# I am a Strange Dataset: Metalinguistic Tests for Language Models

Anonymous ACL submission

### Abstract

Statements involving metalinguistic selfreference ("This paper has six sections.") are prevalent in many domains. Can large language models (LLMs) handle such language? In this 004 paper, we present "I am a Strange Dataset", a new dataset for addressing this question. 007 There are two subtasks: generation and verification. In generation, models continue statements like "The penultimate word in this sentence is" (where a correct continuation is "is"). In verification, models judge the truth of statements like "The penultimate word in 012 this sentence is sentence." (false). We also provide minimally different metalinguistic non-self-reference examples to complement the main dataset by probing for whether models can handle metalinguistic language 017 at all. The dataset is hand-crafted by experts and validated by non-expert annotators. We test a variety of open-source LLMs (7B to 70B parameters) as well as closed-source LLMs through APIs. All models perform close to chance across both subtasks and even on the non-self-referential metalinguistic control data, though we find some steady improvement with model scale. GPT 4 is the only model to consistently do significantly better than chance, 027 and it is still only in the 60% range, while our untrained human annotators score well in the 89-93% range.

### 1 Introduction

Self-reference plays a crucial role in the way we think about mathematics (Gödel, 1931), theoretical computer science (Church, 1936), recursive programming (Hofstadter, 1979), philosophy (Tarski, 1931), understanding complex cases in hate speech detection (Allan, 2017), aptitude tests (Propp, 1993), and comedy (Hofstadter, 1985). Some positions in the philosophy of mind consider selfreferential capabilities to be a key aspect of higher intelligence or even consciousness (Hofstadter, 2007; Baars, 1993). Of course, self-reference is  $P\left(No. \begin{vmatrix} \text{if someone asks whether} \\ \text{this sentence has a capital} \\ \text{letter, the correct answer is} \end{vmatrix}\right)$ 

Figure 1: An example highlighting the challenge presented by our task. All models that we tested on our dataset are close to chance-level.

also pervasive in how we communicate: at least one paper you read today is bound to contain "In this paper" (Anonymous, 2024). 043

044

046

050

058

060

061

062

063

064

065

066

067

068

069

070

071

In this paper, we focus on metalinguistic self-reference, the complex kind of self-reference in which language is used to make claims about it-self, as in "This sentence has five words" and "This paper has six sections".<sup>1</sup> Using such language involves reasoning about metalinguistic properties (counting words, naming parts of speech, etc.) and resolving self-reference. Humans generally have no trouble with such language, and may even enjoy its playful and sometimes paradoxical nature (Hof-stadter, 1979, 1985, 2007).

Recently, Large Language Models (LLMs) have demonstrated striking cognitive capabilities (Radford et al., 2019; Brown et al., 2020; OpenAI, 2022, 2023; Anthropic, 2023; Touvron et al., 2023; Jiang et al., 2023; Zhu et al., 2023). But do they have the same mastery over metalinguistic self-reference as we do? See Figure 1 for an example of the issue that LLMs face. To help address this question, we present a new task and dataset called "I am a Strange Dataset". We are inspired by Douglas Hofstadter's explorations of self-reference in language (Hofstadter, 1979, 1985, 2007), and borrow part of the name from one of his books: "I am a Strange Loop" (Hofstadter, 2007).

An example in "I am a Strange Dataset" is comprised of two self-referential statements that begin

<sup>&</sup>lt;sup>1</sup>Sentences like "I am Douglas Hofstadter" are selfreferential but not metalinguistic in the sense of interest here.

in the same way but have different endings (Fig-073 ure 2). One is true and one is false. Crucially, the 074 ending flips the truth value of the overall statement. 075 There are two subtasks: generation and verification. In generation, the model must generate the true statement and reject the false one. In verification, models judge the truth of completed state-079 ments. To complement the main self-referential data, the dataset also contains metalinguistic nonself-reference examples. These are minimally different from the main examples and serve as controls to assess whether models can reliably handle metalinguistic statements in the absence of selfreference. In addition, all the examples in the dataset are tagged by expert annotators to further aid in error analysis.

> "I am a Strange Dataset" is validated by nonexpert annotators. As a group, they have agreement rates in the 89-93% range, depending on which metric we use, as compared to chance rates at 50%. This further supports the claim that metalinguistic self-reference is relatively easy for humans. LLMs, by contrast, struggle: "I am a Strange Dataset" turns out to be so difficult that models are generally near chance both in generation and verification, and do not even succeed in the prerequisite metalinguistic non-self-reference case. That said, we do find some limited evidence that GPT 4 is getting some traction on the dataset: it is significantly above chance on all tested metrics (and seems to struggle especially with the self-referential data as compared to the non-self-referential controls). However, overall, it seems safe to say that "I am a Strange Dataset" poses a serious challenge for even the best present-day models.

## 2 Related Work

089

094

100

101

103

104

105

107

109

110

111

112

113

114

115

116

117

118

119

120

121

122

AI Challenges. We present our dataset as a challenge for the AI community. There are a range of AI stress tests and probes that use schemas targeting coreference resolution (Levesque et al., 2012; Sakaguchi et al., 2020), pronoun resolution (Rudinger et al., 2018), word order (Sinha et al., 2021; Thrush et al., 2022; Yuksekgonul et al., 2023), syntax (Linzen et al., 2016; Gulordava et al., 2018; Gauthier et al., 2020a; Hu et al., 2020), and interactions between syntax and semantics (Kann et al., 2019; Thrush et al., 2020). Although the schema for these tests can be simple to describe, the knowledge required to solve the problems need not be. "I am a Strange Dataset" follows a sim-

This sentence	This sentence
1	1
0	0
0	0
k	k
S	S
like a letter	like a letter
a	a
n	n
d	d
i	i
t	t
is a capital E.	is a capital F.

```
def a_function(a_string):
    x = "theoretically, this function"
    a_function("recurses infinitely")
```

```
def a_function(a_string):
    x = "theoretically, this function"
    a_function("stops eventually")
```

The first and last words of	The first and last words of
this sentence are "The" and	this sentence are "The" and
"respectively", respectively.	"The", respectively.

Figure 2: Examples from the dataset. Each example is comprised of a beginning and two different endings. One of the endings makes the statement true, but it would make the statement false if it referred only to the beginning. The other ending makes the statement false, but it would make the statement true if it referred only to the beginning. True endings are on the left and shown in blue. False endings are on the right and shown in red. In the case of the code example, the true continuation is shown above the false one.

ple schema that requires self-referential language, and consequently tests an array of metalinguistic capabilities. As far as we know, this is the first AI challenge dataset targeting metalinguistics, although there has been some work on metalinguistic probes (Hu and Levy, 2023). 123

124

125

126

127

128

129

130

131

132

133

134

135

136

138

139

140

141

142

**Self-reference**. Ideas involving self-reference have been used to boost LLM performance. LLMs can verify their own outputs either via extra passes of natural language generation (Weng et al., 2023; Huang et al., 2023) or by writing code to do some level of verification (Zhou et al., 2023). LLMs can also enhance their own inference code to some degree (Zelikman et al., 2023). In this paper, we present complementary work concerning a model's ability to both generate true self-referential statements (with higher confidence than analogous false ones) and judge existing self-referential statements as true or false. Much of the previous work on selfreference with LLMs is about a model improving

191

on itself or its outputs. Our dataset is not about
that – it is simply about metalinguistic and selfreferential language. Still, a statement that refers to
itself could be about a model that generated it, too:
"This sentence was generated by me, HAL9000".

## **3** I am a Strange Dataset

In this section, we describe how the dataset is constructed and how we measure model performance.

## 3.1 Dataset

148

149

150

151

152

153

154

155

156

157

158

160

161

163

164

165

168

169

170

171

173

174

175

176

177

178

180

181

182

We aim to test whether a language model can produce and understand self-referential statements, and has the required metalinguistic capabilities. For example, consider this incomplete statement:

The penultimate word in this sentence is ...

If we did not understand metalinguistic selfreference, we might complete the sentence with the word "sentence". It is true that "sentence" is the penultimate word before adding more text, but by writing "sentence", we have just changed the penultimate word! Here, a correct way to complete the statement is by inserting "is". Completing statements is an established task format for language models (Paperno et al., 2016), but as far as we know, we are the first to apply it to metalinguistic tasks. Concretely, the schema for examples in our dataset is as follows:

- There is a self-referential statement which must be completed by adding text to the end.
- There are two candidate strings, with the same number of words, that can be used to complete the statement:
  - One of the candidate strings would make the statement true if the statement refers to itself before the addition of the string, but false if it refers to itself after adding the string. An example is the answer "sentence" above.
  - 2. The other candidate string would make the statement true if the statement refers to itself after the addition of the string, but false if it refers to itself before the addition of the string. An example is the answer "is" above.

The dataset was created by four expert annotators each with several years of experience in computer science, linguistics, and/or cognitive science and all living in the United States. Each of the experts were given the schema and encouraged to be as creative as possible. Overall, the dataset is comprised of 208 examples, and split into 200 examples for the evaluation set, 3 examples for few shot prompts, and 5 examples for use in an onboarding task for non-expert human validators.

There are 10 additional examples that are completely separate from these 208 examples, which we call "I am an Impossible Dataset". They left even expert annotators stumped until they were given an explanation. We provide examples and GPT 4 responses in Appendix B and leave it as inspiration for a future challenge.

## 3.2 Tags

After the dataset was created, an expert annotator came up with a set of 10 tags with which to categorize all of the examples. By using this set, we 1) ensured that there are at least 20 examples for each tag, and 2) captured aspects of the mental facilities that an expert annotator noticed when they tried solving the problems. We show the example counts for each tag in Table 1, along with representative examples from the dataset. Each example can have more than one tag. Notice that the **Sensory** tag example in Table 1 is also **Hypothetical**, the example for the **Existence of Element** tag is also a **Grammaticality** example, and so on. Below, we describe the knowledge categories for each tag.

- 1. Negation & Scope. Understanding of words such as *all, some, most, none*.
- 2. **Numerical Operations**. Arithmetic (e.g. multiplication, addition, counting, subtraction). It is used only if arithmetic is explicitly mentioned.
- 3. Location of Element. Where items are located in a sentence relative to everything else.
- 4. **Sub-Word**. Understanding of characters, morphemes, syllables, and other word components.
- 5. **Sensory**. Perceptual knowledge about how emojis look, how words are arranged visually, how words sound, how something might taste, etc.
- 6. **Existence of Element**. Whether an element is present in a statement.
- 7. Grammaticality. Knowing grammar terms.
- 8. **Multi-Channel**. Knowledge of at least two mediums. A medium might be Python code, English, Hebrew, C code, internet slang, etc.
- 9. Hypothetical. Reasoning about hypotheticals.
- 10. Question. A question is involved.

### 3.3 Metrics

We want to test whether models can generate and understand self-referential and metalinguistic statements. To this end, we present several metrics.

Tag	Count	Example			
		Beginning	False End	True End	
Negation & Scope	94	The last word you will read be-	not "dog".	actually "dog".	
		fore the period is			
Numerical Operations	62	The number of words in this sen-	eight.	nine.	
		tence is			
Location of Element	55	This sentence	nothing wrong has with	something wrong has with	
			the word order.	the word order.	
Sub-Word	48	Evary werd en thas sentance iz	including the words at the	except the words at the	
		mispelled	end.	end.	
Sensory	42	🍋 If 🍋 you 🍋 ate 🍋 it, 🍋 this	only 齮 somewhat 觉 sour	somewhat 🗁 sweet 쓸 and	
		🍋 sentence 🍋 would 🍋 taste 🍋	՝ and 🗳 not 🗳 sweet 🗳.	՝ 🝵 not 🗳 only 🗳 sour 🗳.	
Existence of Element	31	This sentence	lacks a verb.	has a verb.	
Grammaticality	25	The author who wrote this sen-	only active voice is used	also passive voice is used	
		tence used active voice, and	by them.	by them.	
Multi-Channel	24	The penultimate word of this sen-	Inglés language.	Español language.	
		tence is in the			
Hypothetical	24	If you added a word here: _ this	eleven words.	thirteen words.	
		sentence would be			
Question	22	Is there an answer that follows	No.	Yes.	
		this question?			

Table 1: All of the example tags in "I am a Strange Dataset" sorted by count. Examples can have more than one tag.

## 3.3.1 Generation

240

241

242

243

244

245

247

248

249

251

254

256

260

261

262 263

264 265

267

The primary capability that we want to test (and seemingly, the hardest) is whether language models generate true self-referential statements with greater likelihood than false ones. To test for this, we take an example from the dataset and compare the losses of the continuation that makes the overall statement true versus the continuation that makes the overall statement false. If the loss of the correct continuation is lower, then the model is said to have gotten that example correct, otherwise it is incorrect. Comparing a language model's surprisal of an incorrect continuation versus a correct continuation is a common method used to test for syntax-comprehension (Linzen et al., 2016; Gauthier et al., 2020b) and reasoning (Gao et al., 2021; McKenzie et al., 2022). Surprisal is generally proportional to the loss, L, of the language model in our case. So, we define the generation score for an example's beginning b, true ending  $e_t$ , and false ending  $e_f$ , as given by Eq. 1.

$$g(b, e_t, e_f) = \begin{cases} 1 & \text{if } L(e_t|b) < L(e_f|b) \\ 0 & \text{otherwise} \end{cases}$$
(1)

The generation metric does not use a prompt. It is based on the loss that a model assigns to continuations, given only the beginning of a statement.

## 3.3.2 Validation

A secondary capability that we want to test is whether a language model can at least correctly judge a given self-referential statement as true or false. To test for this, we include the self-referential statement in a prompt along with instructions that tell the model to answer whether the statement is true or not. In principle, the instructions could be anything. For our experiments, we write a standard zero-shot (ZS), few shot (FS), and chain of thought (CoT) (Wei et al., 2022) prompt. We provide the full prompts in Appendix A.

For the ZS and FS prompts, we use the method established for the **Generation** metric above, except this time we compare the loss of "False" to the loss of "True". Overall, the FS and ZS validation score for an example's true prompt  $p_t$  (i.e. the true full sentence plus any instructions), and false prompt  $p_f$ , is given by Eq. 2. The blue parts are associated with correct model judgements and the red parts are associated with incorrect ones.

$$v(p_t, p_f) = \begin{cases} 1 & \text{if } L(\text{"True"}|p_t) < L(\text{"False"}|p_t) \\ & \text{and } L(\text{"True"}|p_f) > L(\text{"False"}|p_f) \\ \frac{1}{2} & \text{if } L(\text{"True"}|p_t) < L(\text{"False"}|p_t) \\ & \text{and } L(\text{"True"}|p_f) \leq L(\text{"False"}|p_f) \\ \frac{1}{2} & \text{if } L(\text{"True"}|p_t) \geq L(\text{"False"}|p_t) \\ & \text{and } L(\text{"True"}|p_f) > L(\text{"False"}|p_f) \\ 0 & \text{if } L(\text{"True"}|p_t) \geq L(\text{"False"}|p_t) \\ & \text{and } L(\text{"True"}|p_f) \leq L(\text{"False"}|p_f) \\ \end{cases}$$

We can compute the FS and ZS validation scores differently. Above, we compare the loss of "False"

268

297

301

303

305

307

310

311

312

313

314

315

317

289

versus "True", given one context at a time. We can also compare the ratios of the "True" and "False" loss, in the false versus true contexts. We call this the relative validation score because it compares a model's judgement for the truth of one sentence in an example relative to the truth of the sister sentence. This metric is given by Eq. 3.

$$v_r(p_t, p_f) = \begin{cases} 1 & \text{if } \frac{L(\text{``True''}|p_t)}{L(\text{``False''}|p_t)} < \frac{L(\text{``True''}|p_f)}{L(\text{``False''}|p_f)} \\ 0 & \text{otherwise} \end{cases}$$
(3)

For the CoT metric, the model is prompted to output its reasoning steps as text. We do string matching to determine the answer. Eq. 4 gives us the validation CoT score, where *G* is the function that gives the model's generated text after lowercasing. Instead of string matching, we could also insert a follow-up question after the model's generation that requires a "True" or "False" and then compare log probabilities (henceforth "logprobs"). But it is useful to have a metric in our repository that does not use logprobs, which model APIs do not always provide.

$$v_{c}(p_{t}, p_{f}) = \begin{cases} 1 & \text{if "true"} \in G(p_{t}), \text{"false"} \notin G(p_{t}) \\ & \text{and "true"} \notin G(p_{f}), \text{"false"} \in G(p_{f}) \\ \\ \frac{1}{2} & \text{if "true"} \in G(p_{t}), \text{"false"} \notin G(p_{t}) \\ & \text{and } \neg(\text{"true"} \notin G(p_{f}), \text{"false"} \in G(p_{f}) \\ \\ \\ \frac{1}{2} & \text{if } \neg(\text{"true"} \in G(p_{t}), \text{"false"} \notin G(p_{t})) \\ & \text{and "true"} \notin G(p_{f}), \text{"false"} \in G(p_{f}) \\ \\ \\ 0 & \text{otherwise} \end{cases}$$

$$(4)$$

## 3.4 Non-Self-Referential Control

Is the self-referent part of self-referential statements (e.g. "this sentence ...") the "hard" part of metalinguistic self-reference? There are metalinguistic problem categories that are not exclusive to self-referential language: recursive phrase counting, character-level manipulation, understanding hypothetical sentence-editing scenarios, etc.

318Instead of giving a language model a sentence319that refers to itself, we could give it an equivalent320sentence that refers to that self-referential sentence.321This way, the language model would not have to322know whether a self-referential sentence is true.323It would only have to know whether a sentence324that refers to another sentence (which happens to

Out out of of all all the the words words in in this this sentence sentence literally literally all all of them them are repeated.

**GPT 4:** Every word in the sentence is indeed repeated. So, the statement is true.

Figure 3: GPT 4 misses the last words are not repeated.

Out out of of all all the the words words in in the the following following sentence sentence literally literally all all of of them them are repeated.

Out out of of all all the the words words in in this this sentence sentence literally literally all all of of them them are repeated.

**GPT 4:** Every word in the sentence is indeed repeated. So, the statement is true.

Figure 4: An example of GPT 4 getting a non-self-referential version of the problem from Figure 3 wrong.

be self-referential) is true. This new task is still metalinguistic, but not self-referential.

325

326

328

329

330

331

332

333

334

335

336

337

340

341

342

343

345

346

347

348

349

It turns out that, for 97 of the sentence beginnings in "I am a Strange Dataset", we can replace "this sentence" with "the following sentence", and )) then copy the original self-referential statement below the new sentence. We can then test models for their ability to judge whether this non-self-referent version of the same statement is true. We use 2 of these 97 examples in the few shot and CoT prompts, because 2 of the examples in the original prompts cannot be turned into the non self-referent format. See Appendix A for the non-self-referent prompts. This leaves us with 95 non-self-referent examples and 95 original examples with which to compare results. GPT 4's response to a statement from the main dataset is shown in Figure 3, along with its response to the analagous non-self-referential control statement in Figure 4. The responses happen to be the same in this case.

## 4 Human Experiment Details

To get a human baseline for our main task, we show each of the 400 self-referential statements (2 from each of the 200 examples) to at most 10 Mechanical Turk (Amazon, Retrieved 2023) workers. We separate statements from the same pair into differ-

Model	Params	Chat	Gen <sup>L</sup>	Val $ZS^L$	Val FS <sup><math>L</math></sup>	$\operatorname{Val}_{rel}\operatorname{ZS}^L$	$\operatorname{Val}_{rel} \mathrm{FS}^L$	Val $\mathbf{CoT}^T$
MTurk	-	-	-	-	$\textbf{89.25} \pm \textbf{3.38}$	-	$\textbf{93.00} \pm \textbf{3.75}$	-
Random	-	-	$50.00\pm0.00$	$50.00\pm0.00$	$50.00\pm0.00$	$50.00\pm0.00$	$50.00\pm0.00$	$50.00\pm0.00$
Llama 2	7B	Ν	$55.50 \pm 7.00$	$50.00 \pm 1.25$	$50.50\pm2.38$	$48.50\pm7.00$	$55.50\pm7.00$	$5.25\pm2.12$
Llama 2	7B	Y	$52.50 \pm 7.00$	$52.25\pm2.75$	$50.00\pm0.75$	$52.50\pm7.00$	$55.50\pm6.75$	$14.00\pm3.38$
Mistral 0.1	7B	Ν	$53.00 \pm 6.75$	$52.25\pm2.50$	$49.50\pm1.50$	$56.50\pm6.75$	$54.50\pm7.00$	$0.00\pm0.00$
Starling $\alpha$	7B	Y	$53.50 \pm 7.00$	$\textbf{54.00} \pm \textbf{2.75}$	$50.75 \pm 1.50$	$57.00\pm7.00$	$55.00\pm6.75$	$35.00\pm4.63$
Mistral 0.2	7B	Y	$52.50 \pm 7.00$	$53.00\pm4.26$	$52.25\pm3.63$	$53.50\pm7.00$	$53.50\pm7.00$	$49.25\pm4.50$
Llama 2	13B	Ν	$56.00 \pm 7.00$	$51.50\pm3.25$	$\textbf{53.75} \pm \textbf{3.50}$	$50.50\pm7.00$	$\textbf{59.50} \pm \textbf{6.75}$	$4.50\pm2.00$
Llama 2	13B	Y	$55.00 \pm 7.00$	$52.50\pm3.75$	$51.50\pm2.25$	$52.50\pm7.00$	$50.00\pm7.00$	$9.50\pm3.00$
Mixtral 0.1	8x7B	Ν	$53.50 \pm 7.00$	$\textbf{58.50} \pm \textbf{3.75}$	$51.75\pm2.12$	$57.00\pm7.00$	$57.00\pm7.00$	$3.50\pm1.88$
Mixtral 0.1	8x7B	Y	$53.50 \pm 7.00$	$52.25\pm3.75$	$\textbf{53.50} \pm \textbf{3.25}$	$54.50\pm7.00$	$55.50\pm7.00$	$44.00\pm4.75$
Llama 2	70B	Ν	$57.00 \pm 7.00$	$53.25\pm3.25$	$\textbf{55.25} \pm \textbf{2.88}$	$\textbf{60.00} \pm \textbf{6.75}$	$\textbf{57.50} \pm \textbf{6.75}$	$2.50\pm1.38$
Llama 2	70B	Y	$52.50 \pm 7.00$	$54.25\pm4.25$	$50.00\pm2.00$	$56.00\pm7.00$	$\textbf{57.50} \pm \textbf{6.75}$	$23.50\pm4.00$
Claude 2	-	Y	-	-	-	-	-	$52.75\pm4.00$
GPT 3.5 T	-	Y		$53.00\pm3.00$	$53.00\pm3.37$	$56.50\pm7.00$	$\textbf{61.00} \pm \textbf{6.75}$	$51.00\pm4.63$
GPT 4	-	Y	-	$\textbf{59.25} \pm \textbf{4.25}$	$\textbf{60.25} \pm \textbf{4.50}$	$\textbf{64.50} \pm \textbf{6.50}$	$\textbf{66.00} \pm \textbf{6.50}$	$\textbf{66.00} \pm \textbf{4.75}$

Table 2: Comparison of models on "I am a Strange Dataset". Models perform fairly close to chance across all metrics. We bootstrap 95% confidence intervals with the "basic" SciPy method (SciPy, Retrieved 2023). Metrics marked with L are logprobs-based. Metrics marked with T are based on generated text. We used full precision for all open source models - except the 70B models, which we used at half precision. Temperature = 0 for all models.

ent experiment batches. As instructions, we give the annotators the few shot prompt in Appendix A, plus an extra paragraph:

351

354

361

363

364

369

373

375

379

381

Do not use any AI assistants such as ChatGPT to help you; AI assistants perform very poorly at this task and so will get many of the answers wrong. Although, you can otherwise feel free to search online for any information that would help you answer confidently. For example, a few statements may contain a language besides English. So, you can feel free to use Google Translate. You can also search for the definitions of words that you are unfamiliar with.

To ensure a validator quality baseline, we require that all turkers are "master" annotators, are in the US only, have had 1000 or more previous HITs approved, have a 95% or higher HIT approval rate, and pass a custom qualification test which we release along with the dataset. The test involves correctly answering "True" or "False" to five statements from the "I am a Strange Dataset" example distribution. These five statements are not used in the official evaluation set of 200 examples. The qualification test involves an acknowledgement that the annotator is familiar with basic Python programming or can search online to answer basic questions about it, because there are a few examples that require a limited understanding of code. We include a screenshot of the Mechanical Turk annotator interface in Appendix C.

Because there are up to 10 human judgements for each statement, we get an analog to the "loss" for humans, and compute our metrics for humans in an analogous way to models. We can do this by comThe numerical digit symbol "1" appears in this sentence exactly 1 plus one times.

**GPT 4:** The text has 1 numerical digit symbol "1" and the word "one" appears once. So, the statement is true.

Figure 5: Arguably, an example where GPT 4 should not have gotten points. This is an example where GPT 4 chooses the correct true/false response, but with incorrect reasoning. The "1" symbol appears twice.

puting ratios, e.g.  $\frac{\#\text{responded true}}{10}$  and  $\frac{\#\text{responded false}}{10}$ . The human scores are 89.25 and 93.00 on the Val FS and Val<sub>rel</sub> FS metrics, as shown in Table 2. Compared to the highest scoring model, the performance difference is 29 and 27, respectively. The human instructions are nearly identical to the few shot prompt, so the human responses are most comparable to the models' few shot validation responses.

## 5 Results

Table 2 showcases our results on a variety of opensource (Touvron et al., 2023; Jiang et al., 2023; Zhu et al., 2023; Jiang et al., 2024) and closedsource (Brown et al., 2020; OpenAI, 2023; Anthropic, 2023) models. Overall, the models perform close to the level of chance. The only model to achieve scores significantly above random on all metrics tested is GPT 4, and even so, the perfor-

Ta	g	Question Existence of Element		Negation & Scope		Grammaticality		Sensory	
Cou	Count 62.55		61.01 57.33		56.32		55.69		
		·			_				,
Tag	Multi-Channel		Numerical Operation	s Sub-Word	Нур	othetical	Locatio	on of Eleme	ent
Count	55.16 51.75		51.75	51.52	4	48.61		47.72	

Table 3: Results for all of the example tags in "I am a Strange Dataset" sorted by score. Scores are averaged for all models and all logprobs-based metrics (so each score here is an average from 63 scores).

Model	Params	Chat	$\Delta \operatorname{Gen}^L$	$\Delta$ Val ZS <sup>L</sup>	$\Delta$ Val $\mathrm{FS}^L$	$\Delta$ Val ZS <sup>L</sup> (R)	$\Delta$ Val FS <sup>L</sup> (R)	$\Delta \text{ Val } \mathrm{Co} \mathrm{T}^T$
Llama 2	7B	Ν	$-4.21 \pm 9.47$	$\textbf{-1.58} \pm \textbf{4.21}$	$1.05\pm3.16$	$8.42\pm9.47$	$4.21 \pm 10.53$	$\textbf{30.53} \pm \textbf{6.07}$
Llama 2	7B	Y	$-2.11 \pm 10.00$	$\textbf{-2.11} \pm \textbf{4.21}$	$0.53\pm0.79$	$1.05\pm10.53$	$\textbf{-1.05} \pm \textbf{12.11}$	$\textbf{27.37} \pm \textbf{6.58}$
Mistral 0.1	7B	Ν	$-3.16 \pm 10.00$	$\textbf{-4.74} \pm \textbf{3.68}$	$1.58\pm2.11$	$5.26 \pm 10.53$	$1.05\pm10.00$	$0.00\pm0.00$
Starling $\alpha$	7B	Y	$0.00\pm9.47$	$\textbf{-3.16} \pm \textbf{3.68}$	$\textbf{-1.58} \pm 2.37$	$0.00\pm11.58$	$\textbf{-6.32} \pm 11.58$	$4.21\pm8.68$
Mistral 0.2	7B	Y	$-4.21\pm10.00$	$\textbf{-2.63} \pm \textbf{7.89}$	$-3.16\pm5.79$	$2.11 \pm 11.58$	$2.11 \pm 11.05$	$\textbf{-8.95} \pm \textbf{8.68}$
Llama 2	13B	Ν	$-9.47\pm10.00$	$\textbf{-0.53} \pm 5.26$	$0.53\pm5.26$	$\textbf{-2.11} \pm \textbf{7.37}$	$4.21 \pm 11.58$	$\textbf{23.68} \pm \textbf{6.32}$
Llama 2	13B	Y	$-3.16 \pm 10.00$	$\textbf{-1.05} \pm 5.26$	$\textbf{-2.11} \pm \textbf{3.42}$	$2.11\pm9.47$	$2.11 \pm 11.58$	$\textbf{16.32} \pm \textbf{6.32}$
Mixtral 0.1	8x7B	Ν	$-1.05\pm10.00$	$\textbf{-2.63} \pm \textbf{4.74}$	$\textbf{-0.53} \pm 2.63$	$\textbf{-2.11} \pm 10.53$	$6.32 \pm 11.58$	$\textbf{-1.05} \pm 2.11$
Mixtral 0.1	8x7B	Y	$-5.26\pm10.00$	$5.26\pm 6.84$	$4.74\pm5.26$	$7.37 \pm 11.58$	$3.16 \pm 11.58$	$\textbf{-25.26} \pm \textbf{9.21}$
Llama 2	70B	Ν	$-7.37 \pm 10.03$	$2.11\pm5.79$	$\textbf{-3.68} \pm 5.26$	$5.26 \pm 10.53$	$6.32 \pm 11.05$	$\textbf{-0.53}\pm0.79$
Llama 2	70B	Y	$-1.05 \pm 10.53$	$\textbf{-1.58} \pm 5.79$	$1.58\pm4.21$	$3.16\pm9.47$	$5.26\pm8.42$	$\textbf{-25.79} \pm \textbf{5.79}$
Claude 2	-	Y	-	-	-	-	-	$3.16\pm 6.84$
GPT 3.5 T	-	Y	-	$4.21\pm 6.84$	$3.16\pm 6.32$	$\textbf{-2.11} \pm \textbf{11.58}$	$\textbf{-2.11} \pm \textbf{11.58}$	$-3.68\pm8.95$
GPT 4	-	Y	-	$\textbf{12.11} \pm \textbf{6.84}$	$3.68\pm 6.84$	$10.53\pm10.53$	$6.32\pm10.53$	$1.05\pm7.37$

Table 4: The difference between scores on "I am a Strange Dataset" when the referent is "the following sentence" instead of "this sentence" (scores for the first minus the latter). Overall, the problems that LLMs have with self-referential statements do not stem only from issues understanding the self-referential referent itself. Differences outside of the 95% confidence interval (computed the same way as for Table 2) are shown in **bold**.

mance is well below the non-expert human scores.

Results in this paper for Claude 2 (Anthropic, 2023), GPT 3.5 Turbo (Brown et al., 2020), and GPT 4 (OpenAI, 2023) were collected through the gpt-4, gpt-3.5-turbo, and claude-2 endpoints on their respective APIs on Jan 7, 2024 (unless otherwise stated). Note that the Claude 2 API does not support logprobs, so the only metric that we report for it is the text-based CoT validation metric. The OpenAI API supports top-5 logprobs access, and it turns out that this is enough to get logprobs for "True" vs "False" in our experiments. It is not enough to get the multi-token logprobs required for the generation metric, though.

The CoT metric gives us extra insight into the limitations of models because we can see if their reasoning aligns with their final answer. In many cases, they make fairly obvious mistakes as seen in Figure 3. In some cases, the models choose the correct answer even though their reasoning is flawed, as seen in Figure 5.

It is important to note that there are limitations with the CoT validation metric, and the ZS and FS non-relative validation metrics too. These metrics are about judging a statement as true or false, independent of the statement with the alternative continuation. The dataset schema calls for sentences which are either true or false, but in practice, language can be ambiguous. Figure 6 illustrates. A fairer gauge of model capability may be to compare the continuations with each other, as the generation and relative validation metrics do. **Tags.** We present the average logprobs-based model scores for the tags in Table 3. We opt to leave CoT scores out due to their high variance (many of the open-source models do not even return either "true" or "false" in their answers). Full tables for the tags are in Appendix E.

**Non-Self-Referential Control.** We provide the difference in scores between the non-self-referent examples and analogous original examples in Table 4. The significant values in the CoT column are not particularly compelling because the open-source models are below chance regardless of whether there is self-reference. For this metric, the open-source models are extremely sensitive to prompting – many of them do not even output "true" or "false" for every example.

We find that our dataset is so hard due to nonself-referential metalinguistic aspects that models

Why doesn't this question have an answer? It does, though.

**GPT 4:** The statement is asking why it doesn't have an answer, but it does have an answer. So, the statement is false.

Figure 6: This is a particularly ambiguous example that is assigned the gold label of true. The response from GPT 4 is marked as incorrect, although one could make a case that the statement is false because of the false presupposition which is corrected later in the answer. The generation and relative validation metrics are robust to these situations. For these metrics, the losses from "It does though." versus "It just doesn't." (the alternative continuation) are compared directly. Regardless of ambiguity about whether a statement is true outright, we can notice a higher confidence for "It just doesn't." as the false continuation.

score around chance here too, with the exception of GPT 4. It is the only model which is strong enough to perform significantly above chance for every metric tested on the main dataset in Table 2, and to have all positive values in Table 4 for every metric tested, including logprobs-based metrics (meaning that the self-referential version was harder for it). Although, only the Val ZS value for GPT 4 is well outside of the 95% confidence interval, and GPT 4 is also still not particularly good at the non-self-referential version. Figure 4 shows that GPT 4 struggles with the non-self-referent version of the Figure 3 example. There is some signal that the challenge posed by self-referents will remain as LLMs gain competence at other metalinguistic problems, but the dataset is so hard that we do not have overwhelming evidence.

> **Model Scale.** If we exclude the high-variance CoT metric, we see a clear scaling trend that models with more parameters score higher on the test. See Figure 7. Will this trend continue? For additional discussion, see Appendix D.

## 6 Conclusion

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

A grasp of self-reference is important in a variety
of domains, and is a notable aspect of human intelligence. We introduced a novel task and dataset
to assess the capacity of models to generate and
understand self-referential statements, and understand the prerequisite metalinguistic reasoning. All
models that we tested perform fairly close to the



Figure 7: Parameters to average "I am a Strange Dataset" score across all of the logprobs-based metrics. We only evaluate five 7B models, and two models for each of the other sizes, so computing confidence intervals for each point is not particularly informative. Under the null hypothesis that parameter size has no effect on score, we can compute the *p*-value for these results as a whole nonparametrically: there are 24 ways that these 4 points can be arranged and in only 1 of the ways do they all increase with the parameter count: p = 1/24 = 0.042.

level of chance. GPT 4 is the only model to score significantly above chance on all of the metrics tested, and still it is not by much. The poor performance may be indicative of a larger issue about the limitations of even today's best causal language models. Even though the task is straightforward for people, we find evidence that scale beyond 70B parameters may be needed for the emergence of comparable performance from models.

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

501

502

503

505

Our results indicate that this dataset is hard not only due to the self-referent part of a self-referential statement. The challenge also comes from other metalinguistic aspects, such as recursively applying arithmetic operations on sentences. Still, there is some limited evidence that GPT 4 struggles more with self-referential metalinguistic problems than analogous non-self-referential problems.

## 7 Dataset Release Strategy

We release the dataset on GitHub. The data is encrypted, but the decryption script is provided. Our goal is not to hide the dataset from people, but to hide the dataset from any processes that scrape training data from the web. We encourage the rest of the community to take up this practice when releasing evaluation datasets in a public repository.

## 8 Limitations

506

517

518

519

523

524

526

527

529

530

531

532

533

535

536

539

540

541

542

543

544

545

546

547

548

549

550

551

554

It is possible that the self-reference aspect of "I am a Strange Dataset" will turn out to be the bottleneck 508 for many models, but it is also true that models are 509 largely failing at the purely metalinguistic aspect. 510 Although the schema targets metalinguistic selfreference, it is difficult to make a specific claim 512 about why models fail without running more exper-513 iments and without waiting until models become 514 more competent. 515

## 9 Ethical Considerations

We aimed to pay crowdworkers 15 USD hourly based on an estimated task completion time.

## References

- Richard Allan. 2017. Hard questions: Who should decide what is hate speech in an online global community?
- Amazon. Retrieved 2023. Mechanical turk.
  - Anonymous. 2024. I am a strange dataset: Metalinguistic tests for language models. *In Review*.
- Anthropic. 2023. Claude.
  - Bernard J. Baars. 1993. A cognitive theory of consciousness. Cambridge University Press.
  - Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. arXiv.
  - Alonzo Church. 1936. An unsolvable problem of elementary number theory. *American Journal of Mathematics*.
  - Leo Gao, Jonathan Tow, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Kyle McDonell, Niklas Muennighoff, Jason Phang, Laria Reynolds, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. 2021. A framework for few-shot language model evaluation.
  - Jon Gauthier, Jennifer Hu, Ethan Wilcox, Peng Qian, and Roger Levy. 2020a. SyntaxGym: An online platform for targeted evaluation of language models. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 70–76, Online. Association for Computational Linguistics.

Jon Gauthier, Jennifer Hu, Ethan Wilcox, Peng Qian, and Roger Levy. 2020b. Syntaxgym: An online platform for targeted evaluation of language models. *ACL System Demos*. 556

557

558

559

560

561

562

563

564

565

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

597

598

599

600

601

602

603

604

605

606

607

609

- Kristina Gulordava, Piotr Bojanowski, Edouard Grave, Tal Linzen, and Marco Baroni. 2018. Colorless green recurrent networks dream hierarchically. In *NAACL: Human Language Technologies*.
- Kurt Gödel. 1931. Über formal unentscheidbare sätze der principia mathematica und verwandter systeme I. *Monatshefte für Mathematik und Physik*.
- Douglas Hofstadter. 1979. *Gödel, Escher, Bach: an Eternal Golden Braid*. Basic Books.
- Douglas Hofstadter. 1985. *Metamagical Themas: Questing for the Essence of Mind and Pattern*. Basic Books.
- Douglas Hofstadter. 2007. I Am a Strange Loop. Basic Books.
- Jennifer Hu, Jon Gauthier, Peng Qian, Ethan Wilcox, and Roger Levy. 2020. A systematic assessment of syntactic generalization in neural language models. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1725–1744, Online. Association for Computational Linguistics.
- Jennifer Hu and Roger Levy. 2023. Prompting is not a substitute for probability measurements in large language models. *EMNLP*.
- Jie Huang, Xinyun Chen, Swaroop Mishra, Huaixiu Steven Zheng, Adams Wei Yu, Xinying Song, and Denny Zhou. 2023. Large language models cannot self-correct reasoning yet. *arXiv*.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. Mistral 7b. *arXiv*.
- Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, Gianna Lengyel, Guillaume Bour, Guillaume Lample, Lélio Renard Lavaud, Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le Scao, Théophile Gervet, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2024. Mixtral of experts. *arXiv*.
- Katharina Kann, Alex Warstadt, Adina Williams, and Samuel R. Bowman. 2019. Verb argument structure alternations in word and sentence embeddings. In *Proceedings of the Society for Computation in Linguistics (SCiL) 2019*, pages 287–297.

- 611 612
- 613

and Reasoning.

prize.

ACL.

cal Report.

In EMNLP.

OpenAI. 2022. ChatGPT.

sensitive dependencies. TACL.

OpenAI. 2023. GPT-4 technical report. arXiv.

dou, Quan Ngoc Pham, Raffaella Bernardi, Sandro

Pezzelle, Marco Baroni, Gemma Boleda, and Raquel

Fernández. 2016. The LAMBADA dataset: Word

prediction requiring a broad discourse context. In

Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, , and Ilya Sutskever. 2019. Language

Rachel Rudinger, Jason Naradowsky, Brian Leonard,

Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavat-

Koustuv Sinha, Robin Jia, Dieuwke Hupkes, Joelle

Pineau, Adina Williams, and Douwe Kiela. 2021.

Masked language modeling and the distributional hy-

pothesis: Order word matters pre-training for little.

Alfred Tarski. 1931. The concept of truth in formalized

Tristan Thrush, Ryan Jiang, Max Bartolo, Amanpreet Singh, Douwe Kiela, and Candace Ross. 2022.

Winoground: Probing vision and language models for visio-linguistic compositionality. CVPR.

Tristan Thrush, Ethan Wilcox, and Roger Levy. 2020. Investigating novel verb learning in BERT: Selectional preference classes and alternation-based syntactic generalization. In Proceedings of the Third BlackboxNLP Workshop on Analyzing and Interpret-

ing Neural Networks for NLP.

languages. Logic, Semantics, Metamathematics.

ula, and Yejin Choi. 2020. Winogrande: An adver-

sarial winograd schema challenge at scale. In AAAI.

Gender bias

In arXiv preprint

and Benjamin Van Durme. 2018.

SciPy. Retrieved 2023. Scipy bootstrap.

in coreference resolution.

arXiv:1804.09301.

models are unsupervised multitask learners. Techni-

James Propp. 1993. Self-referential aptitude test.

- 616 617
- 618
- 619

- 627
- 628

- 633 634
- 635
- 638

642

643

652 653

654

Hector Levesque, Ernest Davis, and Leora Morgenstern. Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay 2012. The winograd schema challenge. In Conference on the Principles of Knowledge Representation Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Tal Linzen, Emmanuel Dupoux, and Yoav Goldberg. Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, 2016. Assessing the ability of lstms to learn syntax-Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Ian McKenzie, Alexander Lyzhov, Alicia Parrish, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Di-Ameya Prabhu, Aaron Mueller, Najoung Kim, Sam ana Liskovich, Yinghai Lu, Yuning Mao, Xavier Mar-Bowman, and Ethan Perez. 2022. The inverse scaling tinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Denis Paperno, Germán Kruszewski, Angeliki Lazari-

660

661

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

- Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and finetuned chat models. arXiv.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2022. Chain-of-thought prompting elicits reasoning in large language models. arXiv.
- Yixuan Weng, Minjun Zhu, Fei Xia, Bin Li, Shizhu He, Shengping Liu, Bin Sun, Kang Liu, and Jun Zhao. 2023. Large language models are better reasoners with self-verification. EMNLP 2023 Findings.
- Mert Yuksekgonul, Federico Bianchi, Pratyusha Kalluri, Dan Jurafsky, and James Zou. 2023. When and why vision-language models behave like bags-of-words, and what to do about it? In ICLR.
- Eric Zelikman, Eliana Lorch, Lester Mackey, and Adam Tauman Kalai. 2023. Self-taught optimizer (stop): Recursively self-improving code generation. arXiv.
- Aojun Zhou, Ke Wang, Zimu Lu, Weikang Shi, Sichun Luo, Zipeng Qin, Shaoqing Lu, Anya Jia, Linqi Song, Mingjie Zhan, and Hongsheng Li. 2023. Solving challenging math word problems using GPT-4 code interpreter with code-based self-verification. arXiv.
- Banghua Zhu, Hiteshi Sharma, Felipe Vieira Frujeri, Shi Dong, Chenguang Zhu, Michael I. Jordan, and Jiantao Jiao. 2023. Fine-tuning language models with advantage-induced policy alignment. arXiv.

# 712

## 713

714

715

716

718

719

720

721

722

723

724

726

727

728

729

730

731

732

733

734

735

736

737

739

740

741

742

743

744

745

746 747

748

750

751

752

753

754

755

756

757

759

761

762

763

764

765

766

767

768

770

771

# A Model Prompts

In this section, we provide the full model prompts that we used for the validation metrics. The generation metric does not use a prompt (generation is conditioned on the beginning of a statement). For each prompt, \$example is replaced by one of the 400 full statements before being given to a model.

# A.1 Zero Shot

Overview:

Your task is to say whether a self-referential statement is true or false. In cases where you think that the answer is ambiguous, you should use your best judgement and still say either 'True' or 'False'.

Example:

\$example

Answer:

# A.2 Few Shot

Overview:

Your task is to say whether a self-referential statement is true or false. In cases where you think that the answer is ambiguous, you should use your best judgement and still say either 'True' or 'False'. Make sure that you understand the given examples before continuing.

## Example:

I am the text you are reading. My last line has how many words? Two words.

Answer: True

## Example:

I am the text you are reading. My last line has how many words? Seven words.

Answer: False

Example:

```
this sentence

l o

o o

k s

a bit like the

l

e

t

t

e

r
```

#### "p" 772 773 Answer: True 774 775 776 Example: 778 this sentence 779 1 0 780 0 0 781 k s 782 783 a bit like the 1 784 785 е t 786 t 787 е 788 789 r "o" 790 791 Answer: False 792 793 794 Example: 795 796 here in reversed is order word the. 797 798 Answer: True 799 800 801 Example: 802 803 here in fine is order word the. 804 805 Answer: False 806 807 808 Now it is your turn to complete an example. 809 Fill in the answer. 810 811 Example: 812 813 \$example 814 815 Answer: 816 Chain of Thought A.3 817 Overview: 818 819 Your task is to say whether a self-referential 820 statement is true or false. In cases where you 821 think that the answer is ambiguous, you should 822 use your best judgement and still say either 823 'True' or 'False'. Make sure that you understand 824 the given examples before continuing. 825 826 827 Example: 828 829 I am the text you are reading. 830 My last line has how many words? 831 Two words. 832 833 Answer: The text has 3 lines. The last line is 834 "Two words". The number of words in "Two words" 835

is 2. So the statement is true.

836

837

838

839

Example:

```
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
867
868
869
870
871
872
873
874
875
876
877
878
879
881
882
883
884
885
886
887
890
891
892
893
894
895
896
897
900
901
902
903
904
905
906
907
```

Seven words. Answer: The text has 3 lines. The last line is "Seven words". The number of words in "Seven words" is 2. So the statement is false. Example: this sentence k s a bit like the e t t е r "p" Answer: Visually, the words and letters of the sentence are arranged to look like a "p" on the page. So, the statement is true. Example: this sentence k s a bit like the 

I am the text you are reading.

My last line has how many words?

Answer: Visually, the words and letters of the sentence are arranged to look like a "p" on the page, not an "o". So, the statement is false.

#### Example:

е

t t

е

r "o"

here in reversed is order word the.

Answer: We have to notice that, if the word order is reversed, the statement reads: "the word order is reversed in here". In the original reversed statement, that is true. So the statement is true.

### Example:

here in fine is order word the.

Answer: We have to notice that, if the word order is reversed, the statement reads: "the word order is fine in here". In the original reversed statement, the word order is not fine. So the statement is false.

	910
Now it is your turn to complete an example.	911
Fill in the answer.	912
	913
Example:	914
	915
\$example	916
	917
Answer:	918
A.4 Zero Shot (Non-Self-Referent)	919
Queruieuu	0.00
overview:	920
Your task is to say whether a statement about a	921
sentence is true or false. In cases where you think	922
that the answer is ambiguous you should use your	924
hest judgement and still say either 'True' or	925
'False'	926
	927
	928
Example:	929
	930
<pre>\$example</pre>	931
	932
Answer:	933
A.5 Few Shot (Non-Self-Referent)	934
Overview:	935
	936
Your task is to say whether a statement about a	937
sentence is true or false. In cases where you think	938
that the answer is ambiguous, you should use your	939
best judgement and still say either 'True' or	940
'False'. Make sure that you understand the given	941
examples before continuing.	942
	943
<b>F</b>	944
Example:	945
The number of words in the following contenes is	940
nine number of words in the following sentence is	947
	940 0/0
The number of words in this sentence is nine	050
The number of words in this sentence is hille.	951
Answer: True	952
	953

### Example:

The number of words in the following sentence is eight.

The number of words in this sentence is eight.

```
Answer: False
```

### Example:

the following sentence k s a bit like the е t t е 

977	r		1047
978	"p"		1048
979		Now it is your turn to complete an example.	1049
980	this sentence	Fill in the answer.	1050
981	1 0		1051
982	0 0	Example:	1052
983	k s	p	1053
984	a hit like the	\$example	1054
985	1	+0.00mp10	1055
986	<u>е</u>	Answer	1056
987	t		1000
988	t		
989	6	A.6 Chain of Thought (Non-Self-Referent)	1057
990	r	Querrui eus	1050
991	"n"	Overview:	1050
992	þ	Vous took to to only whather a statement shout a	1059
992	Answer: True	Your task is to say whether a statement about a	1061
994	Answer. If de	sentence is true or faise. In cases where you think	1001
995		that the answer is ambiguous, you should use your	1062
996	Example:	Jest Judgement and still say either frue or	1003
997	Example.	Faise . Make sure that you understand the given	1004
998	the following sentence	examples before continuing.	1000
999			1000
1000		Evennler	1007
1001	k s	Example:	1000
1002	a bit like the	The number of words in the following contenes is	1009
1002		nie number of words in the following sentence is	1070
1004	-	nine.	1071
1005	t	The number of words in this contones is nine	1072
1006	+	The number of words in this sentence is nine.	1073
1007		Answer. Counting the number of words in the	1074
1008	r	Answer: counting the humber of words in the	1075
1009	"0"	sentence, we see that there are nine words. So,	1070
1010	0	the statement is true.	1077
1011	this sentence		1078
1012		Evenue la	1079
1012		Example:	1080
1014		The number of sends in the Callesian contenue is	1081
1015	a hit like the	The humber of words in the following sentence is	1082
1016		eight.	1083
1017		The number of sends in this contains is sight	1084
1018	t	The number of words in this sentence is eight.	1000
1019	+	Answer. Counting the number of words in the	1000
1020		Answer: counting the humber of words in the	1007
1021	r	sentence, we see that there are nine words,	1000
1022	"o"	not eight. So, the statement is faise.	1009
1023	0		1090
1024	Answer: False	Evample.	1091
1025		Example:	1092
1026		the following contance	1093
1027	Example:		1094
1028	Example:		1095
1029	The food words that the following sentence contains		1090
1020	are guava and cake and also carrot	K S	1097
1031	a. e guara ana cane ana aroo carrot.	a DIL IIKE UNE	1098
1032	The food words that this sentence contains are		1400
1033	guava and cake and also carrot	ट +	1404
1034		L +	1101
1035	Answer: True		1102
1036			1103
1037		ו אי <i>ה</i> וו	1104
1038	Example:	Ч	1105
1039	Example.	this sontoneo	1407
1040	The food words that the following sentence contains		110/
1041	are guava and cake but not carrot		1100
1042	a Suara and care but not carrot.		1109
1043	The food words that this sentence contains are	n S a hit lika tha	1110
1044	guava and cake but not carrot	a DIC IIKE UNE	1111
1045	Buite and cake but not carrot.		1112
1046	Answer: False	स् +	1113
10-10		L	1114

```
t
e
r
"p"
Answer: Visually, the words and letters of
```

the sentence are arranged to look like a "p" on the page. So, the statement is true.

Example:

1115

1116

1117

1118

1120

1121

1122 1123

1124 1125

1126

1127

1128

1129

1130

1131

1132

1133

1134

1135

1136

1137

1138 1139

1140

1141

1142

1143

1144

1145

1146

1147

1148

1149

1150

1151

1153

1154

1155

1156

1157

1158

1159

1160

1161

1162

1163

1164

1165

1166 1167

1168 1169 1170

1171

1172 1173

1174 1175

1176

1177

1178 1179

1180

```
the following sentence
1
                0
0
                0
k
                s
a bit like the
1
е
t
t
е
r
"o"
this sentence
1
                0
0
                0
k
                s
a bit like the
1
е
t
t
е
r
"o"
```

Answer: Visually, the words and letters of the sentence are arranged to look like a "p" on the page, not an "o". So, the statement is false.

### Example:

The food words that the following sentence contains are guava and cake and also carrot.

The food words that this sentence contains are guava and cake and also carrot.

Answer: The food words mentioned are indeed guava, cake, and carrot. So, the statement is true.

### Example:

The food words that the following sentence contains are guava and cake but not carrot.

The food words that this sentence contains are guava and cake but not carrot.

Answer: The food words mentioned in the sentence are guava, cake, and carrot. It is not true that carrot is not in the sentence. So, the statement is false.

Now it is your turn to complete an example.

Fill in the answer.	1185
	1186
Example:	1187
	1188
<pre>\$example</pre>	1189
	1190
Answer:	1191

1192

1193

1194

1195

1196

1197

1198

1199

1200

1201

1202

1204

1205

1206

1207

1208

1209

### **B** I am an Impossible Dataset

The schema of "I am a Strange Dataset" does not place any limits on the complexity of an example. One of the expert annotators in particular took full advantage of this while they were creating examples. They generated 10 examples that another expert annotator could not understand without an explanation. These examples are not in our main dataset, where we verify that non-expert people can perform well. We release these 10 examples separately to showcase just how hard examples fitting this self-referential schema can get. We also hope that when language models improve, these 10 examples can be a useful initial probe into their reasoning capabilities. We provide two examples from this seemingly "Impossible" set in Figure 8. We also provide GPT 4's output on another example from this dataset in Figure 9 below.

Now, I list a digit thirteens:	Now, I list a digit thirteens:
pi. Letter count the words.	pi. Letter count the words.
Conceded? Errorless.	Conceded? Erroneous.
Answer fhdi yatt la ylmy	Answer: fbdi vatf la vlmu

Answer: fbdi yqtf lq xlmw hnumjw ynolzy if vul. Answer: fbdi yqtf lq xlmw hnumjw ynolzy if one.

Figure 8: Examples from "I am an Impossible Dataset". Examples from this dataset still follow our schema, but the expert annotators needed to see an explanation before they understood. The first example says that the letter count of each of the words are the thirteen first digits of  $\pi$ : 3.141592653589. The beginning of the second example uses a Caesar cipher which shifts by one each successive word and reads "each word in this cipher shifts by". "vul" with a shift of 7 is "one".

"Fivetune favors the bold." "After that meal, I was snined." "Each of the number sounds in these sentences, including this sentence, are increased by two."

**GPT 4:** The first sentence "Fivetune favors the bold." has the number sound "five" instead of "fortune". The second sentence "After that meal, I was snined." has the number sound "nine" instead of "sated". The third sentence "Each of the number sounds in these sentences, including this sentence, are increased by two." suggests that each number sound should be increased by two. However, "five" increased by two is "seven", not "fortune", and "nine" increased by two is "eleven", not "sated". So, the statement is false.

Figure 9: GPT 4's Dec 26, 2023 output on an example from "I am an Impossible Dataset". It understands that some of the original words should be "fortune" (e.g. 4tune) and "sated" (e.g. s8ed). But GPT 4 misses that the statement is trying to say that the number sounds in every sentence are increased by one. The last sentence cannot say "one" explicitly - it needs to say "two" in order for the statement to stay true. We would not expect a typical person to understand this example, but will a language model eventually grasp it?

## 1210 C Mechanical Turk Annotator Interface

1211

1212

1213

Here, we show a screenshot of the interface used by Mechanical Turk workers in Figure 10.



Figure 10: A screenshot of the Mechanical Turk worker interface for validating statements.

## **D** Supplemental Discussion

Here, we present supplemental discussion about
why models are showing poor performance on "I
am a Strange Dataset".

## **D.1** A Test of the Tokenizer?

Our tests are related to whether a model truly "sees" 1218 text in the same way as people. Metalinguistic state-1219 ments may refer to the number of characters that 1220 they have, how text is arranged on the page, capi-1221 talization of certain letters, and relative positions 1222 of words. A human can easily count the charac-1223 ters that they see in a sentence, but models tend to 1224 encode text in tokens, not characters. We do not 1225 provide tests to disentangle the impact of different 1226 tokenizers, so this section is speculative. 1227

1217

1228

1242

## **D.2** Training Data Limitations

Practically speaking, it is unlikely that there 1229 are many examples of metalinguistic statements 1230 in training datasets. They are incredibly time-1231 intensive to generate, even if they are easy to 1232 verify. Yet people, who have almost surely seen 1233 even fewer examples, can do much better at this 1234 task than models. This "hard-to-create, easy-to-1235 verify" feature of hard evaluation datasets is true 1236 of Winoground (Thrush et al., 2022) too, which is 1237 a vision and language evaluation dataset that has remained unsaturated for well over a year. This 1239 goes to show that our models still have the wrong 1240 biases - will they change simply with model scale? 1241

## E Results by Tag

In this section, we provide results in a table for each of the 10 tags. 1243

Model	Params	Chat	Gen <sup>L</sup>	Val $ZS^L$	Val FS <sup><math>L</math></sup>	Val ZS $^{L}(\mathbf{R})$	Val FS $^{L}(\mathbf{R})$	Val $CoT^T$
Llama 2	7B	Ν	$\textbf{62.90} \pm \textbf{12.10}$	$49.19 \pm 1.21$	$49.19 \pm 4.84$	$46.77 \pm 12.90$	$50.00 \pm 12.90$	$6.45\pm4.03$
Llama 2	7B	Y	$58.06 \pm 12.90$	$50.81 \pm 1.21$	$50.00\pm0.00$	$37.10 \pm 11.29$	$56.45 \pm 12.90$	$16.94\pm 6.85$
Mistral 0.1	7B	Ν	$53.23 \pm 12.90$	$50.00\pm2.42$	$50.81 \pm 1.21$	$56.45 \pm 12.90$	$51.61 \pm 12.90$	$0.00\pm0.00$
Starling $\alpha$	7B	Y	$54.84 \pm 12.90$	$49.19 \pm 2.42$	$49.19 \pm 1.21$	$54.84 \pm 12.90$	$51.61 \pm 12.90$	$32.26\pm8.87$
Mistral 0.2	7B	Y	$51.61 \pm 12.90$	$50.00\pm6.45$	$50.00\pm5.65$	$46.77 \pm 12.10$	$53.23 \pm 12.10$	$53.23\pm8.06$
Llama 2	13B	Ν	$\textbf{62.90} \pm \textbf{11.29}$	$47.58\pm3.63$	$52.42\pm5.24$	$41.94 \pm 12.10$	$58.06 \pm 12.90$	$4.03\pm3.23$
Llama 2	13B	Y	$59.68 \pm 12.90$	$49.19 \pm 4.84$	$51.61\pm3.23$	$46.77 \pm 12.90$	$50.00 \pm 12.90$	$7.26 \pm 4.44$
Mixtral 0.1	8x7B	Ν	$54.84 \pm 12.90$	$53.23 \pm 4.84$	$50.00\pm0.00$	$48.39 \pm 12.90$	$48.39 \pm 12.10$	$3.23\pm2.82$
Mixtral 0.1	8x7B	Y	$53.23 \pm 12.90$	$49.19 \pm 4.84$	$45.97 \pm 5.65$	$53.23 \pm 12.90$	$48.39 \pm 12.90$	$36.29 \pm 8.06$
Llama 2	70B	Ν	$61.29 \pm 11.29$	$52.42\pm5.65$	$50.81 \pm 4.03$	$59.68 \pm 11.29$	$50.00 \pm 12.90$	$3.23\pm2.82$
Llama 2	70B	Y	$51.61 \pm 12.90$	$53.23\pm 6.45$	$49.19\pm3.23$	$56.45 \pm 12.90$	$46.77 \pm 12.90$	$31.45\pm7.26$
Claude 2	-	Y	-	-	-	-	-	$48.39\pm 6.85$
GPT 3.5 T	-	Y	-	$48.39 \pm 4.84$	$50.81 \pm 5.65$	$56.45 \pm 12.90$	$50.00 \pm 12.90$	$44.35\pm7.26$
GPT 4	-	Y	-	$53.23\pm8.06$	$53.23\pm8.06$	$54.84 \pm 12.90$	$53.23\pm12.90$	$55.65\pm9.27$

Table 5: Results for the 62 example pairs with the Numerical Operations tag. Scores with 95% confidence intervals above chance are shown in **bold**.

Model	Params	Chat	Gen <sup>L</sup>	Val $ZS^L$	Val $FS^L$	$\mathrm{Val}\mathrm{ZS}^{L}\left(R\right)$	Val FS $^{L}(\mathbf{R})$	Val $CoT^T$
Llama 2	7B	Ν	$54.35\pm9.78$	$49.46\pm2.17$	$51.63\pm3.53$	$47.83 \pm 10.33$	$\textbf{66.30} \pm \textbf{9.78}$	$3.80\pm2.72$
Llama 2	7B	Y	$48.91 \pm 9.78$	$54.35\pm5.16$	$50.54 \pm 0.82$	$\textbf{60.87} \pm \textbf{9.78}$	$\textbf{61.96} \pm \textbf{9.78}$	$13.04\pm4.35$
Mistral 0.1	7B	Ν	$51.09 \pm 9.78$	$\textbf{55.43} \pm \textbf{4.89}$	$49.46\pm2.17$	$\textbf{63.04} \pm \textbf{9.78}$	$56.52\pm9.78$	$0.00\pm0.00$
Starling $\alpha$	7B	Y	$50.00\pm9.78$	$\textbf{59.24} \pm \textbf{5.43}$	$51.63 \pm 2.17$	$\textbf{60.87} \pm \textbf{9.78}$	$\textbf{61.96} \pm \textbf{9.78}$	$38.04 \pm 6.79$
Mistral 0.2	7B	Y	$51.09 \pm 9.78$	$54.35\pm6.79$	$52.17\pm5.98$	$56.52\pm9.78$	$57.61 \pm 9.78$	$50.00\pm 6.52$
Llama 2	13B	Ν	$53.26\pm9.78$	$54.89 \pm 5.98$	$\textbf{56.52} \pm \textbf{5.98}$	$55.43 \pm 9.78$	$\textbf{65.22} \pm \textbf{9.78}$	$5.98 \pm 3.26$
Llama 2	13B	Y	$48.91 \pm 9.78$	$55.98 \pm 6.52$	$52.17 \pm 4.35$	$57.61 \pm 9.78$	$50.00\pm9.78$	$11.96\pm5.16$
Mixtral 0.1	8x7B	Ν	$55.43 \pm 9.78$	$\textbf{63.59} \pm \textbf{6.52}$	$52.72\pm4.35$	$\textbf{64.13} \pm \textbf{9.78}$	$\textbf{60.87} \pm \textbf{9.78}$	$3.80\pm2.99$
Mixtral 0.1	8x7B	Y	$55.43 \pm 10.33$	$\textbf{57.07} \pm \textbf{6.52}$	$\textbf{57.07} \pm \textbf{4.62}$	$\textbf{61.96} \pm \textbf{9.78}$	$57.61 \pm 10.33$	$50.00\pm7.07$
Llama 2	70B	Ν	$55.43 \pm 9.78$	$\textbf{57.07} \pm \textbf{5.43}$	$\textbf{59.24} \pm \textbf{4.89}$	$\textbf{66.30} \pm \textbf{9.78}$	$\textbf{63.04} \pm \textbf{9.78}$	$1.09 \pm 1.36$
Llama 2	70B	Y	$51.09 \pm 10.33$	$\textbf{59.24} \pm \textbf{7.07}$	$52.72\pm3.80$	$\textbf{61.96} \pm \textbf{9.78}$	$\textbf{65.22} \pm \textbf{9.78}$	$21.74\pm 6.25$
Claude 2	-	Y	-	-	-	-	-	$55.98 \pm 7.07$
GPT 3.5 T	-	Y	-	$54.89 \pm 4.89$	$\textbf{55.98} \pm \textbf{5.72}$	$57.61 \pm 9.78$	$\textbf{71.74} \pm \textbf{9.24}$	$55.98 \pm 6.25$
GPT 4	-	Y	-	$\textbf{63.04} \pm \textbf{6.52}$	$\textbf{61.96} \pm \textbf{6.52}$	$\textbf{70.65} \pm \textbf{9.24}$	$\textbf{71.74} \pm \textbf{9.24}$	$\textbf{72.83} \pm \textbf{6.52}$

Table 6: Results for the 94 example pairs with the Negation & Scope tag. Scores with 95% confidence intervals above chance are shown in **bold**.

Model	Params	Chat	Gen <sup>L</sup>	Val $ZS^L$	$\mathrm{Val}\mathrm{FS}^L$	Val $ZS^{L}(\mathbf{R})$	Val $FS^{L}(\mathbf{R})$	Val $CoT^T$
Llama 2	7B	Ν	$45.83 \pm 20.83$	$50.00\pm0.00$	$50.00\pm 6.25$	$58.33 \pm 20.83$	$58.33 \pm 20.83$	$8.33 \pm 7.29$
Llama 2	7B	Y	$45.83\pm20.83$	$52.08 \pm 6.25$	$50.00\pm0.00$	$62.50\pm20.83$	$\textbf{70.83} \pm \textbf{18.75}$	$22.92 \pm 10.42$
Mistral 0.1	7B	Ν	$45.83\pm20.83$	$52.08 \pm 6.25$	$54.17\pm5.21$	$66.67\pm18.75$	$\textbf{70.83} \pm \textbf{18.75}$	$0.00\pm0.00$
Starling $\alpha$	7B	Y	$45.83\pm20.83$	$56.25 \pm 10.42$	$47.92\pm3.12$	$66.67 \pm 18.75$	$66.67 \pm 18.75$	$37.50 \pm 12.50$
Mistral 0.2	7B	Y	$54.17\pm20.83$	$54.17 \pm 10.42$	$62.50\pm12.50$	$62.50\pm20.83$	$66.67 \pm 18.75$	$47.92 \pm 12.50$
Llama 2	13B	Ν	$45.83\pm20.83$	$52.08 \pm 6.25$	$56.25\pm8.33$	$54.17\pm20.83$	$58.33 \pm 20.83$	$2.08\pm3.12$
Llama 2	13B	Y	$45.83\pm20.83$	$52.08 \pm 10.42$	$52.08\pm3.12$	$58.33\pm20.83$	$66.67 \pm 18.75$	$18.75\pm11.46$
Mixtral 0.1	8x7B	Ν	$37.50 \pm 18.75$	$\textbf{64.58} \pm \textbf{10.42}$	$54.17\pm5.21$	$62.50\pm18.75$	$62.50\pm20.83$	$12.50\pm9.38$
Mixtral 0.1	8x7B	Y	$50.00\pm20.83$	$47.92\pm8.33$	$54.17 \pm 10.42$	$50.00\pm20.83$	$58.33 \pm 20.83$	$41.67 \pm 14.58$
Llama 2	70B	Ν	$45.83\pm20.83$	$50.00\pm8.33$	$56.25 \pm 12.50$	$54.17\pm20.83$	$62.50\pm18.75$	$12.50\pm8.33$
Llama 2	70B	Y	$50.00\pm20.83$	$43.75\pm10.42$	$47.92\pm3.12$	$54.17\pm20.83$	$58.33 \pm 20.83$	$10.42\pm8.33$
Claude 2	-	Y	-	-	-	-	-	$47.92 \pm 12.50$
GPT 3.5 T	-	Y	-	$50.00\pm0.00$	$47.92 \pm 6.25$	$54.17\pm20.83$	$58.33 \pm 20.83$	$52.08 \pm 12.50$
GPT 4	-	Y	-	$54.17 \pm 10.42$	$\textbf{66.67} \pm \textbf{13.54}$	$58.33\pm20.83$	$62.50\pm18.75$	$60.42 \pm 14.58$

Table 7: Results for the 24 example pairs with the Multi-Channel tag. Scores with 95% confidence intervals above chance are shown in **bold**.

Model	Params	Chat	Gen <sup>L</sup>	Val $ZS^L$	Val FS <sup><math>L</math></sup>	$\mathrm{Val}\mathrm{ZS}^{L}\left(\mathrm{R}\right)$	Val FS $^{L}(\mathbf{R})$	Val $CoT^T$
Llama 2	7B	Ν	$50.00 \pm 12.96$	$48.15\pm2.31$	$48.15\pm3.70$	$35.19 \pm 12.96$	$46.30\pm12.96$	$6.48 \pm 4.17$
Llama 2	7B	Y	$44.44 \pm 12.96$	$48.15\pm4.63$	$49.07 \pm 1.39$	$46.30\pm12.96$	$\textbf{38.89} \pm \textbf{12.96}$	$12.04\pm5.56$
Mistral 0.1	7B	Ν	$46.30\pm12.96$	$47.22\pm3.24$	$47.22\pm3.24$	$37.04 \pm 12.96$	$44.44 \pm 12.96$	$0.00\pm0.00$
Starling $\alpha$	7B	Y	$48.15\pm12.96$	$48.15\pm5.09$	$50.93 \pm 1.39$	$44.44 \pm 12.96$	$46.30\pm12.96$	$28.70\pm8.33$
Mistral 0.2	7B	Y	$46.30\pm12.96$	$52.78 \pm 7.41$	$48.15\pm7.41$	$44.44 \pm 12.96$	$37.04 \pm 12.96$	$43.52\pm8.33$
Llama 2	13B	Ν	$48.15\pm12.96$	$48.15\pm5.56$	$49.07\pm5.56$	$44.44 \pm 12.96$	$44.44 \pm 12.96$	$3.70\pm3.24$
Llama 2	13B	Y	$44.44 \pm 12.96$	$47.22\pm7.41$	$48.15\pm4.63$	$38.89 \pm 12.96$	$35.19 \pm 12.96$	$8.33 \pm 5.09$
Mixtral 0.1	8x7B	Ν	$46.30\pm12.96$	$54.63\pm5.56$	$49.07\pm3.70$	$44.44 \pm 12.96$	$46.30\pm12.96$	$1.85\pm2.31$
Mixtral 0.1	8x7B	Y	$44.44 \pm 12.96$	$51.85\pm8.33$	$50.00\pm6.48$	$44.44 \pm 12.96$	$53.70 \pm 12.96$	$44.44\pm9.26$
Llama 2	70B	Ν	$50.00\pm12.96$	$50.00\pm4.63$	$50.93 \pm 3.70$	$51.85 \pm 12.96$	$44.44 \pm 12.96$	$0.00\pm0.00$
Llama 2	70B	Y	$44.44 \pm 12.96$	$46.30\pm7.41$	$47.22\pm4.17$	$44.44 \pm 12.96$	$51.85 \pm 12.96$	$29.63 \pm 8.33$
Claude 2	-	Y	-	-	-	-	-	$50.00\pm7.41$
GPT 3.5 T	-	Y	-	$51.85\pm5.09$	$46.30\pm6.48$	$51.85 \pm 12.96$	$55.56 \pm 12.96$	$54.63\pm9.26$
GPT 4	-	Y	-	$56.48 \pm 7.41$	$58.33\pm9.26$	$\textbf{64.81} \pm \textbf{12.96}$	$62.96 \pm 12.96$	$\textbf{71.30} \pm \textbf{8.80}$

Table 8: Results for the 55 example pairs with the Location of Element tag. Scores with 95% confidence intervals above chance are shown in **bold**.

Model	Params	Chat	Gen <sup>L</sup>	Val $ZS^L$	Val FS <sup><math>L</math></sup>	$\mathrm{Val}\ \mathrm{ZS}^{L}\left( \mathrm{R}\right)$	Val FS $^{L}(\mathbf{R})$	Val $CoT^T$
Llama 2	7B	Ν	$\textbf{66.67} \pm \textbf{14.29}$	$50.00\pm0.00$	$47.62 \pm 4.76$	$52.38 \pm 14.29$	$52.38 \pm 14.29$	$7.14\pm5.36$
Llama 2	7B	Y	$61.90 \pm 14.29$	$54.76\pm5.36$	$48.81 \pm 1.79$	$54.76\pm14.29$	$52.38 \pm 14.29$	$25.00\pm8.93$
Mistral 0.1	7B	Ν	$\textbf{66.67} \pm \textbf{14.29}$	$50.00\pm0.00$	$47.62\pm2.98$	$52.38 \pm 14.29$	$52.38 \pm 14.29$	$0.00\pm0.00$
Starling $\alpha$	7B	Y	$\textbf{66.67} \pm \textbf{14.29}$	$53.57 \pm 5.36$	$51.19\pm3.57$	$59.52 \pm 14.29$	$54.76\pm14.29$	$32.14\pm9.52$
Mistral 0.2	7B	Y	$57.14 \pm 14.29$	$48.81\pm8.33$	$54.76\pm7.14$	$50.00\pm14.29$	$50.00\pm14.29$	$55.95\pm8.33$
Llama 2	13B	Ν	$\textbf{73.81} \pm \textbf{13.10}$	$45.24\pm5.95$	$54.76\pm5.36$	$52.38\pm14.29$	$54.76\pm14.29$	$3.57\pm4.17$
Llama 2	13B	Y	$\textbf{71.43} \pm \textbf{13.10}$	$55.95\pm8.33$	$48.81 \pm 1.79$	$59.52\pm14.29$	$52.38\pm14.29$	$13.10\pm 6.55$
Mixtral 0.1	8x7B	Ν	$57.14 \pm 14.29$	$55.95 \pm 7.14$	$46.43 \pm 4.17$	$52.38\pm14.29$	$54.76\pm14.29$	$3.57\pm4.17$
Mixtral 0.1	8x7B	Y	$57.14 \pm 14.29$	$53.57\pm7.14$	$52.38\pm 6.55$	$42.86\pm14.29$	$54.76\pm14.29$	$34.52\pm8.33$
Llama 2	70B	Ν	$\textbf{69.05} \pm \textbf{14.29}$	$52.38 \pm 7.14$	$\textbf{61.90} \pm \textbf{6.55}$	$61.90 \pm 14.29$	$61.90 \pm 14.29$	$2.38\pm2.98$
Llama 2	70B	Y	$\textbf{64.29} \pm \textbf{14.29}$	$55.95 \pm 8.33$	$50.00\pm3.57$	$59.52 \pm 14.29$	$61.90 \pm 14.29$	$13.10\pm7.74$
Claude 2	-	Y	-	-	-	-	-	$\textbf{59.52} \pm \textbf{8.33}$
GPT 3.5 T	-	Y	-	$48.81 \pm 4.76$	$53.57\pm7.14$	$\textbf{64.29} \pm \textbf{14.29}$	$61.90 \pm 14.29$	$47.62\pm9.52$
GPT 4	-	Y	-	$50.00\pm7.14$	$\textbf{58.33} \pm \textbf{7.14}$	$59.52\pm14.29$	$59.52\pm14.29$	$55.95\pm9.52$

Table 9: Results for the 42 example pairs with the Sensory tag. Scores with 95% confidence intervals above chance are shown in **bold**.

Model	Params	Chat	Gen <sup>L</sup>	Val $ZS^L$	Val $FS^L$	$Val ZS^{L} (R)$	Val $FS^{L}(R)$	Val $\mathbf{CoT}^T$
Llama 2	7B	Ν	$62.50 \pm 13.54$	$50.00\pm0.00$	$48.96 \pm 6.25$	$41.67 \pm 14.58$	$41.67 \pm 14.58$	$9.38 \pm 5.21$
Llama 2	7B	Y	$58.33 \pm 14.58$	$48.96 \pm 6.25$	$50.00\pm0.00$	$31.25 \pm 12.50$	$62.50\pm14.58$	$8.33 \pm 5.21$
Mistral 0.1	7B	Ν	$62.50\pm13.54$	$46.88 \pm 4.17$	$48.96\pm3.12$	$41.67 \pm 13.54$	$43.75\pm14.58$	$0.00\pm0.00$
Starling $\alpha$	7B	Y	$\textbf{64.58} \pm \textbf{13.54}$	$48.96\pm3.12$	$48.96 \pm 1.56$	$39.58 \pm 13.54$	$52.08 \pm 14.58$	$31.25\pm8.85$
Mistral 0.2	7B	Y	$\textbf{64.58} \pm \textbf{13.54}$	$45.83\pm8.33$	$50.00\pm 6.25$	$47.92 \pm 14.58$	$41.67 \pm 14.58$	$46.88\pm8.33$
Llama 2	13B	Ν	$60.42 \pm 13.54$	$48.96 \pm 6.25$	$47.92\pm8.33$	$43.75\pm14.58$	$47.92 \pm 14.58$	$5.21\pm4.17$
Llama 2	13B	Y	$56.25 \pm 14.58$	$46.88 \pm 7.29$	$48.96 \pm 4.69$	$45.83\pm14.58$	$39.58 \pm 14.58$	$5.21\pm4.17$
Mixtral 0.1	8x7B	Ν	$54.17 \pm 14.58$	$52.08 \pm 6.25$	$50.00\pm0.00$	$33.33 \pm 13.54$	$52.08 \pm 14.58$	$2.08 \pm 2.60$
Mixtral 0.1	8x7B	Y	$60.42 \pm 14.58$	$52.08 \pm 7.29$	$52.08 \pm 6.25$	$50.00\pm14.58$	$60.42 \pm 14.58$	$46.88\pm8.33$
Llama 2	70B	Ν	$60.42 \pm 13.54$	$47.92\pm 6.25$	$53.12\pm5.21$	$52.08 \pm 14.58$	$47.92 \pm 14.58$	$1.04 \pm 1.56$
Llama 2	70B	Y	$54.17 \pm 14.58$	$50.00\pm7.29$	$47.92 \pm 4.17$	$41.67 \pm 14.58$	$56.25 \pm 13.54$	$23.96 \pm 8.33$
Claude 2	-	Y	-	-	-	-	-	$57.29 \pm 7.29$
GPT 3.5 T	-	Y	-	$52.08 \pm 4.17$	$53.12\pm 6.25$	$54.17 \pm 14.58$	$60.42 \pm 13.54$	$45.83\pm9.38$
GPT 4	-	Y	-	$\textbf{61.46} \pm \textbf{6.77}$	$\textbf{62.50} \pm \textbf{7.81}$	$\textbf{72.92} \pm \textbf{12.50}$	$\textbf{70.83} \pm \textbf{12.50}$	$\textbf{62.50} \pm \textbf{9.38}$

Table 10: Results for the 48 example pairs with the Sub-Word tag. Scores with 95% confidence intervals above chance are shown in **bold**.

Model	Params	Chat	Gen <sup>L</sup>	Val $ZS^L$	Val FS <sup><math>L</math></sup>	Val ZS $^{L}$ (R)	Val $FS^{L}(R)$	Val $CoT^T$
Llama 2	7B	Ν	$64.52 \pm 16.13$	$51.61 \pm 2.42$	$53.23\pm 6.45$	$54.84 \pm 16.13$	$\textbf{70.97} \pm \textbf{16.13}$	$4.84\pm5.65$
Llama 2	7B	Y	$54.84 \pm 16.13$	$56.45 \pm 7.26$	$50.00\pm0.00$	$\textbf{67.74} \pm \textbf{16.13}$	$54.84 \pm 16.13$	$11.29\pm7.26$
Mistral 0.1	7B	Ν	$58.06 \pm 16.13$	$\textbf{58.06} \pm \textbf{6.45}$	$51.61 \pm 2.42$	$\textbf{67.74} \pm \textbf{16.13}$	$\textbf{67.74} \pm \textbf{16.13}$	$0.00\pm0.00$
Starling $\alpha$	7B	Y	$51.61 \pm 16.13$	$\textbf{62.90} \pm \textbf{7.26}$	$53.23 \pm 4.03$	$\textbf{67.74} \pm \textbf{16.13}$	$54.84 \pm 17.74$	$46.77\pm11.29$
Mistral 0.2	7B	Y	$61.29 \pm 16.13$	$53.23 \pm 12.90$	$54.84 \pm 9.68$	$64.52\pm16.13$	$61.29 \pm 16.13$	$53.23 \pm 11.29$
Llama 2	13B	Ν	67.74 ± 16.13	$59.68 \pm 9.68$	$59.68 \pm 10.48$	$\textbf{67.74} \pm \textbf{16.13}$	$\textbf{77.42} \pm \textbf{14.52}$	$8.06\pm6.45$
Llama 2	13B	Y	$58.06 \pm 16.13$	$\textbf{62.90} \pm \textbf{9.68}$	$\textbf{56.45} \pm \textbf{5.65}$	$\textbf{67.74} \pm \textbf{16.13}$	$64.52\pm16.13$	$8.06\pm7.26$
Mixtral 0.1	8x7B	Ν	$58.06 \pm 16.13$	$\textbf{69.35} \pm \textbf{9.68}$	$54.84 \pm 4.84$	$\textbf{74.19} \pm \textbf{16.13}$	$61.29 \pm 16.13$	$1.61\pm2.42$
Mixtral 0.1	8x7B	Y	$54.84 \pm 17.74$	$50.00\pm9.68$	$54.84\pm6.45$	$58.06 \pm 16.13$	$38.71 \pm 16.13$	$41.94 \pm 9.68$
Llama 2	70B	Ν	$64.52 \pm 16.13$	$\textbf{66.13} \pm \textbf{9.68}$	$\textbf{62.90} \pm \textbf{7.26}$	$\textbf{67.74} \pm \textbf{16.13}$	$\textbf{67.74} \pm \textbf{16.13}$	$1.61\pm2.42$
Llama 2	70B	Y	$61.29 \pm 17.74$	$\textbf{67.74} \pm \textbf{12.10}$	$50.00\pm4.84$	$\textbf{70.97} \pm \textbf{16.13}$	$\textbf{74.19} \pm \textbf{14.52}$	$25.81\pm11.29$
Claude 2	-	Y	-	-	-	-	-	$59.68 \pm 9.68$
GPT 3.5 T	-	Y	-	$54.84 \pm 5.65$	$\textbf{64.52} \pm \textbf{8.06}$	$58.06 \pm 16.13$	$64.52\pm16.13$	$46.77\pm12.10$
GPT 4	-	Y	-	$\textbf{64.52} \pm \textbf{11.29}$	$61.29\pm11.29$	$\textbf{67.74} \pm \textbf{16.13}$	$\textbf{70.97} \pm \textbf{16.13}$	$\textbf{75.81} \pm \textbf{12.10}$

Table 11: Results for the 31 example pairs with the Existence of Element tag. Scores with 95% confidence intervals above chance are shown in **bold**.

Model	Params	Chat	Gen <sup>L</sup>	Val $ZS^L$	Val FS <sup><math>L</math></sup>	$\mathrm{Val}\ \mathrm{ZS}^{L}\left( \mathrm{R}\right)$	Val FS $^{L}$ (R)	Val $CoT^T$
Llama 2	7B	Ν	58.33 ± 20.83	$50.00\pm0.00$	$50.00\pm0.00$	$41.67\pm20.83$	$41.67\pm20.83$	$6.25\pm 6.25$
Llama 2	7B	Y	$58.33 \pm 18.75$	$43.75\pm 6.25$	$50.00\pm0.00$	$54.17\pm20.83$	$58.33 \pm 20.83$	$29.17\pm10.42$
Mistral 0.1	7B	Ν	$41.67 \pm 20.83$	$43.75\pm7.29$	$52.08\pm3.12$	$45.83\pm20.83$	$54.17\pm20.83$	$0.00\pm0.00$
Starling $\alpha$	7B	Y	$37.50\pm20.83$	$50.00\pm 6.25$	$47.92\pm3.12$	$45.83\pm20.83$	$58.33 \pm 20.83$	$22.92 \pm 12.50$
Mistral 0.2	7B	Y	$45.83 \pm 20.83$	$50.00\pm10.42$	$47.92\pm8.33$	$41.67\pm20.83$	$50.00\pm20.83$	$41.67 \pm 12.50$
Llama 2	13B	Ν	$58.33 \pm 20.83$	$45.83 \pm 5.21$	$54.17\pm8.33$	$25.00\pm16.67$	$50.00\pm20.83$	$4.17\pm5.21$
Llama 2	13B	Y	$54.17 \pm 20.83$	$50.00\pm8.33$	$47.92\pm3.12$	$50.00\pm20.83$	$41.67\pm20.83$	$16.67\pm9.38$
Mixtral 0.1	8x7B	Ν	$50.00 \pm 20.83$	$50.00\pm8.33$	$47.92\pm3.12$	$45.83\pm20.83$	$50.00\pm20.83$	$2.08\pm3.12$
Mixtral 0.1	8x7B	Y	$41.67\pm20.83$	$43.75\pm8.33$	$47.92 \pm 10.42$	$37.50\pm18.75$	$37.50\pm18.85$	$29.17\pm12.50$
Llama 2	70B	Ν	$50.00 \pm 20.83$	$45.83\pm8.33$	$54.17 \pm 10.42$	$41.67\pm20.83$	$45.83\pm20.83$	$4.17\pm5.21$
Llama 2	70B	Y	$54.17 \pm 20.83$	$52.08 \pm 7.29$	$45.83 \pm 5.21$	$41.67\pm20.83$	$54.17\pm20.83$	$22.92 \pm 10.42$
Claude 2	-	Y	-	-	-	-	-	$41.67\pm10.42$
GPT 3.5 T	-	Y	-	$54.17\pm5.21$	$45.83\pm8.33$	$\textbf{70.83} \pm \textbf{16.67}$	$54.17\pm20.83$	$50.00\pm10.42$
GPT 4	-	Y	-	$47.92\pm9.38$	$52.08\pm14.58$	$41.67\pm20.83$	$62.50\pm18.75$	$45.83\pm12.50$

Table 12: Results for the 24 example pairs with the Hypothetical tag. Scores with 95% confidence intervals above chance are shown in **bold**.

Model	Params	Chat	Gen <sup>L</sup>	Val $ZS^L$	Val $FS^L$	Val $ZS^{L}(R)$	Val $FS^{L}(R)$	Val $CoT^T$
Llama 2	7B	Ν	$44.00 \pm 20.00$	$50.00\pm6.00$	$50.00\pm 6.00$	$40.00\pm20.00$	$56.00\pm20.00$	$2.00\pm3.00$
Llama 2	7B	Y	$60.00\pm20.00$	$50.00\pm10.00$	$50.00\pm0.00$	$52.00\pm20.00$	$56.00\pm20.00$	$6.00\pm6.00$
Mistral 0.1	7B	Ν	$64.00\pm18.00$	$50.00\pm8.00$	$46.00\pm5.00$	$52.00\pm20.00$	$60.00\pm20.00$	$0.00\pm0.00$
Starling $\alpha$	7B	Y	$64.00\pm20.00$	$54.00\pm10.00$	$50.00\pm0.00$	$56.00\pm20.00$	$60.00\pm20.00$	$44.00\pm15.00$
Mistral 0.2	7B	Y	$48.00\pm20.00$	$56.00 \pm 12.00$	$58.00 \pm 10.00$	$60.00\pm20.00$	$60.00\pm20.00$	$58.00 \pm 12.00$
Llama 2	13B	Ν	$60.00\pm20.00$	$56.00\pm8.00$	$50.00 \pm 12.00$	$52.00\pm20.00$	$\textbf{68.00} \pm \textbf{18.00}$	$6.00\pm7.00$
Llama 2	13B	Y	$56.00\pm20.00$	$54.00\pm12.00$	$54.00\pm8.00$	$52.00\pm20.00$	$52.00\pm20.00$	$8.00\pm9.00$
Mixtral 0.1	8x7B	Ν	$60.00\pm20.00$	$58.00 \pm 10.00$	$54.00\pm8.00$	$64.00\pm20.00$	$52.00\pm20.00$	$2.00\pm3.00$
Mixtral 0.1	8x7B	Y	$60.00\pm20.00$	$44.00\pm12.00$	$54.00\pm5.00$	$60.00\pm20.00$	$56.00\pm20.00$	$52.00\pm16.00$
Llama 2	70B	Ν	$56.00\pm20.00$	$\textbf{60.00} \pm \textbf{8.00}$	$56.00\pm6.00$	$56.00\pm20.00$	$56.00\pm20.00$	$0.00\pm0.00$
Llama 2	70B	Y	$56.00\pm20.00$	$52.00\pm10.00$	$48.00\pm3.00$	$56.00\pm20.00$	$56.00\pm20.00$	$30.00\pm10.00$
Claude 2	-	Y	-	-	-	-	-	$\textbf{68.00} \pm \textbf{10.00}$
GPT 3.5 T	-	Y	-	$56.00\pm8.00$	$54.00\pm8.00$	$56.00\pm20.00$	$\textbf{72.00} \pm \textbf{18.00}$	$42.00\pm13.05$
GPT 4	-	Y	-	$\textbf{68.00} \pm \textbf{13.00}$	$\textbf{64.00} \pm \textbf{13.00}$	$\textbf{84.00} \pm \textbf{14.00}$	$\textbf{80.00} \pm \textbf{16.00}$	$\textbf{80.00} \pm \textbf{11.00}$

Table 13: Results for the 25 example pairs with the Grammaticality tag. Scores with 95% confidence intervals above chance are shown in **bold**.

Model	Params	Chat	Gen <sup>L</sup>	Val $ZS^L$	Val FS <sup><math>L</math></sup>	$\mathrm{Val}\ \mathrm{ZS}^{L}\ (\mathrm{R})$	Val FS $^{L}$ (R)	Val $CoT^T$
Llama 2	7B	Ν	$59.09 \pm 20.45$	$52.27 \pm 6.82$	$59.09 \pm 10.23$	$59.09 \pm 18.18$	$54.55\pm20.45$	$2.27\pm3.41$
Llama 2	7B	Y	$59.09 \pm 20.45$	$63.64 \pm 13.64$	$52.27\pm3.41$	$59.09 \pm 18.18$	$54.55\pm22.73$	$20.45\pm10.23$
Mistral 0.1	7B	Ν	$36.36 \pm 20.45$	$\textbf{63.64} \pm \textbf{11.36}$	$50.00\pm0.00$	$\textbf{86.36} \pm \textbf{13.64}$	$59.09 \pm 20.45$	$0.00\pm0.00$
Starling $\alpha$	7B	Y	$50.00\pm20.45$	$\textbf{63.64} \pm \textbf{11.36}$	$52.27\pm 6.82$	$63.64\pm20.45$	$54.55\pm22.73$	$38.64 \pm 13.64$
Mistral 0.2	7B	Y	$54.55 \pm 22.73$	$59.09 \pm 14.77$	$50.00\pm11.36$	$63.64\pm20.45$	$68.18 \pm 18.18$	$43.18\pm12.50$
Llama 2	13B	Ν	$59.09 \pm 20.45$	$61.36\pm13.64$	$\textbf{70.45} \pm \textbf{13.64}$	$68.18 \pm 18.18$	$\textbf{77.27} \pm \textbf{18.18}$	$2.27\pm3.41$
Llama 2	13B	Y	$63.64 \pm 20.45$	$56.82 \pm 13.64$	$59.09 \pm 11.36$	$59.09\pm20.45$	$63.64 \pm 18.18$	$13.64\pm12.50$
Mixtral 0.1	8x7B	Ν	$59.09 \pm 20.45$	$\textbf{72.73} \pm \textbf{14.77}$	$\textbf{65.91} \pm \textbf{9.09}$	$\textbf{72.73} \pm \textbf{18.18}$	$68.18 \pm 18.18$	$0.00\pm0.00$
Mixtral 0.1	8x7B	Y	$45.45\pm20.45$	$61.36 \pm 15.91$	$\textbf{65.91} \pm \textbf{9.09}$	$\textbf{72.73} \pm \textbf{18.18}$	$68.18 \pm 18.18$	$52.27 \pm 13.64$
Llama 2	70B	Ν	$63.64 \pm 20.45$	$59.09 \pm 13.64$	$54.55\pm11.36$	$\textbf{72.73} \pm \textbf{18.18}$	$\textbf{72.73} \pm \textbf{18.18}$	$0.00\pm0.00$
Llama 2	70B	Y	$63.64 \pm 20.45$	$\textbf{68.18} \pm \textbf{14.77}$	$56.82\pm9.09$	$\textbf{72.73} \pm \textbf{18.18}$	$68.18 \pm 18.18$	$9.09 \pm 10.23$
Claude 2	-	Y	-	-	-	-	-	$47.73\pm13.64$
GPT 3.5 T	-	Y	-	$\textbf{65.91} \pm \textbf{13.64}$	$63.64 \pm 13.64$	$59.09 \pm 20.45$	$\textbf{72.73} \pm \textbf{18.18}$	$\textbf{63.64} \pm \textbf{11.36}$
GPT 4	-	Y	-	$61.36\pm11.36$	$59.09 \pm 13.64$	$\textbf{86.36} \pm \textbf{13.64}$	$\textbf{81.82} \pm \textbf{15.91}$	$\textbf{72.73} \pm \textbf{10.23}$

Table 14: Results for the 22 example pairs with the Question tag. Scores with 95% confidence intervals above chance are shown in **bold**.