

---

# Investigating the interaction of linguistic and mathematical reasoning in language models using multilingual number puzzles

---

**Antara Raaghavi Bhattacharya**  
Harvard University  
Cambridge, MA, USA  
antara@alumni.harvard.edu

**Isabel Papadimitriou**  
University of British Columbia  
Vancouver, Canada  
isabel.papadimitriou@ubc.ca

**Kathryn Davidson**  
Harvard University  
Cambridge, MA, USA  
kathryndavidson@fas.harvard.edu

**David Alvarez-Melis**  
Harvard University  
Cambridge, MA, USA  
dam@seas.harvard.edu

## Abstract

Across languages, numeral systems vary widely in how they construct and combine numbers. While humans consistently learn to navigate this diversity, large language models (LLMs) struggle with linguistic-mathematical puzzles involving cross-linguistic numeral systems, which humans can learn to solve successfully. We investigate *why* this task is difficult for LLMs through a series of experiments that untangle the linguistic and mathematical aspects of numbers in language. Our experiments establish that models cannot consistently solve such problems unless the mathematical operations in the problems are explicitly marked using known symbols (+, ×, etc, as in “twenty + three”). In further ablation studies, we probe how individual parameters of numeral construction and combination affect performance. While humans use their linguistic understanding of numbers to make inferences about the implicit compositional structure of numerals, LLMs seem to lack this notion of implicit numeral structure. We conclude that the ability to flexibly infer compositional rules from implicit patterns in human-scale data remains an open challenge for current reasoning models.

## 1 Introduction

Language models reason and solve problems using language. What is the connection (and the integration) between their linguistic systems and their impressive reasoning abilities? To investigate this question, we run a suite of experiments to analyze how language models solve puzzles about diverse linguistic number systems. People represent numbers through language, using rule-based systems that are simultaneously linguistic and mathematical [Ifrah, 2000, Dehaene, 2011, Carey, 2004, Le Corre and Carey, 2007, Ionin and Matushansky, 2006, Hammarström, 2010, Comrie, 2011]. Unlike most mathematical reasoning problems, where the mathematical operators are explicit, a numeral system contains implicit operations for describing numerals, and there is considerable variety in how this is done across the world’s languages. For example, French *vingt-neuf* (20 + 9), Bengali *untirish* (30 − 1), Tamil *irupatti onpatu* ( $2 \times 10 + (10 - 1)$ ), and Birom *bākūrū bībā nā vè tūjūn* ( $2 \times 12 + 5$ ) all evaluate to the Hindu-Arabic numeral 29.

We investigate the capabilities of language models to solve puzzles about linguistic number systems, drawn from linguistics competitions (Linguistics Olympiads) where high-school students have to

Variable	Operator			Example
	Explicitness	Familiarity	Type	
Single character	Implicit	-	-	A B
	Explicit	Familiar	Symbol	A + B
	Explicit	Unfamiliar	Symbol	A $\alpha$ B
	Explicit	Unfamiliar	Word	A <i>xebrut</i> B
Multi-character	Implicit	-	-	gbaifi pagig
	Explicit	Familiar	Symbol	gbaifi + pagig
...	...	...	...	

Table 1: A demonstration of the experimental conditions for our explicit operators experiment. We add explicit operators to our base IMPLICIT problems, using both familiar symbols for addition/multiplication/subtraction and unfamiliar symbols and words to symbolize the operation.

reason through data about unknown languages and explain the linguistic rules governing the data [Derzhanski and Payne, 2010]. While language models approach human performance on several language-based benchmarks [Hendrycks et al., 2020, Kojima et al., 2022, Beguš et al., 2023], and recent reasoning models deliberately optimized for logical and mathematical reasoning show remarkable performance improvements for many structured mathematical reasoning tasks [Zhong et al., 2024, Jaech et al., 2024], LLMs perform extremely poorly at solving linguistic-mathematical puzzles about systems of numbers in different languages [Derzhanski and Veneva, 2018, Bean et al., 2024].

**Why do language models fail to solve these problems at the intersection of language and math** — what specifically causes this failure? And how much of this failure is due to the linguistic vs. the mathematical aspects of the problem?

We present a method to systematically isolate individual parameters of number construction and combination and investigate how they affect language model performance. We establish that most individual mathematical features (like base) do not hinder the ability of sufficiently advanced language models to solve such problems. However, unless **the mathematical operations in a problem are made explicit through familiar symbols (+,  $\times$ , etc.)**, models cannot consistently solve the problem. This indicates that, at least within the domain of linguistic-mathematical problems, models cannot infer the compositional structure of numerals like humans can, or sufficiently abstract notions like operators. We discuss our findings in the broader context of human language, concluding that flexible, adaptive use of language across domains appears to remain challenging for LLMs.

## 2 Background

### 2.1 Linguistic and cognitive connections

People acquire systems of number representation as part of learning language, and are consequently able to construct arbitrary numerals using the rules that they learn. Although the system of rules may be language-specific, the general framework of numeral construction and combination is a fundamental cognitive ability [Hurford, 1987, Feigenson et al., 2004]. Performing mathematics in a symbolic sense requires explicit instruction (e.g. a child would not inherently know what + connotes), but once this symbolic meaning has been learned, people can generalize it to apply to any numbers [Sarnecka et al., 2015].

Numeral operations in language can be marked both explicitly (e.g. *und* in German *einundzwanzig*) and implicitly (as in English *twenty-one*), with larger numerals often using a combination of implicit and explicit operations (*five hundred and one* =  $5 \times 100 + 1$ ). Even when operations are implicit, people can understand and infer the cross-linguistic compositional structure of numerals [Ionin and Matushansky, 2006]. In a linguistics contest, a high-school student would not need to know any mathematical concepts beyond basic arithmetic to reason through number system problems and infer the rules needed to solve them. The challenge lies instead in whether models can learn and infer such rules from limited data — a characteristic capacity of humans acquiring language.

## 2.2 Mathematical ability in language models

Recent language models seem to display strong numerical understanding and processing abilities if presented with purely mathematical problems in standard formats [Yang et al., 2024], particularly for small numbers and simple mathematical operations (of the kinds used in linguistics contest problems). Current reasoning models appear to perform well at arithmetic and algebra, math word problems [Ahn et al., 2024], and difficult mathematical contest questions equivalent to advanced college-level math problems [Fang et al., 2024, Chervonyi et al., 2025], although their problem-solving ability is sometimes inconsistent [McCoy et al., 2023, Shojaee et al., 2025]. If such models are unable to solve linguistic-mathematical problems involving much simpler mathematics, and introducing linguistic structure into the problem causes their reasoning ability to break down, this indicates limitations in the *scope* of their reasoning — models may be unable to apply their reasoning flexibly across domains in the ways that humans do.

## 3 Methods

**Models.** We used OpenAI o1-mini [Jaech et al., 2024] and DeepSeek-R1-distill-Qwen-7B [Guo et al., 2025] reasoning models to conduct our experiments, querying o1-mini via the API and running DeepSeek locally. All code and data used for our experiments are available at <https://github.com/antara-raaghavi/multilingual-number-puzzles>.

We additionally queried an instruction-tuned model (qwen-2-7b-vl-instruct) and a base model (llama-3.1-8B), both of which had an accuracy of 0 across all conditions that we test. These models almost always generated longer text answers without numbers rather than the simple numerical answer required, and were hence excluded from our analyses.

**Data.** We obtained data for linguistics olympiad problems from two publicly available datasets: LingOly [Bean et al., 2024] and Linguini [Sánchez et al., 2024], filtering both datasets for problems tagged as “number systems”. After filtering, we had 15 problems from the LingOly and 8 problems from the Linguini dataset. Not every problem in the dataset could be standardized in the ways that our experiments required. The entire dataset was thus manually evaluated for suitable problems, and 10 problems were chosen for evaluation, all in distinct languages (see Appendix G). These problems spanned a range of difficulty from the first round of the UK Linguistics Olympiad to the International Linguistics Olympiad (most challenging).

## 4 Experiments

### 4.1 The effect of explicit operators in problems

Since so many of the mathematical operators in numeral structure are implicit (eg, in English we say ‘twenty three’ to mean ‘twenty + three’), our first experiment investigates how this implicit structure affects how models solve the problems. To do this, we standardize and convert the 10 existing linguistic number system problems to mathematical problems, and vary how explicit the operators are, as shown in Table 1.

First, we standardize all problems to control for model tokenization and task-external knowledge effects: we identify all meaningful morphemes, standardize all phonological changes, and replace them with dummy words as described in detail in Appendix B. This standardized version of each problem is what we call the IMPLICIT setting, since the mathematical operations are largely implicit, as they are in language. Taking these IMPLICIT problems as our baselines, we then make the operators explicit in three ways: 1) as the familiar mathematical operator symbols that perform the operation (eg, ‘+’ for addition), 2) as symbols that are unfamiliar for performing that operation, and 3) as whole words sampled from the tokenizer. A full example prompt with a puzzle in four variations is provided in Appendix C.

We present our results in Figure 1. In all cases, the presence of explicit operations with familiar symbols yields significant improvements over the default IMPLICIT condition (o1-mini performs at ceiling). In the multi-character setting (more linguistic), models perform better on average in the IMPLICIT condition than in the case with an explicit operator as an unfamiliar random word (vid. Figure 4). It is likely harder to differentiate between function words (operators) and number

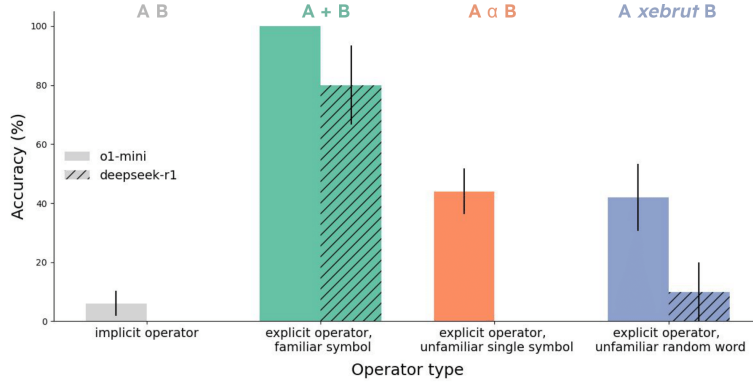


Figure 1: **Making operators explicit significantly improves performance.** Results for explicit operator experiments, for the single-character variable case (see Appendix C Figure 4 for multi-character variables). Making operators explicit shows performance improvement over the IMPLICIT condition, but this is only substantially and reliably the case when the operator is made explicit with a familiar symbol like “+”. Error bars = standard error of the mean. 10 problems, 5 iterations / problem. A detailed error analysis is provided in Appendix D.

words (numerals) in this setting — a finding consistent with work that shows human solvers also find a problem to be more difficult when the operator word is explicit but unfamiliar [Derzhanski and Veneva, 2018]. Overall, our results demonstrate that it is difficult for models to reason about the abstract idea that linguistic quantities might contain *operators*, if the operators are not explicitly provided using familiar symbols.

## 4.2 Providing contextual information

Our first experiment showed that in the absence of problem-specific instructions, when given a linguistic-mathematical problem directly, LLMs struggle to solve it unless the operations are both explicit and familiar. This leaves open the question of whether providing additional problem-specific information would affect the model performance. We thus modulate the context of the problem in three different ways. We query the same four problem variants as described in Table 1, additionally providing the following contextual information:

**Language:** “Here is a puzzle *based on numbers in the {language} language*.”

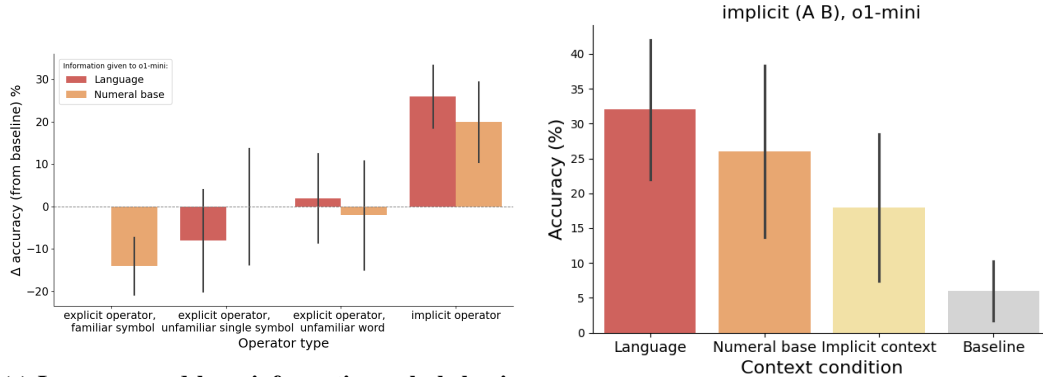
**Base:** “Here is a puzzle *based on numbers in a language that uses a base- $\{n\}$  numeral system*.”

**Implicit operations:** “Here is a puzzle *based on numbers in a language. In this language, numbers may be constructed through implicit operations like addition (twenty-nine =  $20 + 9$ ) or multiplication (five hundred =  $5 \times 100$ )*.” [only for IMPLICIT condition]

We compare these to the baseline results from Section 4.1 for o1-mini, presenting our results in Figure 2a. In cases *other* than the implicit operator condition, the model seems to recognize the problem as requiring a more mathematical kind of reasoning, so providing linguistic information seems to confuse the model and average performance is worse. However, in the implicit operator (A B) condition, model performance improves significantly, perhaps because the setting of the problem is less overtly mathematical. In Figure 2b, we show that providing information about the implicit reasoning needed is not as significant a boost as activating knowledge about the specific language.

## 4.3 Ablations: constructed minimal-pair problems

In order to ensure that it is the difference in operators (as opposed to other features of the numeral system) that explains the models’ inability to solve these problems, we performed an ablation study to test whether models could handle other aspects of numeral construction and combination. Our experiment is inspired by the notion of a linguistic *minimal pair*, a pair of linguistic items that differ in exactly one meaningful element. We construct minimal pairs of simple, synthetic number system problems, where every element is the same except for one specific parameter that differs between two paired problems.



(a) **Language and base information only helps in the IMPLICIT case.** Effect of adding language or numeral base information, plotted as a difference from the baseline values in Figure 1 for o1-mini. In cases with explicit operators, conflating overtly mathematical and linguistic information appears to confuse the models.

(b) **Extra information improves performance on IMPLICIT problems (A B).** Information about implicitness is helpful, but not as much as more direct information like the problem language. Error bars = standard error of the mean. 5 iterations / problem.

We tested five major parameters of numeral systems: numeral representation (symbolic numeral glyphs vs. numeral words), ordering (right-to-left vs. left-to-right), and combination (additive vs. subtractive). As the degenerate case, we compare whether the system is a numeral system or not (i.e. is just a regular linguistic system). We detail our specific test setup and results in Appendix E.

In all cases, GPT-4 and o1-mini could solve the template problems. It thus appears that most basic “building blocks” of number systems (e.g. the base of the system, the order of numerals, etc.) did not affect model performance in isolation, but the models consistently fail to solve number problems that involve constructing and combining complex numerals.

## 5 Discussion and Conclusions

We study the entanglement between linguistic and numeric knowledge in language models, focusing on the ability of models to use mathematical reasoning in problems that display the implicit numerical structure in language. In the setting of these linguistic-mathematical puzzles, we show that the overtness and familiarity of operators affects the performance of language models, although many humans are able to understand how numeral systems work and hence solve the problems without needing specified operators. However, a broader study with different controls and parameter settings remains open for future work. Since all our evaluation was standardized and closed-form, we welcome research on open-ended evaluation of reasoning task responses. Current language models seem to display some level of emergent modular structure [Teehan et al., 2022, Lepori et al., 2023] — perhaps linguistic and mathematical tasks activate separate circuits or subspaces in models, and understanding the ways in which reasoning fine-tuning and reinforcement learning interacts with linguistic pretraining is another promising avenue for future research. Investigating such questions enriches our understanding of both computational and human approaches to representing numbers in language. The ability to understand language and abstract rule-governed systems is a fascinating aspect of human intelligence, and we hope that our research provides some insight into the understanding of this remarkable human trait.

## Acknowledgments and Disclosure of Funding

The authors gratefully thank Tom McCoy and Kaden Holladay for helpful discussions in the initial stages of this project. This work has been made possible in part by a gift from the Chan Zuckerberg Initiative Foundation to establish the Kempner Institute for the Study of Natural and Artificial Intelligence at Harvard University.

## References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. GPT-4 Technical Report. *arXiv preprint arXiv:2303.08774*, 2023.
- Janice Ahn, Rishu Verma, Renze Lou, Di Liu, Rui Zhang, and Wenpeng Yin. Large Language Models for Mathematical Reasoning: Progresses and Challenges. *arXiv preprint arXiv:2402.00157*, 2024.
- Andrew M Bean, Simi Hellsten, Harry Mayne, Jabez Magomere, Ethan A Chi, Ryan Chi, Scott A Hale, and Hannah Rose Kirk. Lingoly: A benchmark of olympiad-level linguistic reasoning puzzles in low-resource and extinct languages. *arXiv preprint arXiv:2406.06196*, 2024.
- Gašper Beguš, Maksymilian Dąbkowski, and Ryan Rhodes. Large linguistic models: Analyzing theoretical linguistic abilities of LLMs. *arXiv preprint arXiv:2305.00948*, 2023.
- Susan Carey. Bootstrapping & the origin of concepts. *Daedalus*, 133(1):59–68, 2004.
- Yuri Chervonyi, Trieu H Trinh, Miroslav Olšák, Xiaomeng Yang, Hoang Nguyen, Marcelo Menegali, Junehyuk Jung, Vikas Verma, Quoc V Le, and Thang Luong. Gold-medalist Performance in Solving Olympiad Geometry with AlphaGeometry2. *arXiv preprint arXiv:2502.03544*, 2025.
- Bernard Comrie. Typology of numeral systems. *Numeral types and changes worldwide. Trends in Linguistics. Studies and monographs*, 118, 2011.
- Stanislas Dehaene. *The number sense: How the mind creates mathematics*. Oxford University Press USA, 2011.
- Ivan Derzhanski and Thomas Payne. The Linguistics Olympiads: Academic competitions in linguistics for secondary school students. *Linguistics at school: language awareness in primary and secondary education*, pages 213–26, 2010.
- Ivan Derzhanski and Milena Veneva. Linguistic Problems on Number Names. In *Proceedings of the Third International Conference on Computational Linguistics in Bulgaria (CLIB 2018)*, pages 169–176, 2018.
- Meng Fang, Xiangpeng Wan, Fei Lu, Fei Xing, and Kai Zou. MathOdyssey: Benchmarking Mathematical Problem-solving Skills in Large Language Models using Odyssey Math Data. *arXiv preprint arXiv:2406.18321*, 2024.
- Lisa Feigenson, Stanislas Dehaene, and Elizabeth Spelke. Core systems of number. *Trends in cognitive sciences*, 8(7):307–314, 2004.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. DeepSeek-R1: Incentivizing reasoning capability in LLMs via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- Harald Hammarström. Rarities in numeral systems. *Rethinking universals: How rarities affect linguistic theory*, 45:11–53, 2010.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*, 2020.
- James R Hurford. *Language and number: The emergence of a cognitive system*. B. Blackwell, 1987.
- Georges Ifrah. *The Universal History of Numbers*. Harvill London, 2000.
- Tania Ionin and Ora Matushansky. The composition of complex cardinals. *Journal of semantics*, 23(4):315–360, 2006.
- Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, et al. OpenAI o1 System Card. *arXiv preprint arXiv:2412.16720*, 2024.

- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large Language Models are Zero-Shot Reasoners. *Advances in neural information processing systems*, 35:22199–22213, 2022.
- Mathieu Le Corre and Susan Carey. One, two, three, four, nothing more: An investigation of the conceptual sources of the verbal counting principles. *Cognition*, 105(2):395–438, 2007.
- Michael Lepori, Thomas Serre, and Ellie Pavlick. Break it down: Evidence for structural compositionality in neural networks. *Advances in Neural Information Processing Systems*, 36:42623–42660, 2023.
- R Thomas McCoy, Shunyu Yao, Dan Friedman, Matthew Hardy, and Thomas L Griffiths. Embers of Autoregression: Understanding Large Language Models Through the Problem They are Trained to Solve. *arXiv preprint arXiv:2309.13638*, 2023.
- Eduardo Sánchez, Belen Alastruey, Christophe Ropers, Pontus Stenetorp, Mikel Artetxe, and Marta R Costa-jussà. Linguini: A benchmark for language-agnostic linguistic reasoning. *arXiv preprint arXiv:2409.12126*, 2024.
- Barbara W Sarnecka, Meghan C Goldman, and Emily B Slusser. *How counting leads to children’s first representations of exact, large numbers*. Oxford University Press, 2015.
- Parshin Shojaei, Iman Mirzadeh, Keivan Alizadeh, Maxwell Horton, Samy Bengio, and Mehrdad Farajtabar. The illusion of thinking: Understanding the strengths and limitations of reasoning models via the lens of problem complexity. *arXiv preprint arXiv:2506.06941*, 2025.
- Ryan Teehan, Miruna Clinciu, Oleg Serikov, Eliza Szczechla, Natasha Seelam, Shachar Mirkin, and Aaron Gokaslan. Emergent Structures and Training Dynamics in Large Language Models. In Angela Fan, Suzana Ilic, Thomas Wolf, and Matthias Gallé, editors, *Proceedings of BigScience Episode #5 – Workshop on Challenges & Perspectives in Creating Large Language Models*, pages 146–159. Association for Computational Linguistics, 2022. doi: 10.18653/v1/2022.bigscience-1.11.
- Haotong Yang, Yi Hu, Shijia Kang, Zhouchen Lin, and Muhan Zhang. Number Cookbook: Number Understanding of Language Models and How to Improve It. *arXiv preprint arXiv:2411.03766*, 2024.
- Tianyang Zhong, Zhengliang Liu, Yi Pan, Yutong Zhang, Yifan Zhou, Shizhe Liang, Zihao Wu, Yanjun Lyu, Peng Shu, Xiaowei Yu, et al. Evaluation of OpenAI o1: Opportunities and challenges of AGI. *arXiv preprint arXiv:2409.18486*, 2024.

## Appendices

### A Limitations

We acknowledge the possibility that our results are explained by limitations in the training data and the small size of our dataset, as language models often equal human performance on benchmarks for which they have large quantities of similar-enough training data [Achiam et al., 2023]. Perhaps an LLM trained on a massive corpus of linguistic number system problems would be able to solve new, previously unseen number system problems. But the data today are far too limited for such an approach, and crucially, a human solver who is familiar with existing number system problems can generalize to unseen problems extremely well! Even a human solver who is unfamiliar with existing number system problems can in theory solve any problem they are provided just by logically reasoning. Importantly, we note that although this may not be true of the *average* human, when comparing the top end of humans with the top-performing current language models, it is clear that intuiting rules from human-scale data is still challenging for LLMs.

### B Randomization strategies for task-external knowledge and tokenization handling

In this section, we address the specific changes we make across linguistic number system problems to convert them into templates suitable for our dataset. In order to truly test whether the model is *solving* a problem, it should not be affected by factors external to the problem, such as flawed tokenization or the usage of memorized knowledge external to the provided task.<sup>1</sup>

In order to remediate this, in the single-letter token setting, we separated all characters by whitespaces to ensure correct tokenization. In the multi-token setting, we identified all meaningful morphemes in the problems and standardized them to remove any phonological changes, such that every morpheme had exactly one surface representation. We separated every meaningful morpheme with whitespaces, and mapped each morpheme to a randomly generated multi-token “dummy word” for each iteration of each experiment. We created each of these “dummy words” by randomly sampling short tokens ( $\text{length} \leq 3$ ) from the language models’ respective tokenizer vocabularies, and concatenating tokens together to create unfamiliar words.

For tokenizers which use schemes like byte-pair encoding, any input string will get mapped to some sequence of tokens that are present in the vocabulary, so there is no situation in which the model will see an unknown token. Since the dummy words themselves have no meaning, the model cannot directly draw on task-external linguistic information to solve the presented problems. For simplicity we restricted the random draw to those containing only romanized (Latin alphabet) characters. We also excluded tokens that contained any numeral symbols from 0-9, to ensure that the mathematical correctness of the problems was not affected.

---

<sup>1</sup>Memorized knowledge would also help a human solver, but people are much less likely to know the number systems of different (particularly low-resource) languages. Although linguistics olympiad contestants might know more number systems than the average person, there are over 7,000 human languages, so the probability of knowing a specific system is low. Moreover, since LLM training corpora scrape large portions of the internet, the breadth of their memorized knowledge far exceeds that of an average human.



## C Multi-character-variable results from Exp 1

We provide an example of our four variations of the puzzle in Table 2. To query all four variants, we used the same prompt “Here is a puzzle. Can you solve it? Please output only the answer (in place of the ??) and nothing else!”.

Explicit + familiar	
$(\text{masaad} \times \text{pagig}) + \text{masaad} + \text{opbob} = 31$	
$(\text{masaad} \times \text{pagig}) + \text{masaad} + \text{buylen} = 26$	
$\vdots$	
$(\text{ajssci} \times \text{pagig}) + (\text{ajssci} \times \text{kould}) = 50$	
$(\text{innops} \times \text{pagig}) + \text{innops} + \text{opbob} = ??$	
Implicit	
$\text{masaad pagig nge masaad opbob} = 31$	
$\text{masaad pagig nge masaad buylen} = 26$	
$\vdots$	
$\text{ajssci pagig nge ajssci kould} = 50$	
$\text{innops pagig nge innops opbob} = ??$	
Explicit + unfamiliar (Greek)	
$(\text{masaad } \beta \text{ pagig}) \alpha \text{ masaad } \alpha \text{ opbob} = 31$	
$(\text{masaad } \beta \text{ pagig}) \alpha \text{ masaad } \alpha \text{ buylen} = 26$	
$\vdots$	
$(\text{ajssci } \beta \text{ pagig}) \alpha (\text{ajssci } \beta \text{ kould}) = 50$	
$(\text{innops } \beta \text{ pagig}) \alpha \text{ innops } \alpha \text{ opbob} = ??$	
Explicit + unfamiliar (random)	
$(\text{masaad hibcat pagig}) \text{xebrut masaad xebrut opbob} = 31$	
$(\text{masaad hibcat pagig}) \text{xebrut masaad xebrut buylen} = 26$	
$\vdots$	
$(\text{ajssci hibcat pagig}) \text{xebrut } (\text{ajssci hibcat kould}) = 50$	
$(\text{innops hibcat pagig}) \text{xebrut innops xebrut opbob} = ??$	

Table 2: Example of four problem variants in the multi-character setting, corresponding to Drehu (IOL 2010) dataset problem in Figure 3.

**Problem #2 (20 points).** Given are Drehu numerals in alphabetical order and their values in ascending order:

*caatr nge caako, caatr nge caangömen, caatr nge caaqaihano,  
ekaatr nge ekengömen, kõniatr nge kõniko, kõniatr nge kõnipt,  
kõniatr nge kõniqaihano, lueatr nge lue, lueatr nge luako, lueatr nge luepi*

26, 31, 36, 42, 50, 52, 73, 75, 78, 89

(a) Determine the correct correspondences.

(b) Write in numerals:

*kõniatr nge eke + caatr nge luepi = ekaatr nge ekako  
luengömen + luako = ekeqaihano*

(c) Write out in Drehu: 21, 48, 83.

△ The Drehu language belongs to the Austronesian language family. It is spoken by approx. 10 000 people on Lifu Island to the east of New Caledonia. *c* = *ch* in *church*; *ng* = *ng* in *hang*; *ö* = French *eu* or German *ö*; *q* is a voiceless *w* (as *wh* in Scottish or Southern American *which*); *tr* ≈ English *t* in *art*, uttered with the tip of the tongue turned back.

—Ksenia Gilyarova

Figure 3: Drehu (IOL 2010) problem

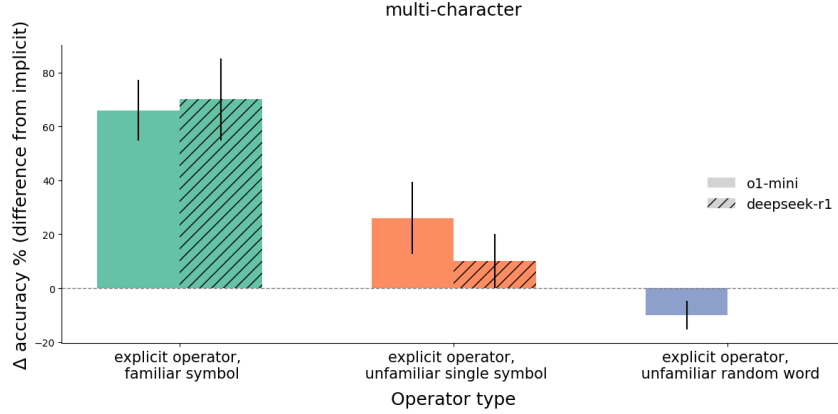


Figure 4: Both o1-mini and DeepSeek struggle with the explicit-unfamiliar condition (o1 shows negative improvement, DeepSeek shows 0%) in the multi-character setting. Error bars = standard error of the mean. 5 iterations / problem tested for 10 problems.

## D Error analysis

We observe some common patterns of error in the model responses. For the three problems which involved squares and cubes of numbers, when the operators were not explicit and familiar, o1-mini almost always responded by pattern-matching (e.g. providing another square/cube number) instead of solving the problem, as seen in Table 3. o1-mini also reproduced a number given in the input question as the answer in several cases (11 for the multi-token condition, and 3 for the single-token condition, across 150 trials) when the operators were not explicit and familiar.

Condition	Single	Multi
explicit symbol	8	4
explicit random word	0	8
implicit	14	5

Table 3: Incorrect pattern-matched square / cube answers (out of 15 possible trials)

Further, when o1-mini answered a problem incorrectly, its responses were often inconsistent across the five trials of that problem. Notably, in 50% of single-character cases lacking explicit and familiar operators, all five responses were distinct and incorrect. This further shows that performance appears to depend on the presence of explicit operator cues; in their absence, o1-mini does not reliably solve the problem.

## E Ablations: minimal pair experiment details

Parameter	GPT-3.5	GPT-4	o1-mini
Numeral system vs. not $AB = \text{fifty one} \mid AB = \text{big bird}$	✓	✓	✓
Typed vs. glyph $AB = \text{fifty one} \mid AB = 51$	✓	✓	✓
Order $L \rightarrow R$ vs. $L \leftarrow R$ $AB = 51 \mid BA = 51$	✓	✓	✓
Additive vs. subtractive $AB = 27 \mid AB = 27$ $(20 + 7) \mid (30 - 3)$	✗	✓	✓
Base of the numeral system*	✗	✓	✓

Table 4: Minimal pair results: GPT-4 and o1-mini solve all paradigms, GPT-3.5-turbo struggles with numeral base and combination. Further data on testing all bases 4-19 linked in Appendix F Table 5.

$L \rightarrow R$	$L \leftarrow R$
$A B = 51$	$B A = 51$
$A C = 57$	$C A = 57$
$D B = 41$	$B D = 41$
$D C = ??$	$C D = ??$

Figure 5: Example of full minimal pair template problem, for the Order parameter, where we varied whether digits are read left-to-right or right-to-left.

## F Base experiment

In order to understand whether sufficiently advanced language models would show performance that was invariant to changes in the base, we conducted a more fine-grained minimal pair experiment into the effect of numeral base on problem performance. Here, the solver would see the Hindu-Arabic numerals corresponding to the English base-10 representation of the numbers, because the problem was presented in English. But the unknown symbols corresponded to the numbers as expressed in a different base, as shown in Figure 6.

We conducted two different versions of this experiment. First, we mapped the unknown symbols to the single-character whitespaced  $A$ ,  $B$ ,  $C$ , and  $D$  tokens, as in Figure 6. In the second version, each of the four unknown symbols ( $A$ ,  $B$ ,  $C$ ,  $D$ ) was instead represented by a corresponding random token drawn from the tokenizer vocabulary, to ensure that the context of the specific tokens  $A$ ,  $B$ ,  $C$ , and  $D$  was not influencing our results.

We tested four increasingly sophisticated GPT models (GPT-3.5-turbo, GPT-4, GPT-4o, and o1-mini) on both versions of the experiment and provide results in Table 5. GPT-4o and o1-mini solved all problems in both conditions, displaying performance that was robust to the base of the problem.

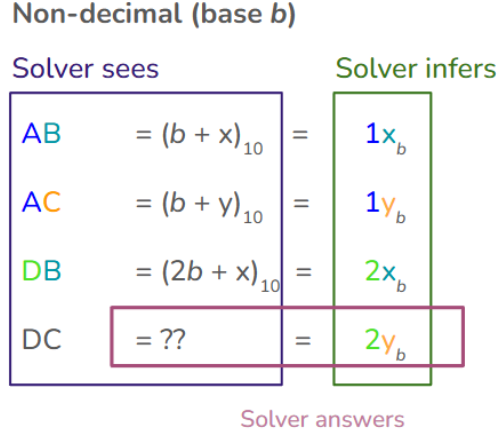


Figure 6: Setup for base experiment

Base	GPT-3.5-turbo		GPT-4		GPT-4o		o1-mini	
	ABCD	Random	ABCD	Random	ABCD	Random	ABCD	Random
4	✗	✓	✓	✓	✓	✓	✓	✓
5	✗	✗	✓	✓	✓	✓	✓	✓
6	✓	✗	✓	✓	✓	✓	✓	✓
7	✓	✓	✓	✓	✓	✓	✓	✓
8	✗	✗	✓	✓	✓	✓	✓	✓
9	✓	✓	✓	✓	✓	✓	✓	✓
11	✓	✓	✓	✓	✓	✓	✓	✓
12	✗	✓	✗	✓	✓	✓	✓	✓
13	✗	✗	✓	✓	✓	✓	✓	✓
14	✗	✗	✓	✓	✓	✓	✓	✓
15	✗	✓	✓	✓	✓	✓	✓	✓
16	✓	✗	✗	✓	✓	✓	✓	✓
17	✗	✗	✓	✓	✓	✓	✓	✓
18	✓	✓	✓	✓	✓	✓	✓	✓
19	✓	✗	✓	✓	✓	✓	✓	✓

Table 5: Base experiment results: GPT-4o / o1-mini solve every problem, regardless of randomization

## G Table of languages

Language	ISO code	Base	Level
Drehu	dhv	20	IOL
Georgian	kat	20	UKLO R1
Gumatj	gnn	5	UKLO R1
Ndom	nqm	6	IOL
Ngkolmpu	kcd	6	UKLO R1
Northern Pame	pmq	8	UKLO R1
Umbu-Ungu	ubu	24	IOL
Waorani	auc	5	UKLO R1
Yoruba	yor	20	UKLO R2
Yup'ik	esu	20	UKLO R2

Table 6: Languages and problem features in final dataset (after removing/standardizing phenomena)

We detail the 10 problems that we used for our analyses. The problems range in difficulty from the first and second rounds of the UK Linguistics Olympiad (UKLO R1 and R2) to the International Linguistics Olympiad, which typically has the most challenging problems.

## H Performance breakdown per language

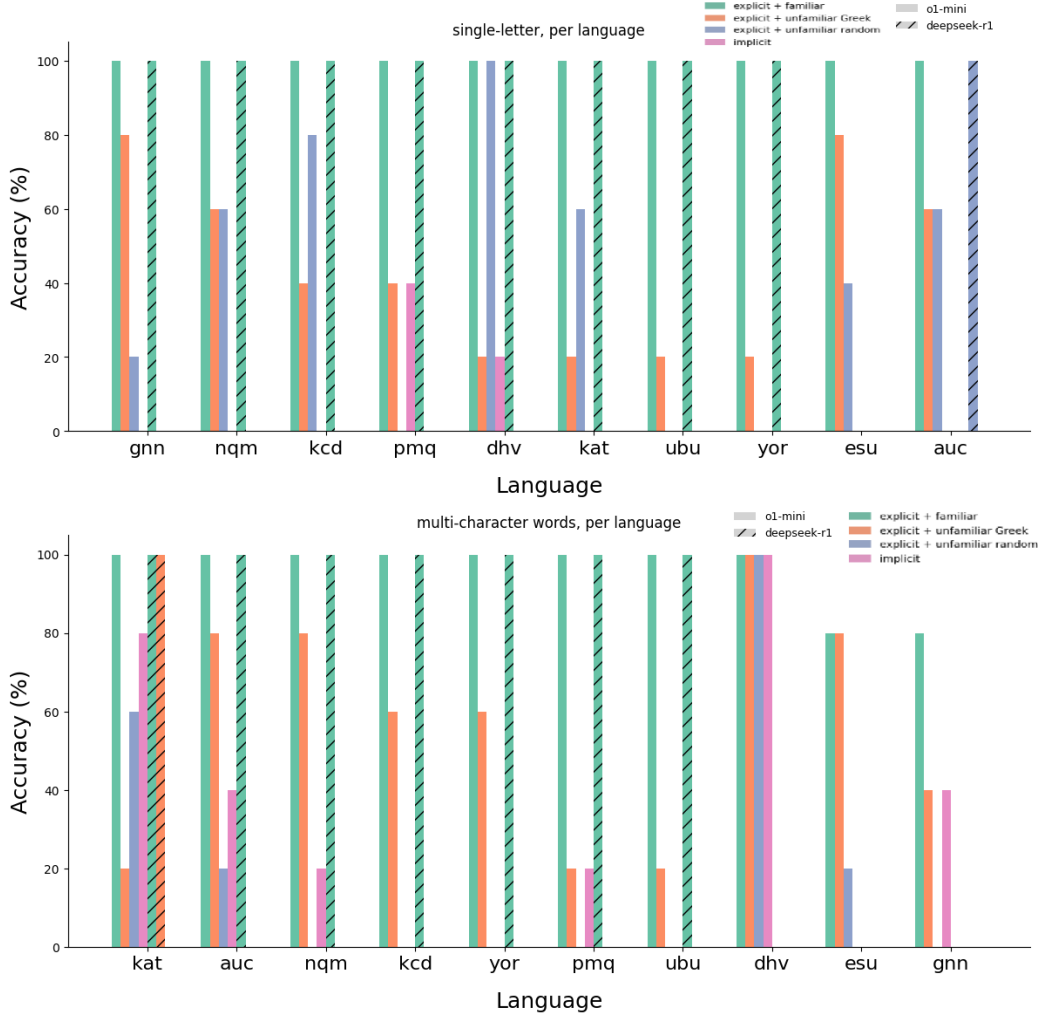


Figure 7: **Results per language, (a) single-character (b) multi-character: performance varies significantly by problem and operator type.** Note that Georgian (kat) and Drehu (dhv) are the two easiest problems in our controlled dataset, as we standardize away the phonological change and randomized numeral ordering (which human solvers find most difficult), leaving straightforward vigesimal-decimal systems like French, which models have likely had exposure to.