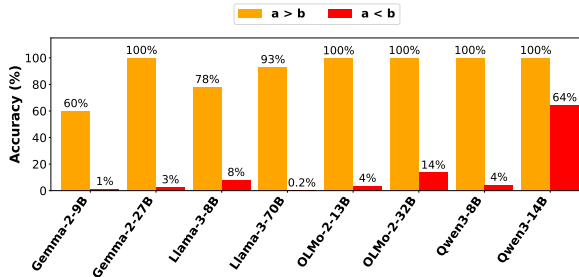


Figure 2: Zero-shot performance on subtraction ($a - b$) across operand and prompt variants. Pretrained LLMs perform subtraction well when $a > b$ but fail almost completely when $a < b$, showing a strong asymmetry.

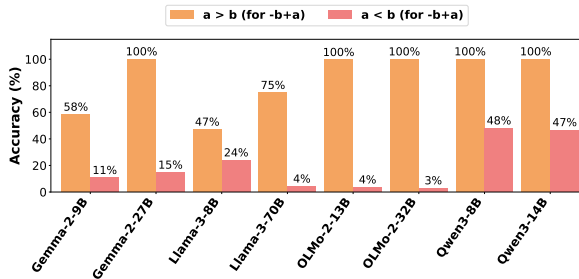


Figure 3: Zero-shot performance on the $-b+a$ input pairs. This plot shows the same asymmetry as Figure 2.

To confirm if this asymmetry is specific to the $a-b$ template or a general difficulty with negative-

signed results, we analyze the $-b+a$ input type. We use the same original operands a and b to generate this data. We find an identical trend for $-b+a$ as shown in Figure 3. This confirms that the models' failure is not due to the subtraction operation itself, but rather a systemic difficulty in producing a negative integer as the final answer.

3.3 What errors do LLMs make when $a < b$?

To isolate the source of errors on negative-result tasks, we measure the numerical accuracy without the negative sign. Figure 4 shows that accuracy jumps significantly under this relaxed metric for both the $a-b$ and $-b+a$ input pairs for all the LLMs.

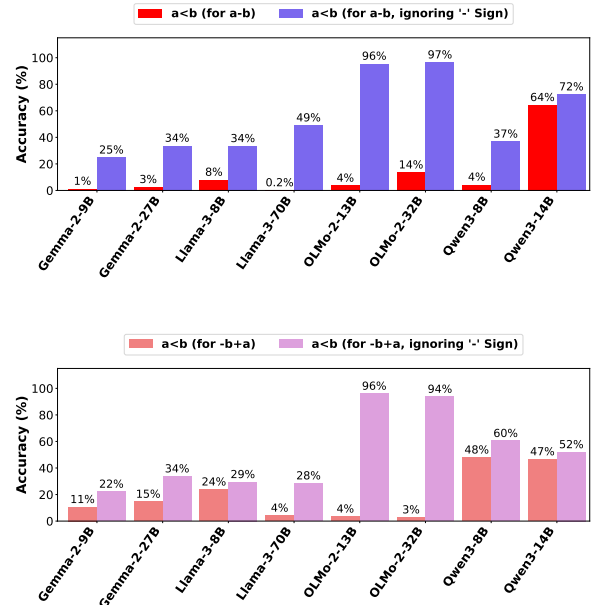


Figure 4: Zero-shot performance on $a-b$ (above) and $-b+a$ (below). The gap between the accuracy with and without the '-' shows that pretrained LLMs often compute the correct magnitude but omit the negative sign.

This common error pattern indicates that LLMs

may compute the correct magnitude but systematically omit the negative sign.

3.4 Did LLMs simply forget to include the sign?

The increase in accuracy when ignoring the negative sign (Figure 4) raises the question of whether LLMs internally represent the sign even if they fail to output it. To test this, we trained linear probes on the final layer activations of three representative LLMs: Gemma-2 9B, Llama-3.1-8B, and Qwen3-8B. The probing task was deliberately simple: predict whether the result of a subtraction should be positive or negative.

The probes achieved near-perfect accuracy across all settings. For instance, 100% on Gemma-2 9B and Qwen3-8B, and above 99% on Llama-3.1-8B. These numbers are averaged over five independent runs, with standard deviations below 0.1, confirming the stability of the results.

Consequently, while LLMs often omit the negative sign in their textual outputs, their hidden states consistently distinguish between positive and negative outcomes. This indicates a disconnect between representation and generation: LLMs reflect internally that there is a difference between positive and negative results, but this knowledge is not faithfully transferred to the decoding stage.

For all the results shown in Section 3.1, 3.2, 3.3, and 3.4, we observe the same trend for the multi-token input numbers. We provide the results for the same in Appendix A.

4 Can we make LLMs better for $a < b$?

4.1 In-context learning with pretrained LLMs

Literature suggests that LLMs can benefit from in-context examples (Brown et al., 2020; Agarwal et al., 2024). Therefore, we evaluate pretrained LLMs with 3, 5, and 10-shot in-context examples to assess whether this improves subtraction performance for the case of $a < b$.

As shown in Table 2, few-shot prompting leads to mixed results. Some LLMs exhibit moderate gains, but the improvements are neither uniform nor consistent across families or sizes. For instance, Llama-3.1-8B improves substantially from 8.1% to 31.5% accuracy (with the negative sign) when moving from zero- to five-shot prompting, suggesting limited but non-trivial in-context learning. Qwen3-14B also benefits initially, though its performance plateaus beyond three shots. Con-

versely, larger LLMs such as Llama-3.1-70B and Gemma-2-27B remain unstable, showing marginal or inconsistent gains with additional examples.

Across nearly all LLMs, accuracy is considerably higher when the negative sign is ignored, often exceeding 90%. This indicates that LLMs frequently compute the correct magnitude but omit the sign, consistent with the analysis in Section 4.

In summary, while in-context examples can occasionally help pretrained LLMs handle negative results, the overall effect is modest and inconsistent. Complete results for all shot settings (3, 5, and 10) are provided in Appendix B.

LLM	0-shot		5-shot	
	w (-)	w/o (-)	w (-)	w/o (-)
Gemma-2-9B	1.33	24.89	0.00	95.56
Gemma-2-27B	2.67	33.78	0.00	92.00
Llama-3-8B	8.13	33.51	31.46	87.72
Llama-3-70B	0.22	49.11	1.52	90.71
OLMo-2-13B	3.70	95.55	18.97	99.09
OLMo-2-32B	13.65	96.55	5.09	94.51
Qwen3-8B	4.44	36.89	4.89	96.44
Qwen3-14B	64.44	72.44	55.11	90.67

Table 2: Few-shot prompting accuracy (%) of pretrained LLMs with (-) sign and without the (-) sign.

4.2 Instruction-tuned LLMs

The instruction-tuned LLMs report strong performance on various math benchmarks such as MATH and GSM8k. Consequently, we test the instruction-tuned variants of our pretrained LLMs on this subtraction subtask. Our results suggest that almost all instruction-tuned LLMs reach above 90% accuracy (Table 3), including LLMs that fail completely in their pretrained variants.

We speculate that these gains come from the instruction fine-tuning stage. Although it is not possible to investigate the instruction-tuning data for all the LLMs, it is possible for OLMo-2. We investigate OLMo-2 instruction-tuning dataset and find that it contains MATH (Hendrycks et al., 2021) training set, GSM8k (Cobbe et al., 2021) training set and Tülu 3 (Lambert et al., 2024) dataset. All these datasets contain subtraction data including input cases where $a < b$. For the same, we speculate that OLMo-2 variants benefit from this additional data seen during instruction fine-tuning (for exact fine-tuning prompts refer to Appendix C).

LLM	Pretrained	Instruction-Tuned
Gemma-2-9B	1.33	100.00
Gemma-2-27B	2.67	100.00
Llama-3-8B	8.13	91.42
Llama-3-70B	0.22	99.91
OLMo-2-13B	3.70	99.54
OLMo-2-32B	13.65	88.27
Qwen3-8B	4.44	100.00
Qwen3-14B	64.44	100.00

Table 3: Zero-shot accuracy (%) of pretrained and instruction-tuned LLMs. Instruction tuning substantially improves performance on arithmetic queries, often pushing accuracy above 90%.

5 Conclusion

Taken together, our results highlight subtraction as a persistent challenge for current LLMs. Even LLMs that solve addition near perfectly often falter when subtraction yields a negative result, typically omitting the sign. Probing experiments suggest that the relevant information is present in hidden states but does not consistently transfer to the output layer, underscoring a mismatch between representation and generation. Few-shot prompting reduces errors for most pretrained LLMs but remains unstable, while instruction-tuned LLMs consistently perform better in most cases and achieve near-perfect performance. Thus, subtraction can provide a valuable diagnostic lens on how LLMs handle basic arithmetic. For the same, subtraction should be given equal importance as a standard testbed for evaluating numerical reasoning in future LLM research.

Limitations

While we test the LLMs on a reasonably large test, we only ran inference using three seeds. We also could not test any closed-source LLMs since their pretrained versions are not available via APIs.

Acknowledgement

This research is funded in part by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) – Project-ID 471607914 – GRK 2853/1 “Neuroexplicit Models of Language, Vision, and Action”; and Project-ID 389792660 – TRR 248 “Foundations of Perspicuous Software Systems”. The authors gratefully acknowledge the computing time granted by the John von Neumann Institute for Computing (NIC) and provided on the

supercomputer JURECA at Jülich Supercomputing Centre (JSC). The authors gratefully acknowledge the computing time made available to them on the high-performance computer at the NHR Center of TU Dresden. This center is jointly supported by the Federal Ministry of Research, Technology and Space of Germany and the state governments participating in the NHR (www.nhr-verein.de/unsere-partner). The authors would like to thank Alisa Kovtunova, Aleksandra Bakalova, Iza Škrjanec, Sukrut Rao, Yash Sarrof for discussions and feedback on the draft.

References

- Rishabh Agarwal, Avi Singh, Lei M Zhang, Bernd Bohnet, Luis Rosias, Stephanie C.Y. Chan, Biao Zhang, Aleksandra Faust, and Hugo Larochelle. 2024. [Many-shot in-context learning](#). In *ICML 2024 Workshop on In-Context Learning*.
- Mislav Balunović, Jasper Dekoninck, Nikola Jovanović, Ivo Petrov, and Martin T. Vechev. 2025. Mathconstruct: Challenging llm reasoning with constructive proofs. *CoRR*, abs/2502.10197.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D 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 Ziegler, Jeffrey Wu, Clemens Winter, and 12 others. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems. *CoRR*, abs/2110.14168.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, and 1 others. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- 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. In *NeurIPS Datasets and Benchmarks Track*.
- Shima Imani, Liang Du, and Harsh Shrivastava. 2023. [MathPrompter: Mathematical reasoning using large language models](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational*

Linguistics (Volume 5: Industry Track), pages 37–42, Toronto, Canada. Association for Computational Linguistics.

Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. *Advances in neural information processing systems*, 35:22199–22213.

Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with PagedAttention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*.

Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester James V Miranda, Alisa Liu, Nouha Dziri, Shane Lyu, and 1 others. 2024. Tulu 3: Pushing frontiers in open language model post-training. *arXiv preprint arXiv:2411.15124*.

Haoyang Li, Xuejia Chen, Zhanchao Xu, Darian Li, Nicole Hu, Fei Teng, Yiming Li, Luyu Qiu, Chen Jason Zhang, Li Qing, and Lei Chen. 2025. [Exposing numeracy gaps: A benchmark to evaluate fundamental numerical abilities in large language models](#). In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 20004–20026, Vienna, Austria. Association for Computational Linguistics.

Sean Michael McLeish, Arpit Bansal, Alex Stein, Neel Jain, John Kirchenbauer, Brian R. Bartoldson, Bhavya Kailkhura, Abhinav Bhatele, Jonas Geiping, Avi Schwarzschild, and Tom Goldstein. 2024. [Transformers can do arithmetic with the right embeddings](#). In *The 4th Workshop on Mathematical Reasoning and AI at NeurIPS’24*.

Yaniv Nikankin, Anja Reusch, Aaron Mueller, and Yonatan Belinkov. 2025. [Arithmetic without algorithms: Language models solve math with a bag of heuristics](#). In *The Thirteenth International Conference on Learning Representations*.

Team OLMo, Pete Walsh, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Shane Arora, Akshita Bhagia, Yuling Gu, Shengyi Huang, Matt Jordan, and 1 others. 2024. 2 olmo 2 furious. *arXiv preprint arXiv:2501.00656*.

Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, and 1 others. 2024. Gemma 2: Improving open language models at a practical size. *arXiv preprint arXiv:2408.00118*.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, and 1 others. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837.

An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, and 1 others. 2025a. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*.

Haotong Yang, Yi Hu, Shijia Kang, Zhouchen Lin, and Muhan Zhang. 2025b. [Number cookbook: Number understanding of language models and how to improve it](#). In *The Thirteenth International Conference on Learning Representations*.

Tianyi Zhou, Deqing Fu, Vatsal Sharan, and Robin Jia. 2024. [Pre-trained large language models use fourier features to compute addition](#). In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.

A Multi-token performance

We run the inference on LLMs and analyze their performance using numbers that were split into multiple tokens as well. We test up to 3 tokens for a single number. For generating multi-token data, we use the same procedure described in Section 2.1. Here, the operand range is between the maximum number in the tokenizer vocab + 1 to the maximum number in the tokenizer vocab $\times 100$. We describe our findings for these experiments in the following subsections.

A.1 Can LLMs subtract multi-token numbers?

We observe an expected drop in performance for multi-token addition predicted by (Yang et al., 2025b), but the trend of lower subtraction performance across all LLMs still holds for the multi-token subset as well. We provide the results in Figure 5 below.

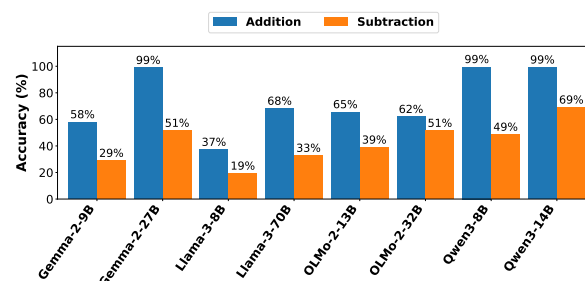


Figure 5: Zero-shot performance on multi-token addition and subtraction problems (averaged across prompt variants). Similar to single-token subtraction, multi-token subtraction is also consistently harder for all the LLMs.

A.2 Multi-token Subtraction with $a > b$ versus $a < b$

The trend of lower performance on the $a < b$ subset compared to the $a > b$ subset across all LLMs holds for the multi-token subset as well. We provide these results in Figure 6.

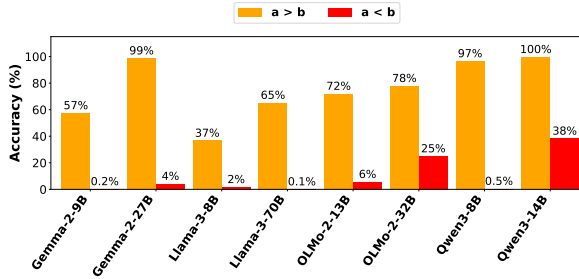


Figure 6: Zero-shot performance on Multi-token subtraction across operand and prompt variants. Pretrained LLMs do not perform multi-token subtraction as well as the single-token subtraction, but the $a > b$ vs $a < b$ performance asymmetry still holds.

A.3 What goes wrong for $a < b$ in the multi-token subset?

The trend of missing ‘-’ sign on the $a < b$ subset across all LLMs holds for the multi-token subset as well. We provide these results in Figure 7.

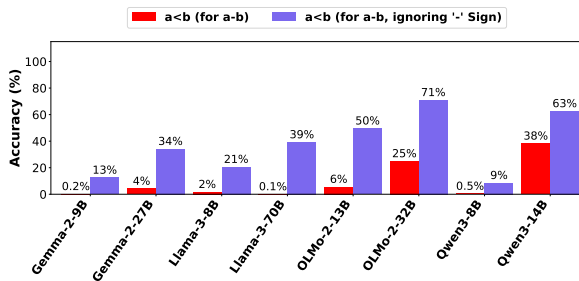


Figure 7: Zero-shot performance on $a < b$ examples for the Multi-token subset. The consistent gap between the accuracy and the accuracy ignoring sign across all the LLMs shows that pretrained LLMs often compute the correct magnitude but fail to output the negative sign.

B Few-shot additional results

C OLMo Instruction-tuning prompt

The prompt used in the Tulu dataset for fine-tuning OLMo models is given in Figure 8.

D Prompt details

Here, we provide all the prompts used in our experiments verbatim. Figure 9 contains the zero-shot variants.

E Negative sign probing

Table 5 contains all the details for each run.

LLM	0-shot		3-shot		5-shot		10-shot	
	w (-)	w/o (-)	w (-)	w/o (-)	w (-)	w/o (-)	w (-)	w/o (-)
Gemma-2-9B	1.33	24.89	0.00	92.00	0.00	95.56	0.00	99.56
Gemma-2-27B	2.67	33.78	0.44	72.89	0.00	92.00	0.00	93.33
Llama-3-8B	8.13	33.51	21.94	85.14	31.46	87.72	29.28	86.08
Llama-3-70B	0.22	49.11	1.49	71.82	1.52	90.71	0.77	97.49
OLMo-2-13B	3.70	95.55	19.92	98.63	18.97	99.09	18.27	99.01
OLMo-2-32B	13.65	96.55	5.21	93.96	5.09	94.51	3.89	82.61
Qwen3-8B	4.44	36.89	4.44	93.78	4.89	96.44	2.22	92.89
Qwen3-14B	64.44	72.44	61.78	80.44	55.11	90.67	15.11	88.89

Table 4: Few-shot prompting accuracy (%) of pretrained LLMs with (-) sign and without the (-) sign

Math Problems (response)
<p>Provide solution to the given math problem.</p> <p>Problem: {math_problem}</p> <p>Note: Provide your solution step-by-step, and end your solution in a new line in the following format:</p> <p>Final Answer: The final answer is <i>final_answer</i>.</p>

Figure 8: Prompt used for math word problems.

1. ### a {op} b = \boxed{{
2. Solve the problem below and give only the simplified answer in \boxed{answer}.
Problem: a {op} b =
Answer: \boxed{
3. Solve the problem by computing the {op}, and report only the final answer in \boxed{answer}.
Problem: a {op} b =
Answer: \boxed{
4. Find the {op} ({op_sign}) for the problem below and respond only with \boxed{answer}.
Problem: a {op} b =
Answer: \boxed{
5. Perform the {op} ({op_sign}) to solve the problem; output only \boxed{answer}.
Problem: a {op} b =
Answer: \boxed{

Figure 9: Zero-shot prompt variants used in our experiments. Here, a and b are numbers; op is either *sum* or *difference*; op_sign is $+$ or $-$. Each prompt requires the LLMs to return the answer inside a boxed expression.

LLM	Probe Type	Run 1	Run 2	Run 3	Run 4	Run 5	Mean \pm Std. Dev.
Gemma-2 9B	Single-token	100.00	100.00	100.00	100.00	100.00	100.00 \pm 0.00
	Multi-token	99.52	99.72	99.60	99.56	99.60	99.60 \pm 0.07
	Combined	99.82	99.70	99.82	99.86	99.66	99.77 \pm 0.08
Llama-3.1 8B	Single-token	99.60	99.48	99.72	99.88	99.64	99.66 \pm 0.13
	Multi-token	99.00	98.88	99.20	99.24	99.16	99.10 \pm 0.14
	Combined	99.26	99.46	99.30	99.32	99.50	99.37 \pm 0.09
Qwen3 8B	Single-token	100.00	100.00	100.00	100.00	100.00	100.00 \pm 0.00
	Multi-token	99.92	99.88	99.92	100.00	100.00	99.94 \pm 0.05
	Combined	99.83	100.00	100.00	99.93	100.00	99.95 \pm 0.07

Table 5: Detailed probe accuracy results across all LLMs and runs. All values are percentages (%).