

A Dual-Task Paradigm to Investigate Sentence Comprehension Strategies in Language Models

Anonymous ACL submission

Abstract

Language models (LMs) behave more like humans when their cognitive resources are restricted, particularly in predicting sentence processing costs such as reading times. However, it remains unclear whether such constraints similarly affect sentence comprehension strategies, and existing methods do not directly target the balance between memory storage and sentence processing, which is central to human working memory. To address this issue, we propose a dual-task paradigm that combines an arithmetic computation task with a sentence comprehension task, such as “The 2 cocktail + blended 3 =...”. Our experiments show that under dual-task conditions, GPT-4o, o3-mini, and o4-mini shift toward plausibility-based comprehension, mirroring humans’ rational inference. Specifically, these models show a greater accuracy gap between plausible sentences (e.g., The cocktail was blended by the bartender) and implausible sentences (e.g., The bartender was blended by the cocktail) in the dual-task condition compared to the single-task conditions. These findings suggest that constraints on the balance between memory and processing resources promote rational inference in LMs. More broadly, they support the view that human-like sentence comprehension fundamentally arises from the allocation of limited cognitive resources.

1 Introduction

Working memory is a cognitive system that temporarily stores and maintains information necessary for processing in an accessible state (Atkinson and Shiffrin, 1971; Baddeley and Hitch, 1974; Baddeley, 2003). It is essential for understanding language in humans.

Comparing the working memory of LMs with that of humans helps us understand what makes sentence comprehension more human-like. Limitation of cognitive resources (analogous to working memory in humans) for sentence comprehension makes

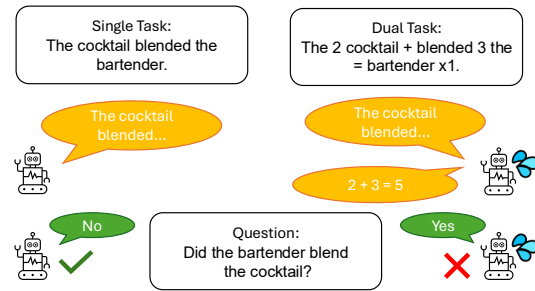


Figure 1: Overview of hypothesis and tasks. This study investigates whether language models, similar to humans, prioritize plausibility over grammar when understanding implausible sentences in a dual-task situation where they simultaneously perform calculation and sentence comprehension.

LMs behave more like humans (Futrell et al., 2020; Hahn et al., 2022; Kuribayashi et al., 2021, 2022; Oh et al., 2022; Oh and Schuler, 2023; Timkey and Linzen, 2023; Wilcox et al., 2025). Specifically, LMs with restricted memory resources better approximate human reading costs, such as reading times. These findings suggest that resource constraints may be a fundamental property of human-like language comprehension.

However, it remains unclear whether the reading strategies of LMs exhibit patterns analogous to those of humans under limited cognitive resources. Previous studies have shown that LMs achieve lower accuracy on complex sentences, suggesting that LMs rely on working-memory-like mechanisms during sentence processing (Amouyal et al., 2025a,b; Irwin et al., 2023). We therefore examine whether LMs adopt human-like comprehension strategies when cognitive resources are constrained.

Existing approaches to constraining cognitive resources in LMs also face methodological limitations. Prior work has mainly manipulated input length or model parameters (Asami and Sugawara,

067	2024; Kuribayashi et al., 2021, 2022; Oh et al.,	to misunderstand implausible sentences with	117
068	2022; Oh and Schuler, 2023; Timkey and Linzen,	passive, dative, or benefactive structures. It	118
069	2023; Wilcox et al., 2025). These methods either	suggests that they rely more on world knowl-	119
070	fail to capture the balance between storage and pro-	edge and superficial word order than on func-	120
071	cessing that characterizes human working memory	tion words under limited cognitive resources.	121
072	or fail to induce a shift in reading strategy within a		
073	single LM.		
074	To address the methodological issue, we propose	(iv) Our findings support the view that resource	122
075	a dual-task paradigm in which models simultane-	constraints are a fundamental property of	123
076	ously solve arithmetic problems and answer com-	human-like language comprehension, extend-	124
077	prehension questions, such as “The 2 cocktail +	ing previous evidence from reading cost to	125
078	blended 3 =...” (see Figure 1). We compare com-	reading strategies.	126
079	prehension accuracy across three conditions: (i)		
080	single task (comprehension of sentences without	2 Related Work	127
081	calculation, (ii) noisy single task (comprehension	2.1 Approach to Constrain the LMs’ Working	128
082	of sentences with embedded calculation, but no	Memory	129
083	concurrent arithmetic solving), and (iii) dual task	There are two major approaches in computational	130
084	(comprehension of sentences with embedded cal-	psycholinguistics for constraining the cognitive re-	131
085	ulation while solving the problems).	sources of LMs: manipulating the input text and	132
086	We focus on one characteristic property of hu-	altering the model’s parameters.	133
087	man reading strategies: limited cognitive resources	For the first approach, Asami and Sugawara	134
088	promote rational inference (Futrell and Gibson,	(2024) manipulates the length of entire sentences	135
089	2017; Gibson et al., 2013, 2016). Rational infer-	and compares the accuracy differences between	136
090	ence involves interpreting literally implausible sen-	plausible and implausible sentences. Their results	137
091	tences (e.g., “the cocktail blended the bartender,”) as	show that longer sentences reduced accuracy for	138
092	plausible meanings consistent with world knowl-	both types, indicating that increased length does	139
093	edge (e.g., “the bartender blended the cocktail.”),	not promote greater reliance on plausibility infor-	140
094	prioritizing plausibility over grammatical structure.	mation. This may be because such manipulation	141
095	Humans are more likely to adopt this strategy under	merely increases processing demands without en-	142
096	high cognitive load, such as during dual-task con-	gaging the balance between storage and processing,	143
097	ditions, than under low load (Ayasse et al., 2021;	which is a central feature of human working mem-	144
098	Rogalsky et al., 2008).	ory. A key function of working memory is its dual	145
099	Taken together, we examine whether LMs adopt	role of storage and processing, such as remember-	146
100	rational strategies under dual-task conditions and,	ing a series of numbers while performing a distract-	147
101	if so, which conditions promote such strategies.	ing task (Atkinson and Shiffrin, 1971; Baddeley	148
102	Our contributions are summarized as follows:	and Hitch, 1974; Baddeley, 2003).	149
103	(i) We propose a dual-task paradigm to test the	Regarding the second approach, studies have	150
104	behavior of LMs under limited cognitive re-	shown that reducing the number of attention heads	151
105	sources. This paradigm allows us to observe	(Timkey and Linzen, 2023), larger perplexity	152
106	the LMs’ function in integrating memory	(Kuribayashi et al., 2021; Oh et al., 2022; Oh	153
107	with sentence processing, analogous to hu-	and Schuler, 2023), limiting context access (Kurib-	154
108	man working memory.	ayashi et al., 2022), and small training data size and	155
109	(ii) We demonstrate that GPT-4o, o3-mini, and	training steps (Wilcox et al., 2025) lead to better	156
110	o4-mini shift their comprehension strategies	prediction of human reading times. However, these	157
111	toward rational inference under the dual-task	manipulations effectively create different model	158
112	condition, showing a larger accuracy gap be-	configurations, each trained for specific tasks, and	159
113	tween plausible and implausible sentences	therefore reflect differences between models rather	160
114	than in the single-task and noisy single-task	than shifts within a single model. This is analo-	161
115	conditions.	gous to comparing human participants with dif-	162
116	(iii) We show that these models are more likely	ferent working memory capacities, rather than ex-	163
		amining how one participant adapts under varying	164
		conditions.	165

To address these limitations, we introduce a dual-task paradigm in which arithmetic expressions are interleaved with sentence words (see Figure 1). This design maintains the need for memory storage while imposing additional processing demands, thereby constraining the working-memory function that balances storage and processing. Furthermore, by manipulating the task rather than the model, our approach sheds light on how a single model shifts its reading strategy under resource constraints.

2.2 Dual-Task Approaches in LMs and Humans

Previous work on multi-task processing in language models has primarily aimed to improve performance or efficiency by enabling models to handle multiple tasks simultaneously (Cheng et al., 2023; Son et al., 2024). These studies focus on optimizing accuracy or speed and are not designed to investigate how cognitive resource limitations affect language comprehension.

In contrast, some studies have attempted to constrain LMs’ working memory using n-back tasks (Kirchner, 1958), with the explicit goal of taxing internal memory resources (Gong et al., 2024; Zhang et al., 2024). While these studies share our objective of probing working memory limitations, n-back tasks primarily engage numerical memory and calculation rather than sentence comprehension.

We build on this latter approach by focusing on working memory in language comprehension. Specifically, our dual-task approach is inspired by human working-memory paradigms, particularly the operation span task (Turner and Engle, 1989). This task requires participants to perform arithmetic operations while memorizing words or letters. In the original version (Turner and Engle, 1989), a mathematical problem followed by a to-be-remembered word is presented, such as “ $(3 \times 4) + 11 = 20?$ BEAR”. Participants first read the problem aloud and judge whether the answer is correct, then read and memorize the following word. After several trials (typically two to six), they are asked to recall the memorized words in the correct order.

3 Methods

3.1 Task

We conduct three types of question-answering tasks. The Dual Task is designed according to the LM’s specifications (where a list of strings is more suitable). The Noisy Single Task is included

to examine whether performance changes are due to the presence of noisy arithmetic expressions or to the additional cognitive demands of the Dual Task. Exact prompts are in Appendix B.

- (i) **Single Task (Single)**: The LM receives a sentence without any embedded arithmetic problems, then answers a comprehension question about the sentence.
- (ii) **Noisy Single Task (Noisy)**: The LM receives a sentence with embedded arithmetic problems but ignores them, then answers a comprehension question about the sentence.
- (iii) **Dual Task (Dual)**: The LM receives a sentence with embedded arithmetic problems, solves the arithmetic problems, and then answers a comprehension question about the sentence.

Following prior work in psycholinguistics, we evaluate language models against well-established human phenomena (Ayasse et al., 2021; Futrell and Gibson, 2017; Gibson et al., 2013, 2016; Rogalsky et al., 2008) without collecting new human data. The validity of this approach lies in pattern-level correspondence: our goal is not to re-estimate human behaviour, but to determine whether LMs exhibit patterns similar to those observed in human behaviour.

3.2 Dataset

We use a subset of stimuli from the GELP dataset (Asami and Sugawara, 2024). The dataset consists of sentence–question pairs. Each sentence includes one premise connected with two propositions.¹

Premises are manipulated by plausibility (**Plausible** / **Implausible**) and construction (**Transitive** / **Passive** / **Dative** / Experiencer Subject (**Exp.Subj.**) / Experiencer Object (**Exp.Obj.**) / Benefactive For (**Ben.For**)), as illustrated in Table 1. Although the original dataset includes eight constructions, two are excluded during preprocessing (see Section 4.2.1 for details).

In addition, to conduct the Dual Task and Noisy Single Task, we add arithmetic expressions to these sentences. Randomly generated computation problems are inserted with identifiers (“x1”, “x2”, “x3”, ...). After the “=”, the corresponding identifier

¹GELP also contains sentences with one or no propositions. We use only sentences with two propositions, corresponding to the high memory-load condition.

Factor	Variable	Premise and stimuli (Implausible except for the top)
Plausibility	Plausible	The bartender blended the cocktail. (Premise)
	Implausible	The cocktail blended the bartender. (Premise)
Construction	Transitive	The cocktail blended the bartender. (Premise)
	Passive	The bartender was blended by the cocktail. (Premise)
	Dative	The chef sent the friend to the gift. (Premise)
	Exp.Subj.	The view missed the traveler. (Premise)
	Exp.Obj.	The researcher encouraged the results. (Premise)
	Ben.For	The uncle bought the nephew for the toy. (Premise)
Task	Single	The cocktail blended the bartender and the intruder cited the patent after the neurologist baffled the hippie. (Stimuli)
	Noisy & Dual	The 5 cocktail + blended 6 the = bartender x5633 and 9 the + authorities 3 agitated = the x5634 organist 6 after + the 8 infantryman = saluted x5635 the 3 pollster. (Stimuli, 1dig.2add.)
Correct Answer	Yes / No	Did the bartender blend the cocktail? (Question)

Table 1: Examples of premises and stimuli depending on factors and variables. Plausibility and Construction display premises, and Task displays stimuli we actually used in the experiment. Abbreviations: Exp.Subj = Experiencer Subject; Exp.Obj. = Experiencer Object; Ben.For = Benefactive For.

string (“x1”, “x2”, “x3”, ...) is appended. Each arithmetic expression is interleaved into the sentence one word at a time. If the end of the sentence is reached in the middle of an arithmetic problem, the remaining calculation problems are not added.

We use ten types of arithmetic problems, varying in both digit length (1, 3, 5, 10, and 30 digits) and the number of addends (two vs. three). They are abbreviated as Xdig.Yadd. (e.g., 1dig.2add. indicates the addition of two one-digit numbers). Example stimuli for some arithmetic types are provided in Appendix A.

All comprehension questions are binary (Yes/No), balanced such that half of the correct answers are Yes and half are No. The final dataset contains 2,560 sentence–question pairs (2 plausibility levels \times 8 constructions \times 160 items).

4 Experiments

4.1 Experimental Setup

We evaluate seven LMs: GPT-4o (OpenAI, 2024), o3-mini, o4-mini², GPT-4.1 (OpenAI, 2024), DeepSeek-V3 (DeepSeek-AI et al., 2025), Llama-3.3 (Grattafiori et al., 2024), and Gemma-3 (Team et al., 2025). The model and prompts are selected based on the following two criteria: (i) accuracy for the implausible condition in the Single Task must

²<https://openai.com/index/introducing-o3-and-o4-mini/>

be at least 70%, and (ii) accuracy for the arithmetic problems in the Dual Task of the 1dig.2add. condition must be at least 80%.³ We set the temperature to 0.0.⁴

4.2 Evaluation Metrics

4.2.1 Preprocessing

Prior to analysis, we filter the data based on model performance on the Single-Task comprehension task and the Dual-Task arithmetic problems.

First, we compute Single-Task accuracy for each construction and plausibility condition and exclude those below 80%. As a result, two constructions (double object and benefactive double object) are excluded for all models, with one additional construction excluded for DeepSeek-V3 and three for Gemma-3. This ensures that analyses include only constructions that models reliably comprehend in the Single Task.

Second, we exclude arithmetic problem types with incorrect answers exceeding 40%. We also remove trials where the arithmetic problem was solved incorrectly or the comprehension response could not be extracted. This filtering step ensures

³We also test GPT-3.5-turbo (<https://platform.openai.com/docs/models>), but it does not meet these criteria.

⁴This operation was restricted for the o3-mini and o4-mini.

that the models analyzed do not adopt strategies that ignore or skip the arithmetic task.

4.2.2 Accuracy of Comprehension Task

We statistically analyze whether the plausibility effect is larger in the Dual Task than in the Single Task and Noisy Single Task, using R (R Core Team, 2025). We use a per-item, non-parametric difference-in-differences procedure. For each item i and task t , we compute mean accuracy \hat{p}_{itp} within each plausibility level p . The within-task plausibility contrast is defined as:

$$\Delta_{it} = \hat{p}_{it,Plausible} - \hat{p}_{it,Implausible} \quad (1)$$

For each item, we then calculate two difference-in-differences contrasts:

$$D_i^{DS} = \Delta_{it,Dual} - \Delta_{it,Single} \quad (2)$$

$$D_i^{DN} = \Delta_{it,Dual} - \Delta_{it,Noisy} \quad (3)$$

Finally, we conduct one-sided Wilcoxon signed-rank tests to assess $H_0 : \text{median}(D) \leq 0$ separately for D_i^{DS} and D_i^{DN} . The null hypothesis H_0 is rejected if the one-sided test is significant at $\alpha = 0.05$.

4.3 Results

Figure 2 shows the mean comprehension accuracy by plausibility, LM, and construction. As seen in the graph, whether the Dual Task promotes rational inference depends on both the LM and the sentence construction. The models can be grouped into three categories as follows.

- (i) GPT-4o is likely to use rational inference in the Dual Task. Four out of six constructions show lower accuracy for implausible sentences in the Dual Task than in either the Single or Noisy Single Tasks.
- (ii) o3-mini and o4-mini show a similar tendency but maintain high accuracy across all conditions (around 100%). Four out of six constructions show significantly lower accuracy for implausible sentences in the Dual Task than in the Single or Noisy Single Tasks.
- (iii) The other models, i.e., GPT-4.1, DeepSeek-V3, Llama-3.3, and Gemma-3, are likely to rely on rational inference in both the Noisy Single and Dual Tasks. These models generally show significantly lower accuracy in the Noisy Single and Dual Tasks than in the

Single Task, but no significant difference between the Noisy Single and Dual Tasks.

In summary, the results suggest that GPT-4o, o3-mini, and o4-mini are more likely to engage in rational inference when cognitive resources are constrained. A consistent trend is observed when analyzing plausibility effects across different arithmetic problems (see Appendix C) and correct answers (see Section 5.1).

5 Analysis

5.1 Effects of Plausibility by Correct Answer

Figure 3 presents the mean accuracy rates of comprehension questions by plausibility, task, LM, and correct answer (Yes or No). When the correct answer was “Yes”, all models show the largest plausibility contrast under the Dual Task condition. This effect was driven by a substantial decrease in accuracy for implausible sentences under the Dual Task. That is, the models often fail to correctly respond “Yes” when asked whether an implausible sentence expressed an implausible meaning, instead responding “No”.

On the other hand, when the correct answer is “No”, GPT-4o, o3-mini, and o4-mini still exhibit tendencies consistent with rational inference, similar to the Yes condition. Other models do not show such a pattern. In summary, consistent with the results in Section 4.3, GPT-4o, o3-mini, and o4-mini reliably demonstrate a shift toward rational inference across both answer types.

5.2 Conditions Where the Implausible Sentences Are Misunderstood

The previous sections show that GPT-4o, o3-mini, and o4-mini tend to shift toward rational inference under the Dual Task. Here, we examine when these models are most likely to misinterpret implausible sentences as plausible under the Dual Task. Figure 4 illustrates the proportion of implausible items that are answered correctly in the Single and Noisy Single Tasks but incorrectly in the Dual Task.

Examining the distribution across constructions in Figure 4 (a), GPT-4o, o3-mini, and o4-mini show the highest error rates for dative and benefactive-for constructions. GPT-4o also frequently fails on passive sentences. These constructions share a key property: when function words (and morphemes) are removed, the remaining word sequence appears semantically plausible. For instance, removing function words from “The bartender was blended



Figure 2: Accuracy of comprehension tasks by plausibility, task, LM, and construction. Some conditions are excluded during preprocessing (see Section 4.2.1). * $p < 0.05$. ** $p < 0.01$. *** $p < 0.001$. Abbreviations: Exp.Subj = Experiencer Subject; Exp.Obj. = Experiencer Object; Ben.For = Benefactive For.

by the cocktail” and “The chef sent the friend to the gift” yields “bartender blend cocktail” and “chef send friend gift” respectively, which could be interpreted as a plausible event. Therefore, it suggests that under resource constraints, these models rely primarily on the superficial word sequence of content words and plausibility derived from world knowledge, rather than function words.

Next, the distribution across the correct answer in Figure 4 (c) shows that errors are more frequent when the correct answer is “No” than when it is “Yes”. This indicates that when asked comprehension questions such as “Did the bartender blend the cocktail?”, the models tend to respond “Yes”, showing a bias toward affirmative answers. This pattern resembles acquiescent, called “yea-saying”, observed in humans during Yes/No question answering tasks (Jackson and Messick, 1958; Knowles

and Condon, 1999). LM’s acquiescent has also been observed (Dentella et al., 2023).

Finally, the calculation graph in Figure 4 (b) shows that all three models fail the most frequently in the 30-digit condition. This suggests that cognitive load increases for the models as the numerical magnitude grows.

6 Discussion

6.1 Limitation of Resources Promotes LMs’ Rational Inference

Although the patterns varied across LMs, our data demonstrate that several models showed a tendency to adopt more rational comprehension strategies under dual-task conditions. This suggests that one of the reasons for errors that prioritize plausibility information over function words is the reduction of

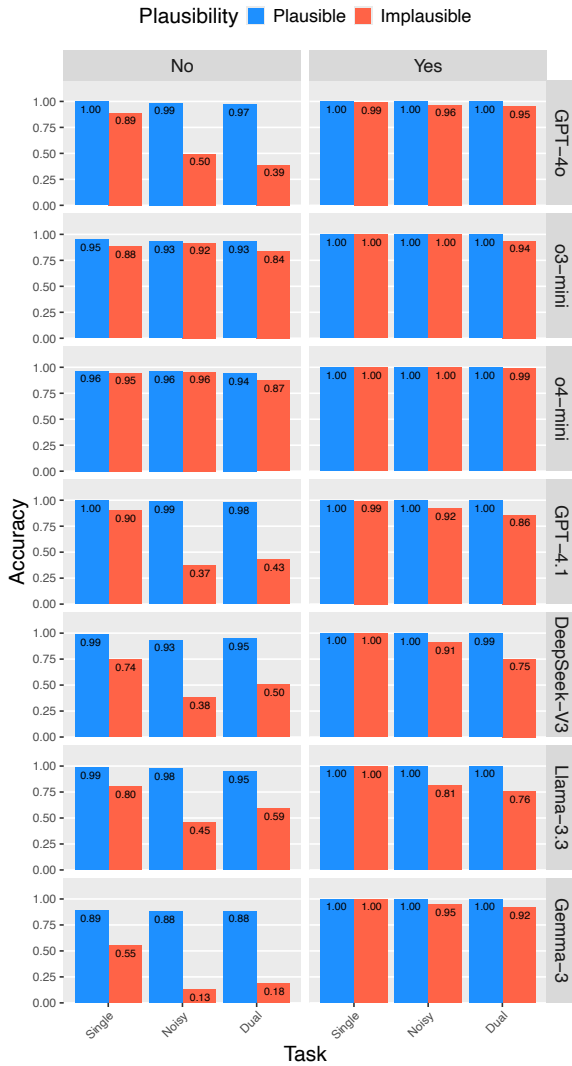


Figure 3: Accuracy of comprehension tasks by plausibility, task, LM, and correct answer.

cognitive resources available for sentence comprehension.

Could it be possible that this plausibility-based shift in the dual task is due to jumbled and longer input? We find that dependence on plausibility increased in the dual task condition compared to the noisy single task condition, even though both included identical arithmetic expressions. Thus, the effect cannot be attributed solely to input complexity. Instead, it reflects the additional cognitive demands imposed by performing two tasks simultaneously, which deplete available resources.

Thus, an additional task is a key to distinguishing the present study from Asami and Sugawara (2024). They manipulate memory load by lengthening sentences: while longer sentences reduced overall accuracy, they did not alter the influence of

plausibility. Therefore, our results suggest that task-induced cognitive load, rather than input length alone, is a critical factor in constraining LMs' cognitive resources. Human working memory is not simply about storage, but about the dynamic interaction between memory access and processing (Atkinson and Shiffrin, 1971; Baddeley and Hitch, 1974; Baddeley, 2003). Therefore, our findings suggest that LMs, like humans, rely not only on short- or long-term memory stores, but also on a working-memory-like mechanism that integrates memory with ongoing computation.

Our results also align with findings from reasoning studies. LMs show heuristic reasoning strategies, called "shortcut solution" (Geirhos et al., 2020; Jia and Liang, 2017; Ko et al., 2021; Tang et al., 2023). They sometimes rely on superficial letter sequences rather than the content of the documents, similar to this study. Furthermore, it has been observed that when cognitive resources become limited, reasoning processing shifts from syntactic interpretation to more superficial comprehension based on word-level associations (Lampinen et al., 2024; Zhang et al., 2024) and degrades some functions, such as safety mechanisms (Upadhyay et al., 2025; Xu et al., 2024). Our work further explores what functions are reduced in sentence comprehension. In particular, our data show that adding an arithmetic computation task depletes the resources and consequently reduces syntactic processing.

6.2 Contributions to Psycholinguistics

Finally, we discuss how our results can contribute to Psycholinguistic theories. Our findings suggest that limiting cognitive resources induces a shift toward rational inference, i.e., a human-like comprehension strategy. More broadly, these results support the hypothesis that human-like sentence understanding fundamentally arises from how limited cognitive resources are allocated.

In human sentence comprehension, behavioral effects are often attributed to memory limitations (Gibson, 1998; Van Dyke and Lewis, 2003, and others). Working memory, particularly its temporary storage used for ongoing processing, has long been considered a crucial factor. However, it remains an open question whether human comprehension behavior can truly be explained solely in terms of working memory capacity.

This study provides supporting evidence from non-human systems, namely LMs, that their be-

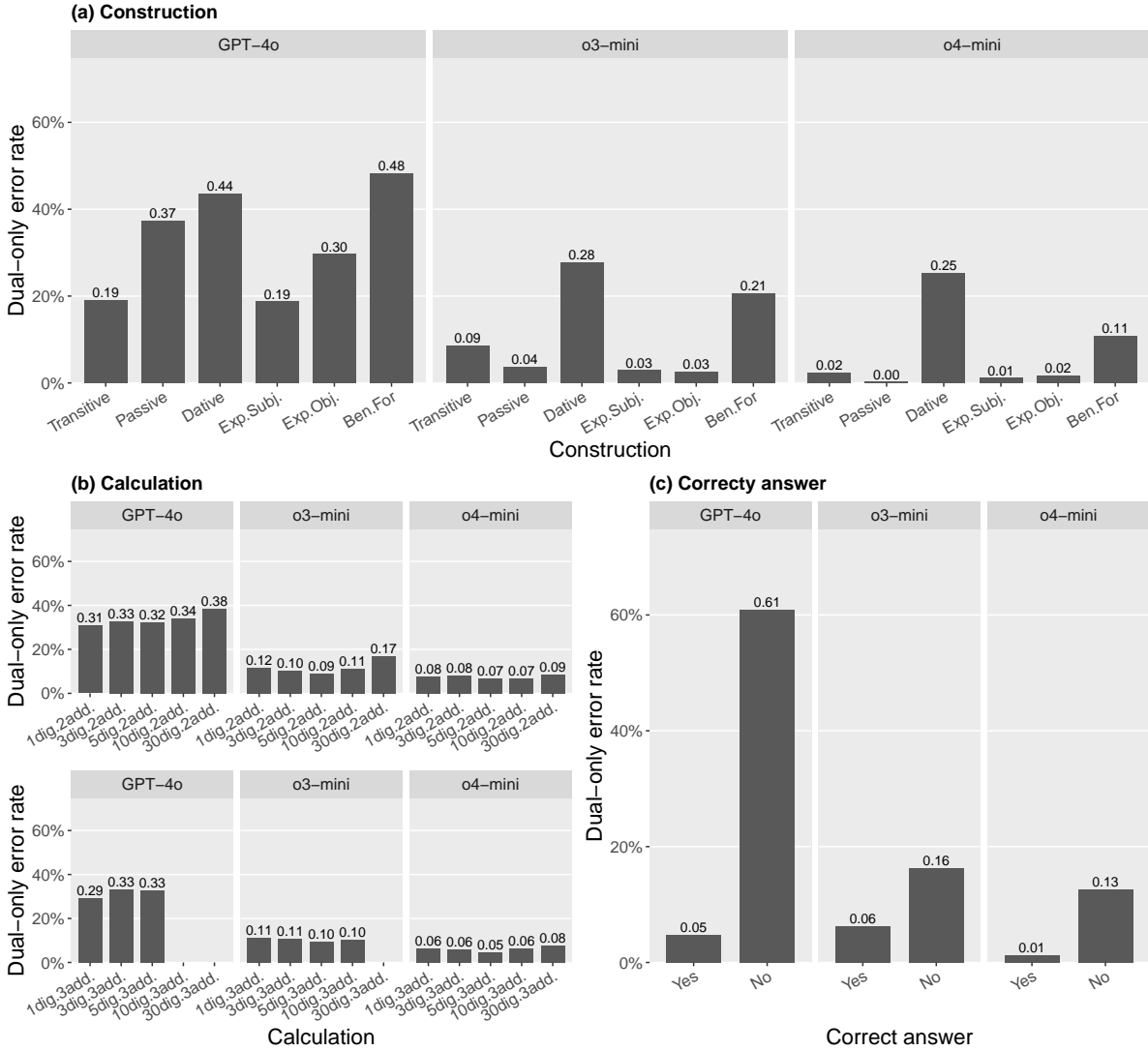


Figure 4: Proportion of implausible sentences correct in single and noisy-single tasks but incorrect in the dual task. Some conditions are excluded during preprocessing (see Section 4.2.1. Abbreviations: Exp.Subj = Experiencer Subject; Exp.Obj. = Experiencer Object; Ben.For = Benefactive For).

506 havior under constrained conditions resembles hu- 522
 507 man behavior. Namely, under the limited resources, 523
 508 LMs and humans (i) reduced reliance on syntactic 524
 509 function, with increased reliance on world knowl- 525
 510 edge and superficial word order (Gibson et al., 526
 511 2016; Rogalsky et al., 2008), and (ii) a bias toward 527
 512 acquiescence (Condon et al., 2006; Knowles 528
 513 and Condon, 1999; Lechner and Rammstedt, 2015; 529
 514 Rammstedt et al., 2023). That is, both humans 530
 515 and LMs degrade certain functions supporting syn- 531
 516 tactic processing and rejection when the cognitive 532
 517 resources are constrained. As a result, they rely 533
 518 more on semantic plausibility. 534

519 Taken together, our results suggest that the un- 535
 520 derlying principle of human sentence comprehen-
 521 sion may lie in the resource limitation of working

memory, leading to strategies adopted to achieve 522
 efficient sentence comprehension (cf. Cognitive 523
 Load Theory by Sweller (1988) and Sweller et al. 524
 (2019)). 525

7 Conclusion 526

527 We implement a dual-task paradigm in which the 528
 529 models simultaneously solve arithmetic problems 530
 and answer comprehension questions. GPT-4o, 531
 o3-mini, and o4-mini shift their comprehension 532
 strategies under cognitive resource limitations to- 533
 ward rational inference, similar to humans. Our 534
 findings suggest that task-induced cognitive load, 535
 rather than input length alone, constrains LMs' cog-
 nitive resources and makes them more human-like.

536 Limitations

537 First, further research is needed to explore what
538 drives LMs’ rational reading strategies under dual-
539 task conditions. While some models demonstrate
540 a clear shift toward plausibility-based comprehen-
541 sion, others do not. Notably, even within the GPT-4
542 family, GPT-4o exhibits rational inference behavior,
543 whereas GPT-4.1 does not, indicating that identify-
544 ing the source of this difference is not straightfor-
545 ward. It remains unclear whether these variations
546 arise from differences in internal architecture, train-
547 ing procedures, or other aspects of the model.⁵
548 Recent studies have proposed that attention mech-
549 anisms in LMs serve as a working-memory-like
550 component in the context of modeling sentence pro-
551 cessing costs (Ryu and Lewis, 2021; Timkey and
552 Linzen, 2023; Yoshida et al., 2025). Thus, future
553 research could extend this line of inquiry to com-
554 prehension strategies, potentially providing new
555 insights into human working memory processes
556 during comprehension.

557 The second limitation of this study is that the
558 results may depend on the choice of baseline and
559 prompt design. Although our Noisy Single Task
560 is designed to isolate the effect of adding a cal-
561 culation task by comparing the Dual Task, other
562 baselines or formatting choices may yield differ-
563 ent outcomes. This is relevant because large LMs
564 are known to be sensitive to prompt formulation
565 (Kojima et al., 2022; Schmidt et al., 2024; Sclar
566 et al., 2024). Thus, systematically varying prompts,
567 including non-semantic interruptions such as de-
568 limiters or formatting tokens instead of numerics,
569 would be an important direction for future work.

570 Finally, direct human–LM comparisons would
571 be highly valuable for more detailed modeling of
572 human sentence comprehension. Although human
573 rational inference in similar comprehension tasks is
574 well established in the psycholinguistic literature,
575 differences between humans and LMs may be sen-
576 sitive to task design or stimulus properties. Import-
577 antly, establishing a robust and well-characterized
578 dual-task paradigm for LMs is a necessary prerequi-
579 site for meaningful human–LM comparison, as pre-
580 mature comparisons risk conflating methodological
581 artifacts with cognitive effects. Our contribution

⁵We examined whether language models acquire a rational inference strategy during training using OLMo (Groeneveld et al., 2024), which provides access to both training data and intermediate training checkpoints. However, even the final models fail to solve the arithmetic problems in the dual-task condition.

582 should therefore be viewed as a first step: intro-
583 ducing a dual-task paradigm for LMs and demon-
584 strating its potential to reveal rational inference
585 behavior.

586 Taken together, future work should therefore en-
587 hance the dual-task paradigm by exploring alter-
588 native baselines and prompt designs. After estab-
589 lishing robust experimental settings, subsequent
590 work can pursue more detailed human–LM com-
591 parisons and deeper investigation into the internal
592 mechanisms underlying working memory in LMs.

593 References

- 594 Samuel Joseph Amouyal, Aya Meltzer-Asscher, and
595 Jonathan Berant. 2025a. [Comparing human and
596 language models sentence processing difficulties on
597 complex structures](#). *Preprint*, arXiv:2510.07141.
- 598 Samuel Joseph Amouyal, Aya Meltzer-Asscher, and
599 Jonathan Berant. 2025b. [When the lm misunder-
600 stood the human chuckled: Analyzing garden path
601 effects in humans and language models](#). *Preprint*,
602 arXiv:2502.09307.
- 603 Daiki Asami and Saku Sugawara. 2024. [What makes
604 language models good-enough?](#) In *Findings of
605 the Association for Computational Linguistics: ACL
606 2024*, pages 15453–15467, Bangkok, Thailand. As-
607 sociation for Computational Linguistics.
- 608 Richard C. Atkinson and Richard M. Shiffrin. 1971. [The control of short-term memory](#). *Scientific Ameri-
609 can*, 225(2):82–91.
- 611 Nicolai D. Ayasse, Alana J. Hodson, and Arthur Wing-
612 field. 2021. [The principle of least effort and compre-
613 hension of spoken sentences by younger and older
614 adults](#). *Frontiers in Psychology*, Volume 12 - 2021.
- 615 Alan Baddeley. 2003. [Working memory and language:
616 An overview](#). *Journal of communication disorders*,
617 36(3):189–208.
- 618 Alan D. Baddeley and Graham Hitch. 1974. [Working
619 memory](#). volume 8 of *Psychology of Learning and
620 Motivation*, pages 47–89. Academic Press.
- 621 Zhoujun Cheng, Jungo Kasai, and Tao Yu. 2023. [Batch
622 prompting: Efficient inference with large language
623 model APIs](#). In *Proceedings of the 2023 Conference
624 on Empirical Methods in Natural Language Process-
625 ing: Industry Track*, pages 792–810, Singapore. As-
626 sociation for Computational Linguistics.
- 627 Lorena Condon, Pere J. Ferrando, and Josep Demestre.
628 2006. [A note on some item characteristics related to
629 acquiescent responding](#). *Personality and Individual
630 Differences*, 40(3):403–407.
- 631 DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, Bingx-
632 uan Wang, Bochao Wu, Chengda Lu, Chenggang

633	Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan,	<i>the Association for Computational Linguistics (Vol-</i>	688
634	Damai Dai, Daya Guo, Dejian Yang, Deli Chen,	<i>ume 1: Long Papers)</i> , pages 15789–15809, Bangkok,	689
635	Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai,	Thailand. Association for Computational Linguistics.	690
636	and 181 others. 2025. Deepseek-v3 technical report .		
637	<i>Preprint</i> , arXiv:2412.19437.		
638	Vittoria Dentella, Fritz Günther, and Evelina Leivada.	Michael Hahn, Richard Futrell, Roger Levy, and Ed-	691
639	2023. Systematic testing of three language mod-	ward Gibson. 2022. A resource-rational model of	692
640	els reveals low language accuracy, absence of re-	human processing of recursive linguistic structure .	693
641	sponse stability, and a yes-response bias . <i>Pro-</i>	<i>ceedings of the National Academy of Sciences</i> ,	694
642	<i>ceedings of the National Academy of Sciences</i> ,	119(43):e2122602119.	695
643	120(51):e2309583120.		
644	Richard Futrell and Edward Gibson. 2017. L2 pro-	Tovah Irwin, Kyra Wilson, and Alec Marantz. 2023.	696
645	cessing as noisy channel language comprehension .	BERT shows garden path effects . In <i>Proceedings</i>	697
646	<i>Bilingualism: Language and Cognition</i> , 20(4):683–	<i>of the 17th Conference of the European Chapter</i>	698
647	684.	<i>of the Association for Computational Linguistics</i> ,	699
		pages 3220–3232, Dubrovnik, Croatia. Association	700
		for Computational Linguistics.	701
648	Richard Futrell, Edward Gibson, and Roger P. Levy.	Douglas N Jackson and Samuel Messick. 1958. Content	702
649	2020. Lossy-context surprisal: An information-	and style in personality assessment . <i>Psychological</i>	703
650	theoretic model of memory effects in sentence pro-	<i>bulletin</i> , 55(4):243.	704
651	cessing . <i>Cognitive Science</i> , 44(3):e12814.		
652	Robert Geirhos, Jörn-Henrik Jacobsen, Claudio	Robin Jia and Percy Liang. 2017. Adversarial exam-	705
653	Michaelis, Richard Zemel, Wieland Brendel,	ples for evaluating reading comprehension systems .	706
654	Matthias Bethge, and Felix A Wichmann. 2020.	<i>Preprint</i> , arXiv:1707.07328.	707
655	Shortcut learning in deep neural networks . <i>Nature</i>		
656	<i>Machine Intelligence</i> , 2(11):665–673.	Wayne K Kirchner. 1958. Age differences in short-term	708
		retention of rapidly changing information . <i>Journal</i>	709
		<i>of experimental psychology</i> , 55(4):352.	710
657	Edward Gibson. 1998. Linguistic complexity: locality	Eric S Knowles and Christopher A Condon. 1999. Why	711
658	of syntactic dependencies . <i>Cognition</i> , 68(1):1–76.	people say "yes": A dual-process theory of acquies-	712
		cence . <i>Journal of Personality and Social Psychology</i> ,	713
659	Edward Gibson, Steven T. Piantadosi, Kimberly Brink,	77(2):379.	714
660	Leon Bergen, Eunice Lim, and Rebecca Saxe. 2013.	Miyoung Ko, Jinhyuk Lee, Hyunjae Kim, Gangwoo	715
661	A noisy-channel account of crosslinguistic word-	Kim, and Jaewoo Kang. 2021. Look at the first sen-	716
662	order variation . <i>Psychological Science</i> , 24(7):1079–	tence: Position bias in question answering . <i>Preprint</i> ,	717
663	1088.	arXiv:2004.14602.	718
664	Edward Gibson, Chaleece Sandberg, Evelina Fedorenko,	Takeshi Kojima, Shixiang (Shane) Gu, Machel Reid, Yu-	719
665	Leon Bergen, and Swathi Kiran. 2016. A rational in-	taka Matsuo, and Yusuke Iwasawa. 2022. Large lan-	720
666	ference approach to aphasic language comprehension .	guage models are zero-shot reasoners . In <i>Advances in</i>	721
667	<i>Aphasiology</i> , 30(11):1341–1360.	<i>Neural Information Processing Systems</i> , volume 35,	722
		pages 22199–22213. Curran Associates, Inc.	723
668	Dongyu Gong, Xingchen Wan, and Dingmin Wang.	Tatsuki Kuribayashi, Yohei Oseki, Ana Brassard, and	724
669	2024. Working memory capacity of chatgpt: An	Kentaro Inui. 2022. Context limitations make neural	725
670	empirical study . <i>Proceedings of the AAAI Confer-</i>	language models more human-like . In <i>Proceedings</i>	726
671	ence on Artificial Intelligence , 38(9):10048–10056.	<i>of the 2022 Conference on Empirical Methods in</i>	727
		<i>Natural Language Processing</i> , pages 10421–10436,	728
672	Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri,	Abu Dhabi, United Arab Emirates. Association for	729
673	Abhinav Pandey, Abhishek Kadian, Ahmad Al-	Computational Linguistics.	730
674	Dahle, Aiesha Letman, Akhil Mathur, Alan Schel-		
675	ten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh	Tatsuki Kuribayashi, Yohei Oseki, Takumi Ito, Ryo	731
676	Goyal, Anthony Hartshorn, Aobo Yang, Archi Mi-	Yoshida, Masayuki Asahara, and Kentaro Inui. 2021.	732
677	tra, Archie Sravankumar, Artem Korenev, Arthur	Lower perplexity is not always human-like . In <i>Pro-</i>	733
678	Hinsvark, and 542 others. 2024. The llama 3 herd of	<i>ceedings of the 59th Annual Meeting of the Associa-</i>	734
679	models . <i>Preprint</i> , arXiv:2407.21783.	<i>tion for Computational Linguistics and the 11th Inter-</i>	735
		<i>national Joint Conference on Natural Language Pro-</i>	736
680	Dirk Groeneveld, Iz Beltagy, Evan Walsh, Akshita	<i>cessing (Volume 1: Long Papers)</i> , pages 5203–5217,	737
681	Bhagia, Rodney Kinney, Oyvind Tafjord, Ananya	Online. Association for Computational Linguistics.	738
682	Jha, Hamish Ivison, Ian Magnusson, Yizhong Wang,		
683	Shane Arora, David Atkinson, Russell Authur, Khy-	Andrew K Lampinen, Ishita Dasgupta, Stephanie C Y	739
684	athi Chandu, Arman Cohan, Jennifer Dumas, Yanai	Chan, Hannah R Sheahan, Antonia Creswell, Dhar-	740
685	Elazar, Yuling Gu, Jack Hessel, and 24 others. 2024.	shan Kumaran, James L McClelland, and Felix	741
686	OLMo: Accelerating the science of language mod-	Hill. 2024. Language models, like humans, show	742
687	els . In <i>Proceedings of the 62nd Annual Meeting of</i>		

743	content effects on reasoning tasks. <i>PNAS Nexus</i> , 3(7):pgae233.	798
744		799
745	Clemens M Lechner and Beatrice Rammstedt. 2015. Cognitive ability, acquiescence, and the structure of personality in a sample of older adults. <i>Psychological assessment</i> , 27(4):1301.	800
746		801
747		802
748		803
749	Byung-Doh Oh, Christian Clark, and William Schuler. 2022. Comparison of structural parsers and neural language models as surprisal estimators. <i>Frontiers in Artificial Intelligence</i> , 5:777963.	804
750		805
751		806
752		807
753	Byung-Doh Oh and William Schuler. 2023. Why does surprisal from larger transformer-based language models provide a poorer fit to human reading times? <i>Transactions of the Association for Computational Linguistics</i> , 11:336–350.	808
754		809
755		810
756		811
757		812
758	OpenAI. 2024. GPT-4 technical report. <i>Preprint</i> , arXiv:2303.08774.	813
759		814
760	R Core Team. 2025. <i>R: A Language and Environment for Statistical Computing</i> . R Foundation for Statistical Computing, Vienna, Austria.	815
761		816
762		817
763	Beatrice Rammstedt, Lena Roemer, and Clemens M. Lechner. 2023. Do simpler item wording and response scales reduce acquiescence in personality inventories? a survey experiment. <i>Personality and Individual Differences</i> , 214:112324.	818
764		819
765		820
766		821
767		822
768	Corianne Rogalsky, William Matchin, and Gregory Hickok. 2008. Broca’s area, sentence comprehension, and working memory: an fmri study. <i>Frontiers in human neuroscience</i> , 2:327.	823
769		824
770		825
771		826
772	Soo Hyun Ryu and Richard Lewis. 2021. Accounting for agreement phenomena in sentence comprehension with transformer language models: Effects of similarity-based interference on surprisal and attention. In <i>Proceedings of the Workshop on Cognitive Modeling and Computational Linguistics</i> , pages 61–71, Online. Association for Computational Linguistics.	827
773		828
774		829
775		830
776		831
777		832
778		833
779		834
780	Douglas C. Schmidt, Jesse Spencer-Smith, Quchen Fu, and Jules White. 2024. Towards a catalog of prompt patterns to enhance the discipline of prompt engineering. <i>Ada Lett.</i> , 43(2):43–51.	835
781		836
782		837
783		838
784	Melanie Sclar, Yejin Choi, Yulia Tsvetkov, and Alane Suhr. 2024. Quantifying language models’ sensitivity to spurious features in prompt design or: How i learned to start worrying about prompt formatting. In <i>The Twelfth International Conference on Learning Representations</i> .	839
785		840
786		841
787		842
788		843
789		844
790	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 <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 5606–5627, Bangkok, Thailand. Association for Computational Linguistics.	845
791		846
792		847
793		848
794		849
795		850
796		851
797		852
	John Sweller. 1988. Cognitive load during problem solving: Effects on learning. <i>Cognitive science</i> , 12(2):257–285.	
	John Sweller, Jeroen JG Van Merriënboer, and Fred Paas. 2019. Cognitive architecture and instructional design: 20 years later. <i>Educational psychology review</i> , 31(2):261–292.	
	Ruixiang Tang, Dehan Kong, Longtao Huang, and Hui Xue. 2023. Large language models can be lazy learners: Analyze shortcuts in in-context learning. In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 4645–4657, Toronto, Canada. Association for Computational Linguistics.	
	Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane Rivière, Louis Rouillard, Thomas Mesnard, Geoffrey Cideron, Jean bastien Grill, Sabela Ramos, Edouard Yvinec, Michelle Casbon, Etienne Pot, Ivo Penchev, and 197 others. 2025. Gemma 3 technical report. <i>Preprint</i> , arXiv:2503.19786.	
	William Timkey and Tal Linzen. 2023. A language model with limited memory capacity captures interference in human sentence processing. In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 8705–8720, Singapore. Association for Computational Linguistics.	
	Marilyn L Turner and Randall W Engle. 1989. Is working memory capacity task dependent? <i>Journal of Memory and Language</i> , 28(2):127–154.	
	Bibek Upadhayay, Vahid Behzadan, and Amin Karbasi. 2025. Working memory attack on LLMs. In <i>ICLR 2025 Workshop on Building Trust in Language Models and Applications</i> .	
	Julie A Van Dyke and Richard L Lewis. 2003. Distinguishing effects of structure and decay on attachment and repair: A cue-based parsing account of recovery from misanalyzed ambiguities. <i>Journal of Memory and Language</i> , 49(3):285–316.	
	Ethan Gotlieb Wilcox, Michael Y. Hu, Aaron Mueller, Alex Warstadt, Leshem Choshen, Chengxu Zhuang, Adina Williams, Ryan Cotterell, and Tal Linzen. 2025. Bigger is not always better: The importance of human-scale language modeling for psycholinguistics. <i>Journal of Memory and Language</i> , 144:104650.	
	Nan Xu, Fei Wang, Ben Zhou, Bangzheng Li, Chaowei Xiao, and Muhao Chen. 2024. Cognitive overload: Jailbreaking large language models with overloaded logical thinking. In <i>Findings of the Association for Computational Linguistics: NAACL 2024</i> , pages 3526–3548, Mexico City, Mexico. Association for Computational Linguistics.	
	Ryo Yoshida, Shinnosuke Isono, Kohei Kajikawa, Taiga Someya, Yushi Sugimoto, and Yohei Oseki. 2025. If attention serves as a cognitive model of human	

memory retrieval, what is the plausible memory representation? In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, page 9795–9812. Association for Computational Linguistics.

Chunhui Zhang, Yiren Jian, Zhongyu Ouyang, and Soroush Vosoughi. 2024. *Working memory identifies reasoning limits in language models*. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 16896–16922, Miami, Florida, USA. Association for Computational Linguistics.

A Example stimuli for the Noisy Single Task and the Dual Task

Table 2 represents example stimuli for the noisy single task and the dual task.

B Prompts

The method to make prompts is as follows. We put instructions of tasks and few-shot examples on the system message, and a sentence and a question on the user message. The few-shot example contains four sets of sentence and question, each of these are plausible or implausible, and question whose answer is Yes or No. Examples of prompts are in Figures 5–7.

C Effect of plausibility by calculation

Figure 8 represents the mean accuracy rates of comprehension questions by plausibility, LM, and calculation. The graph shows the consistent results with the results by construction and correct answer (see Figure 2 and 3): GPT-4o, o3-mini, and o4-mini have the greater difference between plausible and implausible conditions in the dual task than the single task or the noisy single task. DeepSeek-V3 also represents this trend.

D Accuracy for Calculation

The accuracy of arithmetic problems is displayed in Table 3. The table shows that most of the LMs have around 90% accuracy except for Llama-3.3, which has below 80% in any conditions.

Single Task

System Message:

Please read the sentence. Once the sentence ends, answer the yes/no question about the sentence. The question relates to the sentence and requires a yes or no response. Return the question answers as [Yes] or [No]. The following are examples of sentences and questions about it.

Sentence: The boy is eating an apple and a banana in his house.
Question: Is the boy eating an apple?
Answer: [Yes]

Sentence: The teacher put a cup on the table and read his favorite book while listening to music.
Question: Did a book read the teacher?
Answer: [No]

Context: The dog was chased by the man before the woman saw the pianist.
Question: Was the man chased by the dog?
Answer: [No]

Context: The refrigerator painted the carpenter while the mountain laughed at the sandwich and the sky was green.
Question: Did the refrigerator paint the carpenter?
Answer: [Yes]

User Message:

Sentence: The bartender was blended by the cocktail.
Question: Did the bartender blend the cocktail?
Answer:

Figure 5: Example prompt of the single task.

Calculation	Stimuli
1dig.2add.	The 5 cocktail + blended 6 the = bartender x5633 and 9 the + authorities 3 agitated = the x5634 organist 6 after + the 8 infantryman = saluted x5635 the 3 pollster.
1dig.3add.	The 7 cocktail + blended 6 the + bartender 5 and = the x4210 authorities 2 agitated + the 4 organist + after 4 the = infantryman x4211 saluted 4 the + pollster.
3dig.2add.	The 952 cocktail + blended 604 the = bartender x5633 and 793 the + authorities 271 agitated = the x5634 organist 770 after + the 832 infantryman = saluted x5635 the 531 pollster.
3dig.3add.	The 212 cocktail + blended 260 the + bartender 341 and = the x4210 authorities 220 agitated + the 631 organist + after 753 the = infantryman x4211 saluted 264 the + pollster.

Table 2: Example stimuli by arithmetic problem. Context stimuli in the noisy single task and the dual task embed arithmetic problems in sentences.

Noisy Single Task

System Message:

The sentence combines a math problem and a text alternating one word at a time. The math problem is split and placed after each word in the sentence, such as $2 + 8 = x1$. Please read the sentence while ignoring the math problems. Once the sentence ends, answer the yes/no question about the sentence. The question relates to the sentence and requires a yes or no response. Return the question answers as [Yes] or [No]. The following are examples of sentences with math problems and a yes/no question about it.

Sentence: The 2 boy + is 8 eating = an x1 apple 1 and + a 5 banana = in x2 his 3 house.
Question: Is the boy eating an apple?
Answer: [Yes]

Sentence: The 4 teacher + put 9 a + cup 7 on = the x3 table 1 and + read 4 his + favorite 6 book = while x4 listening 5 to + music.
Question: Did a book read the teacher?
Answer: [No]

Context: The 5 dog + was 5 chased = by x6 the 2 man + before 7 the = woman x7 saw 3 the + pianist.
Question: Was the man chased by the dog?
Answer: [No]

Context: The 3 refrigerator + painted 6 the + carpenter 2 while = the x9 mountain 8 laughed + at 1 the + sandwich 7 and = the x10 sky 1 was + green.
Question: Did the refrigerator paint the carpenter?"
Answer: [Yes]

User Message:

Sentence: The 5 cocktail + blended 6 the = bartender x5633 and 9 the + authorities 3 agitated = the x5634 organist 6 after + the 8 infantryman = saluted x5635 the 3 pollster.
Question: Did the bartender blend the cocktail?
Answer:

Figure 6: Example prompt of the noisy single task.

Dual Task

System Message:

The sentence combines a math problem and a text alternating one word at a time. The math problem is split and placed after each word in the sentence, such as $2 + 8 = x1$. Please read the sentence while accurately calculating the math problems. When $x1$, $x2$, $x3...$ appear, output the answer to the math problem carefully. Once the sentence ends, answer the yes/no question about the sentence. The question relates to the sentence and requires a yes or no response. However, prioritize the quality of solving the math problems over answering the sentence. Ensure all math problems are solved correctly, with each one being an addition of two or three numbers. Return the math problems and their answers as tuples, e.g., $(x1, 2 + 8, 10)$, $(x3, 4 + 9 + 7, 20)$. Return the question answers as [Yes] or [No]. The following are examples of sentences with math problems and a yes/no question about it.

Sentence: The 2 boy + is 8 eating = an x1 apple 1 and + a 5 banana = in x2 his 3 house.
Question: Is the boy eating an apple?
Answer for the math problem: $(x1, 2 + 8, 10)$, $(x2, 1 + 5, 6)$
Answer for the question: [Yes]

Sentence: The 4 teacher + put 9 a + cup 7 on = the x3 table 1 and + read 4 his + favorite 6 book = while x4 listening 5 to + music.
Question: Did a book read the teacher?
Answer for the math problem: $(x3, 4 + 9 + 7, 20)$, $(x4, 1 + 4 + 6, 11)$
Answer for the question: [No]

Context: The 5 dog + was 5 chased = by x6 the 2 man + before 7 the = woman x7 saw 3 the + pianist.
Question: Was the man chased by the dog?
Answer for the math problem: $(x6, 5 + 5, 10)$, $(x7, 2 + 7, 9)$
Answer for the question: [No]

Context: The 3 refrigerator + painted 6 the = carpenter x9 while 8 the + mountain 1 laughed = at x10 the 7 sandwich + and 9 the = sky x11 was 2 green.
Question: Did the refrigerator paint the carpenter?
Answer for the math problem: $(x9, 3 + 6, 9)$, $(x10, 8 + 1, 9)$, $(x11, 7 + 9, 16)$
Answer for the question: [Yes]

User Message:

Sentence: The 5 cocktail + blended 6 the = bartender x5633 and 9 the + authorities 3 agitated = the x5634 organist 6 after + the 8 infantryman = saluted x5635 the 3 pollster.
Question: Did the bartender blend the cocktail?
Answer for the math problem:
Answer for the question:

Figure 7: Example prompt of the dual task.

Model	Calculation	Mean	SD
GPT-4o	1dig.2add.	0.99	0.12
GPT-4o	1dig.3add.	1.00	0.00
GPT-4o	3dig.2add.	0.99	0.11
GPT-4o	3dig.3add.	1.00	0.00
GPT-4o	5dig.2add.	0.99	0.10
GPT-4o	5dig.3add.	0.99	0.11
GPT-4o	10dig.2add.	0.92	0.26
GPT-4o	10dig.3add.	0.50	0.50
GPT-4o	30dig.2add.	0.77	0.42
GPT-4o	30dig.3add.	0.04	0.19
<hr/>			
o3-mini	1dig.2add.	0.97	0.17
o3-mini	1dig.3add.	0.98	0.14
o3-mini	3dig.2add.	0.95	0.22
o3-mini	3dig.3add.	0.98	0.14
o3-mini	5dig.2add.	0.95	0.22
o3-mini	5dig.3add.	0.97	0.17
o3-mini	10dig.2add.	0.93	0.25
o3-mini	10dig.3add.	0.93	0.25
o3-mini	30dig.2add.	0.71	0.45
o3-mini	30dig.3add.	0.58	0.49
<hr/>			
o4-mini	1dig.2add.	0.97	0.16
o4-mini	1dig.3add.	0.95	0.21
o4-mini	3dig.2add.	0.97	0.17
o4-mini	3dig.3add.	0.94	0.23
o4-mini	5dig.2add.	0.96	0.20
o4-mini	5dig.3add.	0.93	0.25
o4-mini	10dig.2add.	0.94	0.24
o4-mini	10dig.3add.	0.91	0.29
o4-mini	30dig.2add.	0.82	0.38
o4-mini	30dig.3add.	0.76	0.43
<hr/>			
GPT-4.1	1dig.2add.	1.00	0.00
GPT-4.1	1dig.3add.	1.00	0.05
GPT-4.1	3dig.2add.	1.00	0.00
GPT-4.1	3dig.3add.	1.00	0.01
GPT-4.1	5dig.2add.	1.00	0.04
GPT-4.1	5dig.3add.	0.99	0.11
GPT-4.1	10dig.2add.	0.92	0.28
GPT-4.1	10dig.3add.	0.12	0.33
GPT-4.1	30dig.2add.	0.64	0.48
GPT-4.1	30dig.3add.	0.03	0.18
<hr/>			
DeepSeek-V3	1dig.2add.	0.99	0.11
DeepSeek-V3	1dig.3add.	0.99	0.09
DeepSeek-V3	3dig.2add.	0.99	0.12
DeepSeek-V3	3dig.3add.	0.99	0.08
DeepSeek-V3	5dig.2add.	0.97	0.17
DeepSeek-V3	5dig.3add.	0.99	0.10
DeepSeek-V3	10dig.2add.	0.90	0.30
DeepSeek-V3	10dig.3add.	0.84	0.37
DeepSeek-V3	30dig.2add.	0.91	0.29
DeepSeek-V3	30dig.3add.	0.58	0.49
<hr/>			
Llama-3.3	1dig.2add.	0.82	0.38
Llama-3.3	1dig.3add.	0.73	0.44
Llama-3.3	3dig.2add.	0.76	0.43
Llama-3.3	3dig.3add.	0.50	0.50
Llama-3.3	5dig.2add.	0.78	0.42
Llama-3.3	5dig.3add.	0.42	0.49
Llama-3.3	10dig.2add.	0.66	0.47
Llama-3.3	10dig.3add.	0.26	0.44
Llama-3.3	30dig.2add.	0.24	0.43
Llama-3.3	30dig.3add.	0.01	0.08
<hr/>			
Gemma-3	1dig.2add.	0.89	0.32
Gemma-3	1dig.3add.	0.82	0.39
Gemma-3	3dig.2add.	0.93	0.26
Gemma-3	3dig.3add.	0.83	0.38
Gemma-3	5dig.2add.	0.81	0.39
Gemma-3	5dig.3add.	0.72	0.45
Gemma-3	10dig.2add.	0.56	0.50
Gemma-3	10dig.3add.	0.44	0.50
Gemma-3	30dig.2add.	0.18	0.38
Gemma-3	30dig.3add.	0.00	0.00

Table 3: Accuracy of arithmetic problems. Some conditions are excluded during preprocessing (see Section 4.2.1). SD = standard deviation.

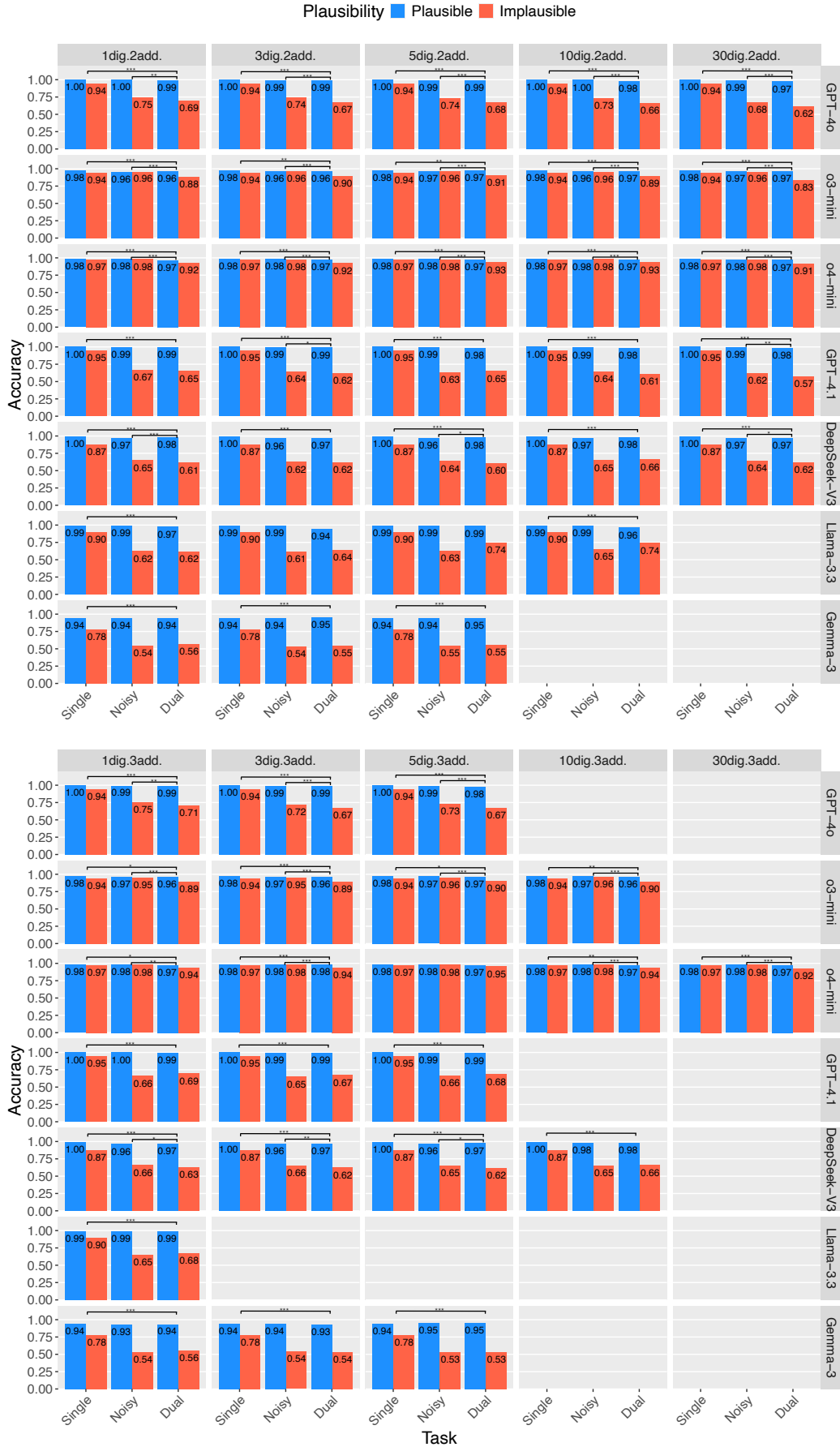


Figure 8: Accuracy rate of comprehension tasks by plausibility, task, LM, and calculation. Some conditions are excluded during preprocessing (see Section 4.2.1). $*p < 0.05$. $**p < 0.01$. $***p < 0.001$.