# 002 003 004 005 006 007 008 009 010 011 012 013 014 015 016 017 020 021 022 023 024 025 026

037

041

001

# Simple Linguistic Inferences of Large Language Models (LLMs): Blind Spots and Blinds

## **Anonymous ACL submission**

#### Abstract

We evaluate LLMs' language understanding capacities on simple inference tasks that most humans find trivial. Specifically, we target (i) grammatically-specified entailments, (ii) premises with evidential adverbs of uncertainty, and (iii) monotonicity entailments. We design evaluation sets for these tasks and conduct experiments in both zero-shot and chainof-thought setups, and with multiple prompts and LLMs. The models exhibit moderate to low performance on these evaluation sets. Subsequent experiments show that embedding the premise in syntactic constructions that should preserve the entailment relations (presupposition triggers) or change them (non-factives), further confuses the models, causing them to either under-predict or over-predict certain entailment labels regardless of the true relation, and often disregarding the nature of the embedding context. Overall these results suggest that, despite LLMs' celebrated language understanding capacity, even the strongest models have blindspots with respect to certain types of entailments, and certain information-packaging structures act as "blinds" overshadowing the semantics of the embedded premise.

# 1 Introduction

LLMs have gained immense popularity thanks to their unprecedented ability to understand user queries and generate fluent seemingly-human responses. At the same time, people constantly report LLMs' failures, anecdotal (Borji, 2023) and systematic, e.g, the lack of reliability and consistency (Shen et al., 2023; Jang and Lukasiewicz, 2023; Plevris et al., 2023), contradictory or unreasonable answers (Zhong et al., 2023), inability to detect false assumptions (Shen et al., 2023), wrong information in prompts (Zuccon and Koopman, 2023), contradictory responses to identical queries (Jang and Lukasiewicz, 2023; Plevris et al., 2023).

However, humans are prone to some failures as

well, e.g., overlooking false assumptions in questions beyond their area of expertise, or failing to find the correct solution to a math problem.

042

043

044

045

046

047

054

057

059

061

062

063

064

065

067

068

069

070

071

073

074

075

076

077

In this work we focus on tasks that are trivial for humans, and do not require any specialized expertise beyond proficiency in English. For example, it is obvious to a human that *Her brother was singing* entails *Someone was singing*, and *Fred's tie is very long* implies *Fred's tie is long*, but not vice versa. However, as shall be seen shortly, LLMs fail to establish such systematic relations correctly. LLMs' errors on such simple tasks are much more indicative of absence of *human-like* text understanding.

We experiment with several types of natural language inferences (NLI), (a.k.a. recognizing textual entailment (Dagan et al., 2005; Bowman et al., 2015)), that are easy for humans, and show that they pose a challenge to LLMs. These NLI tests reveal some of the models' blind spots, and indicate that they are far from a genuine *human-level* understanding. Moreover, some information-packaging structures, where a statement is embedded within a surrounding linguistic context, may act as "blinds" that hinder the semantics of embedded premises, again in contrast to human-like behavior.

A summary of the conditions we test, as well as the trends observed for one model (gpt-3.5-turbo-0301) is available in Figure 1.<sup>1</sup>

We focus on inference types that are solely based on common linguistic phenomena and "trival" world-knowledge such as class membership ("a dog is an animal", "navy blue is a shade of blue"). Specifically, we test LLMs' ability to make three inference types: (i) *Grammatically-specified entailments*, i.e. replacing a constituent of the premise with an indefinite pronoun as *somebody* or *something*. (ii) Premises with *evidential adverbs of un-*

<sup>&</sup>lt;sup>1</sup>Results for other models are presented in Table 2 and discussed in Section 3. While directions and magnitudes differ, all models are far from human level, and do not respond correctly to embedding environments.

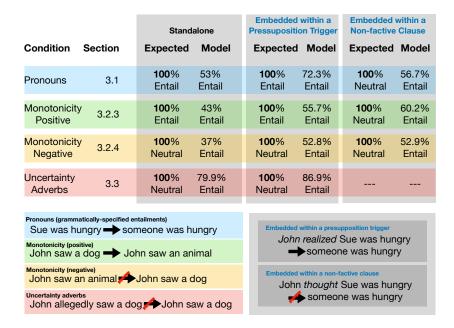


Figure 1: High-level summary of the experiments and results (reported numbers are for gpt-3.5-turbo-0301).

certainty (supposedly, allegedly etc.), that block the entailment of the rest of the clause, and (iii) Monotonicity entailment (see MacCartney and Manning (2008)) of two kinds: upward, i.e. from subsets to supersets ("Jack is a dog" entails "Jack is an animal"), and downward, i.e. from supersets to subsets ("Jack isn't an animal" entails "Jack isn't a dog"). We manually curate test sets for these inference types and experiment with them in a zero-shot setup, observing that LLMs struggle with these phenomena, leading to low accuracy.

090

091

100

101

102

103

104

105

106

107

108

We next check how embedding of the premise in a larger grammatical context affects the prediction. Such embedding can take several forms. Contexts consisting of presupposition triggers (e.g. He realized that [...], They were glad that [...], Something happened before [...]) serve to strengthen the embedded premise, while similarly structured non-factives (e.g. I feel that [...], He imagined that [...]) may cancel it. We experiment with both context types and show that in most cases they affect the LLMs' predictions incorrectly. E.g., ChatGPT<sup>2</sup> has a hard time discerning the two cases, incorrectly treating both as hints towards entailment (for regular prompting) or against it (for chainof-thought prompting). These or similar trends are observed across different prompts and models (GPT-3.5, GPT-4, LLaMA 2). While GPT-4 clearly improves over other LLMs in terms of accuracy and shows some more sound trends, it still leaves

much room for improvement compared to the human baseline.

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

These results demonstrate that state-of-the-art LLMs were unable to learn simple linguistic inferences that humans find trivial: they did not acquire them automatically in pre-training, and also in the process of instruct-tuning or human-feedback tuning. Persistence of the problem across prompts and LLMs implies that this is a systematic issue.

While current LLMs struggle with these trivial phenomena, we hope future models will show improvements. Our methodology and the dataset we release<sup>3</sup> can help track this progress.

## 2 Linguistic phenomena considered

We focus on the following linguistic phenomena.

Gramatically-specified entailments The set of the entailments of any sentence includes so-called grammatically-specified entailments (Wilson and Sperber, 1979), i.e., entailments where a constituent of the premise is substituted with a variable (such as an indefinite pronoun like somebody, something etc.). For instance, the entailments of "You've eaten all my apples" include, among others:

You've eaten all someone's apples.
You've eaten all of something.
You've eaten something.

<sup>2</sup>https://openai.com/chatgpt

<sup>&</sup>lt;sup>3</sup>URL will be provided in camera ready

You've done something.

Someone's eaten all my apples.

**Monotonicity entailments** hold when less specific predicates are substituted with more specific ones, or vice versa. They can be of two types:

• Upward: more specific predicates can be substituted with less specific ones.

*Jack is a dog.*  $\models$  *Jack is an animal.* 

• Downward: less specific predicates can be substituted with more specific ones.

All animals need water.  $\models$  All dogs need water.

Evidential Adverbs "express degrees of certitude with respect to the speaker's subjective perception of the truth of a proposition" (Haumann, 2007). We test LLMs' ability to understand evidential adverbs expressing *uncertainty* (allegedly, purportedly, supposedly etc.). Introducing such adverbs into a clause cancels the entailment of the rest of the clause. E.g., Mike allegedly worked all night does not entail Mike indeed worked all night. The relation between the two statements is neutral.

Presuppositions and Presupposition Triggers Presupposition (Beaver et al., 2021; Jeretic et al., 2020; Parrish et al., 2021) is a type of inference "whose truth is taken for granted in the utterance of a sentence" (Huang, 2011). Below, a presupposes **b** (i.e. if **b** is false, a cannot be felicitously uttered):

- **a**. Jane returned to New York.
- $\models$  **b**. Jane has been to New York before.

Presuppositions are not presented as at-issue content of the utterance, but rather as part of the background, mutually known or assumed by the speaker and the hearer (even if in reality it is not the case). The speaker of **a** does not *inform* the hearer that Jane has been to New York before: she *assumes* it; and if the hearer does not know it, she *accommodates* it upon hearing the utterance (Fintel, 2008).

Presuppositions are normally evoked by constructions or lexical items, called *presupposition triggers* (Karttunen, 2016). In sentence **a** above, the presupposition is triggered by the verb *returned*, from the class of *iterative verbs* which presuppose that the action has happened before. Other iterative verbs are *relearn*, *reread*, *reapply* etc. Presupposition triggers used in this work are factives, temporal and other adverbial clauses and embedded *wh*-questions.

Non-factive Verbs and Expressions (Kiparsky and Kiparsky, 1970), such as *believe*, *claim*, *feel*, *hope*, *suspect*, *think*, **do not entail either truth or falsity of their complements**. For example, given:

- a. Jane thinks that Bill bought bread.
- **b**. Bill bought bread.
- c. Bill didn't buy bread.

Sentence  $\mathbf{a}$  does not entail either  $\mathbf{b}$  or  $\mathbf{c}$ . The relation between  $\mathbf{a}$  and  $\mathbf{b}$  is neutral, and so is the relation between  $\mathbf{a}$  and  $\mathbf{c}$ .

# Presupposition Triggers, Non-Factives and NLI

It is important to note that embedding a premise under a presupposition trigger does not affect the relations between the premise and hypothesis. By contrast, if we embed the premise under a nonfactive, the relation becomes neutral. For example:

- **a**. A balloon hit a light post.
- $\models$  **b**. Something hit a light post

Premise **a** above entails hypothesis **b**. If we embed premise **b** under a presupposition trigger as in **a**?:

- a'. She realized that a balloon hit a light post
- **b**. Something hit a light post

the relation does not change: the new premise **a'** still entails **b**. However, when embedding premise **a** under a non-factive verb:

- a". I suspect a balloon hit a light post
- $\not\models$  **b**. Something hit a light post

the relation becomes neutral: without additional context the new premise **a**" does not entail **b**.

## 3 Main Experiments

Benchmark Design. We construct our benchmark in the form of several groups of NLI pairs, where each group focuses on a single linguistic phenomenon, and all the items within a group share the same label. This allows us to isolate the models' performance on the different conditions, rather than just providing aggregate accuracy numbers, which hide the underlying linguistic trends.

**Zero-shot Setup.** We deliberately choose a zeroshot setup rather than in-context learning, for the following reasons. 1) All tested patterns are easily learnable from examples due to their consistent structure within each group, so predictions based on such learning would be a form of "parroting" rather than "understanding". 2) It is zeroshot capabilities that indicate if an LLM is truly a general-purpose system (Qin et al., 2023). In real-world scenarios, the goal is not necessarily to handle these inferences directly, but rather to use them as a means to an end. Hence, zero-shot capability becomes crucial. 3) Finally and most importantly, our goal was not to answer "can we make an LLM perform well on these simple cases through prompting" but rather to answer "is this semantic picked up by an LLM automatically as part of its training, from text alone" as some works seemingly suggest (e.g., Merrill et al., 2022).

223

231

235

236

240

241

242

245

247

249

251

254

255

259

260

262

263

**Models.** We test a range of OpenAI models: GPT-3.5 (text-davinci-003), ChatGPT (gpt-3.5-turbo-0301 and gpt-3.5-turbo-0613), and GPT4 (gpt-4-0314 and gpt-4-0613). We access these LLMs through OpenAI's API<sup>4</sup> with the default settings<sup>5</sup>, using a single prompt which receives two texts and asks if, given text 1, text 2 is true, false or neutral (see Appendix B for details).<sup>6</sup>

We also evaluate the 70B LLaMA-2 Chat model (Touvron et al., 2023).<sup>7</sup>

Data Collection Methodology. For monotonicity entailment, we used a sample from the MED dataset (Yanaka et al., 2019a) (see §3.2.2, §3.2.3). For grammatically-specified entailments and uncertainty adverbs, the data was created by the first author by manually collecting naturally-occurring sentences from web-pages retrieved by random keywords and selecting a subset of sentences in the page that fit our desired criteria (e.g., for grammatically-specified entailment - sentences containing noun phrases in different positions that could be substituted by indefinite pronouns), and potentially re-writing them (for example, remov-

ing embedding contexts) to obtain simple sentences. We did not consult an LLM during this process. We then derived premise-hypotheses pairs from these sentences as described in §3.2.1 and §3.2.4. For the embedding contexts, we manually curated a set of candidate environments by collecting naturally occurring presupposition triggers and non-factives from textual material on the same web-pages.

264

265

266

267

269

270

271

272

273

274

275

276

277

278

279

281

282

285

289

290

291

292

293

294

295

296

297

298

300

301

302

303

304

305

306

307

308

The labels were assigned based on the linguistic properties of the phenomena (e.g. entailment-cancelling property of uncertainty adverbs). The resulting labels for each data group are confirmed by the majority vote in our human annotation experiment (see Section 3.1) in 99.8% of the cases and by all three annotators unanimously in 95.8% of the cases<sup>8</sup> with a fleiss kappa of 0.94, which corresponds to almost perfect agreement, attesting to the high quality of the data.

## 3.1 Human Performance

To verify that the inference types we consider are indeed trivial for humans, we sample 100 examples from each of the 11 data types detailed in Figure 1 and have each example labeled by three human annotators from Amazon Mechanical Turk, using majority vote to obtain the final label. We obtain an overall accuracy of 97.55%, and a full agreement among the annotators for 91% of the examples. These very high numbers confirm that these inference types pose no challenge for humans. For results grouped by inference type see Table 2. More details are available in Appendix A.

# 3.2 Testing LLMs in an NLI setting

Table 2 (top) summarizes the results. For brevity, we include only one ChatGPT and one GPT-4 version in the table. Results for the other versions can be found in Appendix C.

# 3.2.1 Grammatically-Specified Entailments

**Data:** We manually curated a dataset of 100 pairs with grammatically specified entailments replacing an NP in each premise with an indefinite pronoun in the hypothesis (see Table 1, example (1)). This is a seemingly very easy dataset, trivial for any human, where all gold labels are "ENTAILMENT".

**Results:** All LLMs except GPT4 perform poorly on these items (maximum 53% accuracy). GPT-4

<sup>4</sup>https://openai.com/product

<sup>&</sup>lt;sup>5</sup>For text-davinci-003 we set the temperature at 0.

<sup>&</sup>lt;sup>6</sup>This prompt yields a 71% accuracy on 300-instance sample from the SNLI dataset, consistent with previous SNLI prompting literature (Qin et al., 2023; Wang et al., 2023; Jang and Lukasiewicz, 2023).

<sup>&</sup>lt;sup>7</sup>We access LLaMA 2 through the Replicate API: https://replicate.com. Since the prompt used for the OpenAI models yields very poor accuracy (39%) on SNLI for LLaMA 2, we used a different prompt template, with SNLI accuracy of 61%. We use temperature of 0.01 and top-k=1.

<sup>&</sup>lt;sup>8</sup>These numbers differ from the ones reported in §3.1, because here we do not consider the monotonicity data, which was taken from an existing dataset.

<sup>&</sup>lt;sup>9</sup>All the results reported in this paper are from a single run.

Inference type	Standalone	Under presupposition triggers (premise modified; label unmodified)	Under non-factives (premise modified; label - always neutral)
Grammatically- specified entailments	(1) <b>P:</b> Crown Princess Mary of Denmark has given birth to a healthy baby boy. <b>H:</b> Someone has given birth to a healthy baby boy. <b>L:</b> entailment	(5) <b>P:</b> We are happy that Crown Princess Mary of Denmark has given birth to a healthy baby boy. <b>H:</b> Someone has given birth to a healthy baby boy. <b>L:</b> entailment	(9) <b>P:</b> I hope Crown Princess Mary of Denmark has given birth to a healthy baby boy. <b>H:</b> Someone has given birth to a healthy baby boy. <b>L:</b> neutral
Monotonicity positive	(2) <b>P:</b> She planted blue and purple pansies in the flower bed. <b>H:</b> She planted pansies in the flower bed. <b>L:</b> entailment	(6) <b>P:</b> After she planted blue and purple pansies in the flower bed, she started planting other flowers. <b>H:</b> She planted pansies in the flower bed. <b>L:</b> entailment	<ul> <li>(10) P: I think she planted blue and purple pansies in the flower bed.</li> <li>H: She planted pansies in the flower bed.</li> <li>L: neutral</li> </ul>
Monotonicity negative	(3) P: Susan made a dress for Jill. H: Susan made a long dress for Jill. L: neutral	(7) P: They are aware that Susan made a dress for Jill. H: Susan made a long dress for Jill. L: neutral	(11) P: They believe Susan made a dress for Jill. H: Susan made a long dress for Jill. L: neutral
Adverbs	(4) <b>P:</b> These persons were allegedly inhabiting the home. <b>H:</b> These persons were inhabiting the home. <b>L:</b> neutral	(8) P: The owner was aware that these persons were allegedly inhabiting the home. H: These persons were inhabiting the home. L: neutral	N/A

Table 1: Inference types (P - premise, H - hypothesis, L - label)

versions handle them much better, with 84% and 100% (human-level) accuracy respectively.

## 3.2.2 Positive Monotonicity Entailment

**Data:** We sample 100 positive examples (the gold label is "ENTAILMENT") from the Monotonicity Entailment Dataset (MED) (Yanaka et al., 2019a) - see example (2) in Table 1.

**Results:** Except for GPT-4, all models yield very low accuracy (far below 50). The highest result (by gpt-4-0613) of 70% accuracy is still far below human performance.

# 3.2.3 Negative Monotonicity Entailment

**Data:** We sample 100 negative examples (the gold label is "NEUTRAL") from MED (Yanaka et al., 2019a) - see example (3) in Table 1.

**Results:** Here, the *earlier* GPT-4 version (gpt-4-0314) only slightly underperforms humans (88%). All other models (including gpt-4-0613) have much lower accuracy: 28%- 56%. <sup>10</sup>

# 3.2.4 Adverbs of Uncertainty

**Data:** We manually create a dataset of 100 sentence pairs where the premise contains an uncertainty adverb, while the hypothesis omits it (exam-

ple (4) in Table 1). We apply 9 adverbs (*allegedly*, *hopefully*, *possibly*, *presumably*, *probably*, *purportedly*, *reportedly*, *seemingly*, *supposedly*) to each of the 100 pairs, obtaining 900 examples, 100 per adverb. The gold label for all pairs is "NEUTRAL".

**Results:** Strongest result (70% accuracy) is obtained by the earlier GPT-4 (0314), and is substantially below human performance. All other models, including the newer GPT-4, achieve very low accuracy on this data (4.67% - 38%).

Bottomline: Except for GPT-4, all tested models perform very poorly on all the explored inference types. In many cases the LLMs (especially LLaMA 2) predict mostly the opposite of the expected label ("ENTAILMENT" where "NEUTRAL" is expected and vice versa). While the two GPT-4 versions improve significantly over other LLMs, neither handles all the inference types equally well: the later version (gpt-4-0613) is very good at handling the cases where "ENTAILMENT" is expected, but its predecessor (gpt-4-0314) outperforms it on "NEUTRAL" cases.

# 3.3 Embedding the Premises under Presupposition Triggers

Table 2 (middle) summarizes the results (except for gpt-3.5-turbo-0613 and gpt-4-0314 covered in

<sup>&</sup>lt;sup>10</sup>Similar results on MED are reported by Liu et al. (2023).

Appendix C).

**Data:** Next we modify all the datasets described above by embedding the original premises under presupposition triggers (see examples (5)–(8) in Table 1). We use 23 trigger types: 21 factives, such as *know*, *realize*, *be glad*; embedded wh-questions, such as *This explains* why he came; and adverbial clauses, e.g. After he came, I cooked dinner.

For uncertainty adverbs, we randomly sample 100 examples from our dataset of 900 sentence pairs, and apply each of the 23 triggers to each sampled pair, obtaining 2300 pairs. For other data types (grammatically specified entailments, positive and negative monotonicity entailments) we embed the 100 original premises under 23 triggers, which also results in 2300 pairs per data type. The hypotheses remain unchanged (see Table 1).<sup>11</sup>

Presupposition triggers do not change the relation between the premise and the hypothesis, so the gold label remains unchanged: "ENTAILMENT" for grammatically specified and positive monotonicity entailments; "NEUTRAL" for negative monotonicity entailments and uncertainty adverbs.

**Results:** The OpenAI models show the same pattern across all data types: under presupposition triggers they predict "ENTAILMENT" more *regardless of the correct label*. This improves accuracy where "ENTAILMENT" is expected (grammatically specified and positive monotonicity entailments), but reduces accuracy for "NEUTRAL" data (negative monotonicity, uncertainty adverbs). <sup>12</sup>

LLaMA 2 shows a "healthier" trend: under presupposition triggers it seems to capture all inference types better than in the standalone experiments, predicting more entailment for the "positive" datasets and more "neutral" relations for the "neutral" datasets. However, its accuracy remains consistently low (below 40%) across all data types. Additionally, under presupposition triggers, it consistently predicts more "contradiction" labels for all inference types, which is always wrong.

**Bottomline:** Under presupposition triggers the OpenAI LLMs tend to predict more entailment irrespective of the correct label. LLaMA 2 improves

in this setting, but shows a different wrong trend: predicting more "contradiction" labels in all cases.

# 3.4 Embedding the Premises under Non-Factives

Table 2 (bottom) summarizes the results (gpt-3.5-turbo-0613 and gpt-4-0314 are covered in Appendix C).

**Data:** We modify the original datasets (see §3.2), embedding the premises under non-factives (see examples (9)–(11) in Table 1). For grammatically-specified entailments and both types of monotonicity data, we embed each of the 100 original premises under 23 different non-factives, e.g. *feel, hope, believe* etc., obtaining 2300 pairs for each data type. We omit this experiment for adverbs for semantic reasons: including both an uncertainty adverb and a non-factive into the premise (*I guess he allegedly worked all night.*) results in double expression of uncertainty, creating a tautology.

The hypotheses remain unchanged (see Table 1). Since non-factives cancel entailment, the gold label in all the cases changes to "NEUTRAL".

**Results:** Surprisingly, GPT-3.5 and ChatGPT<sup>13</sup> again predict "ENTAILMENT" more often for all the data groups. <sup>14</sup> This hurts the accuracy for all data groups, since now the correct label is always "NEUTRAL". Both GPT-4 versions show the *correct* trend predicting *less* "ENTAILMENT" and *more* "NEUTRAL" labels under non-factives. Yet, the now-incorrect "ENTAILMENT" label remains very frequent reducing the accuracy.

LLaMA 2 again (as with presupposition triggers) seems more accurate compared to standalone experiments, but a closer look reveals it is likely "right for the wrong reason": the proportion of neutral predictions has changed very little compared to the standalone experiments (maximum by 7%), but now the "NEUTRAL" label is always correct, which seemingly boosts the accuracy.

<sup>&</sup>lt;sup>11</sup>In experiments with both embedding context types we use full data for text-davinci-003, gpt-3.5-turbo-0301 and LLaMA 2, and 100 examples per inference type for other models.

<sup>&</sup>lt;sup>12</sup>The only exception is gpt-3.5-turbo-0613: under presupposition triggers it predicts less entailment (except for negative monotonicity), which in most cases hurts the accuracy (see Appendix C).

<sup>&</sup>lt;sup>13</sup>The newer ChatGPT version, gpt-3.5-turbo-0601, shows the same trends as under presupposition triggers, predicting less entailment (except for negative monotonicity) (see Appendix C). It shows that, just like its predecessors, this LLM is unable to distinguish the two types of embedding contexts.

<sup>&</sup>lt;sup>14</sup>This inclination to overpredict entailment might imply the use of lexical overlap heuristics, given the lexical similarity between premises and hypotheses in our data. However, a closer examination reveals a contrary trend: standalone premises have greater lexical overlap with hypotheses than embedded ones, yet the models predict more entailment for embedded premises. Our evaluation on a subset of HANS (McCoy et al., 2019) also suggest that lexical overlap does not significantly influence ChatGPT towards entailment.

## 3.4.1 Bottomline

Only the GPT-4 models capture the overall trend correctly reducing entailment predictions and increasing neutral predictions under non-factives. However, while, ideally, the entailment predictions should be now reduced to 0, they still remain quite frequent, especially in the supposedly more advanced gpt-4-0613, reducing the models' accuracy.

Across all the experiments described above, GPT-4 models show the best results and the most human-like trends. They are also the only ones consistantly refraining from the always incorrect "CONTRADICTION" label. However, the earlier version (see Appendix C) proves biased towards neutral predictions, while the later one constantly leans towards entailment. Under presupposition triggers they share their predecessors' erroneous trend: predicting more entailment regardless of the correct label. In spite of the favourable dynamics both GPT-4 versions still perform far below the human level on most inference types we explored.

Notably, overall, the newer ChatGPT and GPT-4 versions prove less accurate than the respective earlier ones (see Appendix C).

# 4 Prompt Variations

We also experiment with prompt variations.

**Prompt paraphrasing.** We ask ChatGPT (gpt-3.5-turbo-0301) to rephrase our prompt template, obtaining two templates that we verify to be semantically equivalent w.r.t. the task (details in Appendix D). While accuracies vary between prompts, the scores on all tasks remain low.

Chain-of-Thought Prompting (CoT). Using gpt-3.5-turbo-0301, we investigate CoT prompting (Kojima et al., 2023), and find that it *reverses* the LLM's trend to predict more entailment for embedded premises (see Section 3): now the LLM predicts *less* entailment under presupposition triggers or non-factives, while the number of neutral predictions grows - again, regardless of the correct label, i.e. the reversed trend is equally inaccurate. CoT prompting improves the accuracy, but in a one-sided way: scoring higher on "neutral" test sets, but lower on almost all the "entailment" ones.

Analysis of CoT Results. The CoT technique allows us to explore the model's "reasoning". We manually evaluate a subset of the CoT explanations. In half of the cases (50.9%) both the final decision and the CoT explanation were wrong. In 23.6% a

correct explanation was followed by a correct decision; in 23.6% a wrong explanation was followed by a correct decision. In 1.86% of the cases a correct explanation was followed by a wrong decision. 81% of the cases reflected a correct understanding of the task expressed in the prompt. In half of the cases (49.1%) the CoT mentioned the undelying linguistic phenomena explicitly, but only in half of those (23.6% of the total) reflected their correct understanding and only 14.5% of the time used them as a basis for the final prediction.

The details of the CoT experiments and the manual analysis are available in Appendix E.

# 5 Relation to previous findings

A vast literature evaluated linguistic abilities of earlier generations of pretrained models<sup>15</sup>, including some phenomena we consider here: Jiang and de Marneffe (2019); de Marneffe et al. (2019); Ross and Pavlick (2019); McCoy et al. (2019); Yanaka et al. (2021) (embedded clauses); Jeretic et al. (2020); Parrish et al. (2021); Kabbara and Cheung (2022) (presuppositions and implicatures); Yanaka et al. (2019a,b); Goodwin et al. (2020); Yanaka et al. (2020); Jumelet et al. (2021); Geiger et al. (2020) (monotonicity entailment); Ettinger (2020); Guerin and Chemla (2023) (hypernymy). Some works focused on other phenomena: negation, synonyms, antonyms, tenses, temporal adverbials, active vs. passive voice, symmetric and asymmetric relations, comparatives, superlatives etc. (e.g., Kassner and Schütze, 2020; Hartmann et al., 2021; Hossain et al., 2020, 2022; Ettinger, 2020; Truong et al., 2022; Jang et al., 2022; Ravichander et al., 2022; Ribeiro et al., 2020). Many of these studies make use of the NLI task to evaluate LMs' basic linguistic abilities (e.g., McCoy et al., 2019; Jiang and de Marneffe, 2019; Yanaka et al., 2021; Kabbara and Cheung, 2022; Jeretic et al., 2020; Parrish et al., 2021; Kabbara and Cheung, 2022; Yanaka et al., 2019a,b). We draw inspiration from this line of work.

In the remainder of the section, we focus on work that evaluates SOTA LLMs (starting with GPT-3).

Truong et al. (2023) seek to determine if **negation** remains a challenge in the age of LLMs. <sup>16</sup> They assess GPTneo, GPT-3, and InstructGPT

<sup>&</sup>lt;sup>15</sup>There are comprehensive overviews of such literature, for example, Rogers et al. (2020); Chang and Bergen (2023).

<sup>&</sup>lt;sup>16</sup>In fact, their work covers a range of negation-related phenomena, such as synonym/antonym lexical relations, monotonicity entailment and hypernymy under negation.

		human			GPT-	3.5		ChatGPT		GPT-4			LLaMA 2 70b			
		e	n	c	e	n	c	e	n	c	e	n	С	e	n	c
	pronouns	100.0	0.0	0.0	39.0	56.0	5.0	53.0	34.0	13.0	100.0	0.0	0.0	31.0	61.0	8.0
Standalone	monotonicity positives	89.0	11.0	0.0	25.0	46.0	29.0	43.0	51.0	6.0	<u>70.0</u>	30.0	0.0	36.0	59.0	5.0
Standarone	monotonicity negatives	5.0	94.0	1.0	22.0	28.0	50.0	37.0	42.0	21.0	33.0	<u>56.0</u>	11.0	46.0	35.0	19.0
	uncertainty adverbs	0.0	100.0	0.0	77.6	4.7	17.8	79.9	9.1	11.0	61.0	<u>38.0</u>	1.0	89.9	9.0	1.1
	pronouns	100.0	0.0	0.0	65.4	27.8	6.8	72.3	23.3	4.3	<u>100.0</u>	0.0	0.0	32.0	58.1	9.9
Under presupposition	monotonicity positives	97.0	3.0	0.0	38.0	25.1	36.8	55.7	37.2	7.1	<u>71.0</u>	29.0	0.0	36.9	44.9	18.3
triggers	monotonicity negatives	5.0	95.0	0.0	39.8	19.0	41.2	52.8	29.4	17.8	36.0	<u>53.0</u>	11.0	32.4	38.5	29.0
	uncertainty adverbs (sample)	1.0	99.0	0.0	87.8	2.7	9.6	86.9	7.5	5.6	70.0	<u>30.0</u>	0.0	58.4	29.4	12.3
Under	pronouns	0.0	100.0	0.0	51.0	41.0	8.0	56.7	40.0	3.3	97.0	3.0	0.0	29.7	<u>64.7</u>	5.5
non-factives	monotonicity positives	1.0	99.0	1.0	28.8	27.3	43.9	60.2	30.2	9.7	44.0	<u>56.0</u>	0.0	34.4	52.3	13.3
non-ractives	monotonicity negatives	0.0	100.0	0.0	30.3	17.4	52.3	52.9	22.3	24.7	23.0	<u>69.0</u>	8.0	29.3	42.0	28.7

Table 2: Experiment results across all models. Background indicates the expected label (green - "entailment", blue - "neutral"). The green and red font colors indicate increase or decrease compared to standalone premises. The numbers in **bold** also indicate *accuracy* (e.g. for uncertainty adverbs accuracy coincides with the percentage of neutrals). The **underlined bold** numbers indicate *the best accuracy for a specific inference type* across all models (human baseline excluded). **e** - entailment (%); **n** - neutral (%); **c** - contradiction (%). Model versions shown: GPT-3.5 - text-davinci-003; ChatGPT – gpt-3.5-turbo-0301; GPT-4 - gpt-4-0613. Result for additional GPT models (gpt-3.5-turbo-0613 and gpt-4-0314) available in Appendix C.

across various negation benchmarks with different model sizes and prompts, and discover a number of important limitations, including insensitivity to negation, an inability to capture synonym/antonym lexical relations and to reason under negation. Ye et al. (2023) find that LLMs<sup>17</sup> with CoT-style prompts struggle with lexical negation (words with **negative suffixes**: *implausible*, *unreasonable* etc.). Lorge and Pierrehumbert (2023) explore the capability of LMs (including text-davinci-002) to handle scalar adverbs (specifically, adverbs of frequency, modality and degree). Despite some success they notice many shortcomings: "weak differentiation amongst the semantic classes of adverbs, poor ability to discriminate scalar adverbs from negations", "strong effects of adverb frequencies and lack of generalisation across two logically equivalent entailment constructions". Cong (2022) probes GPT-3 for understanding of scalar implicatures and presuppositions and find that its performance is "mostly at chance". Stengel-Eskin and Van Durme (2022) investigate how GPT3davinci, GPT-Neo, Jurassic, T5 for QA, T0pp handle subject-control sentences. Davis (2022) analyzes how autoregressive models with GPT-like architectures, such as GPT-2 XL, GPT-Neo, GPT-J, GPT-3 (text-davinci-002) handle binding and conclude that they consider antecedents which humans reject as structurally ungrammatical. Focusing on "language illusions", Zhang et al. (2023) find that LMs including text-davinci-003 "have some capacity to process comparative structures", but struggle with multiple negation and NPI licensors.

536

537

538

541

542

544

545

546

547

548

549

550

552

553

554

556

564

565

This overview suggests that for newer LLMs (GPT-3 and later) there is much less research focusing on fundamental linguistic abilities. We add to this line of work, focusing on very simple inference types that to the best of our knowledge were not previously evaluated (grammatically-specified entailments, uncetainty adverbs) as well as evaluating stability under embedding contexts (presupposition-triggers and non-factives).

569

570

571

572

573

574

575

576

577

578

579

580

581

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

600

At the same time, the existing research shows that current pre-training objectives do not necessarily capture the basic linguistic phenomena, thus "new training paradigms are essential to achieve better linguistic competency" (Truong et al., 2023).

## 6 Conclusion

We tested LLMs with several simple linguistic inferences that humans find trivial, and further experimented by embedding the premises within (simple) linguistic contexts that should either preserve or cancel the inference. The models performed poorly in all cases. We conclude that LLMs do not learn entailment semantics "naturally". The persistence of the issue across prompts, models and setups shows that these limitations are robust and this topic merits further systematic investigations. While the results are tested for several specific models, we do expect them to hold more generally. We believe the ability of LLMs to correctly understand simple textual inferences humans find natural is important. We provide a benchmark and methodology that can be used to track the improvement of future models.

<sup>&</sup>lt;sup>17</sup>They test 14 LLMs including GPT-4 and 3.5 and four LLaMA variants.

## 7 Limitations

As a consequence of LLMs' sensitivity to prompt phrasings, there may exist prompts that can potentially modify the reported results. At the same time we agree with Jang and Lukasiewicz (2023) who point out that "improvements with prompt design can be considered another violation of semantic consistency, because the prompts will deliver identical semantic meaning, i.e., task description".

We showed that embedding the premises under presupposition triggers or non-factives affects the models' predictions exhibiting certain patterns. However, our treatment of non-factives is non-exhaustive: for example, we did not try adding negation to the non-factives or using non-factive predicates denoting a high level of uncertainty (*I doubt, I'm skeptical*) or containing negative prefixes (*I an uncertain, I disbelieve*). It's possible that implicit or explicit negation in the embedding predicates may change the LLMs' behavior. Also, while considering presupposition triggers and non-factives, we did not consider other types of clause-embedding predicates (e.g., implicative verbs).

Finally, since LLMs undergo continuous updates, the test results presented here may vary over time, and the closed, black-box nature of the models' training details and training data limits the scope of the conclusions that can be reached. We cannot, for example, know for certain if the low performance on these inference type stems already from the base model and its inability to learn these type of inferences from raw text, or was added during instruct tuning. That said, our data and methodology for benchmarking these capabilities is model-agnostic and remains intact.

# **Ethics Statement**

Annotation To collect annotations for our human baseline (see §3.1) we used Amazon Mechanical Turk (AMT). We employed 5 annotators who showed themselves reliable in handling annotation tasks assigned by us in the past. All annotators had the following qualifications: (1) over 500 completed HITs; (2) 99% approval rate or higher; (3) native English speakers. We paid \$0.10 per example. Labeling each example normally (excluding outliers) took about 21.7 seconds which resulted in a payment of about \$16.5 per hour. Upon completing the task the annotators were given bonuses.

# **Annotation Data Collection and Usage Policy**

Annotators were notified that their annotations are meant for research purposes in the field of Natural Language Processing. The task and collected annotations were objective and included no personal information.

# References

- David I. Beaver, Bart Geurts, and Kristie Denlinger. 2021. Presupposition. In Edward N. Zalta, editor, *The Stanford Encyclopedia of Philosophy*, Spring 2021 edition. Metaphysics Research Lab, Stanford University.
- Ali Borji. 2023. A categorical archive of chatgpt failures
- Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. 2015. A large annotated corpus for learning natural language inference. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 632–642, Lisbon, Portugal. Association for Computational Linguistics.
- Tyler A. Chang and Benjamin K. Bergen. 2023. Language Model Behavior: A Comprehensive Survey. *Computational Linguistics*, pages 1–55.
- Yan Cong. 2022. Psycholinguistic diagnosis of language models' commonsense reasoning. In *Proceedings of the First Workshop on Commonsense Representation and Reasoning (CSRR 2022)*, pages 17–22, Dublin, Ireland. Association for Computational Linguistics.
- Ido Dagan, Oren Glickman, and Bernardo Magnini. 2005. The pascal recognising textual entailment challenge. In *Machine Learning Challenges Workshop*, pages 177–190. Springer.
- Forrest Davis. 2022. Incremental processing of principle B: Mismatches between neural models and humans. In *Proceedings of the 26th Conference on Computational Natural Language Learning (CoNLL)*, pages 144–156, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.
- Marie-Catherine de Marneffe, Mandy Simons, and Judith Tonhauser. 2019. The commitmentbank: Investigating projection in naturally occurring discourse. *Proceedings of Sinn und Bedeutung*, 23(2):107–124.
- Allyson Ettinger. 2020. What BERT Is Not: Lessons from a New Suite of Psycholinguistic Diagnostics for Language Models. *Transactions of the Association for Computational Linguistics*, 8:34–48.
- Kai Fintel. 2008. What is presupposition accommodation, again? *Philosophical Perspectives*, 22:137–170.

Atticus Geiger, Kyle Richardson, and Christopher Potts. 2020. Neural natural language inference models partially embed theories of lexical entailment and negation.

- Emily Goodwin, Koustuv Sinha, and Timothy J. O'Donnell. 2020. Probing linguistic systematicity.
- Nicolas Guerin and Emmanuel Chemla. 2023. It is a bird therefore it is a robin: On BERT's internal consistency between hypernym knowledge and logical words. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 8807–8817, Toronto, Canada. Association for Computational Linguistics.
- Mareike Hartmann, Miryam de Lhoneux, Daniel Hershcovich, Yova Kementchedjhieva, Lukas Nielsen, Chen Qiu, and Anders Søgaard. 2021. A multilingual benchmark for probing negation-awareness with minimal pairs. In *Proceedings of the 25th Conference on Computational Natural Language Learning*, pages 244–257, Online. Association for Computational Linguistics.
- Dagmar Haumann. 2007. Adverb Licensing and Clause Structure in English.
- Md Mosharaf Hossain, Dhivya Chinnappa, and Eduardo Blanco. 2022. An analysis of negation in natural language understanding corpora. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 716–723, Dublin, Ireland. Association for Computational Linguistics.
- Md Mosharaf Hossain, Venelin Kovatchev, Pranoy Dutta, Tiffany Kao, Elizabeth Wei, and Eduardo Blanco. 2020. An analysis of natural language inference benchmarks through the lens of negation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 9106–9118, Online. Association for Computational Linguistics.
- Yan Huang. 2011. *14. Types of inference: entailment, presupposition, and implicature*, pages 397–422. De Gruyter Mouton, Berlin, New York.
- Myeongjun Jang and Thomas Lukasiewicz. 2023. Consistency analysis of chatgpt.
- Myeongjun Jang, Frank Mtumbuka, and Thomas Lukasiewicz. 2022. Beyond distributional hypothesis: Let language models learn meaning-text correspondence. In *Findings of the Association for Computational Linguistics: NAACL 2022*, pages 2030–2042, Seattle, United States. Association for Computational Linguistics.
- Paloma Jeretic, Alex Warstadt, Suvrat Bhooshan, and Adina Williams. 2020. Are natural language inference models IMPPRESsive? Learning IMPlicature and PRESupposition. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8690–8705, Online. Association for Computational Linguistics.

Nanjiang Jiang and Marie-Catherine de Marneffe. 2019. Evaluating BERT for natural language inference: A case study on the CommitmentBank. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 6086–6091, Hong Kong, China. Association for Computational Linguistics.

- Jaap Jumelet, Milica Denic, Jakub Szymanik, Dieuwke Hupkes, and Shane Steinert-Threlkeld. 2021. Language models use monotonicity to assess NPI licensing. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 4958– 4969, Online. Association for Computational Linguistics.
- Jad Kabbara and Jackie Chi Kit Cheung. 2022. Investigating the performance of transformer-based NLI models on presuppositional inferences. In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 779–785, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.
- Lauri Karttunen. 2016. Presupposition: What went wrong? Semantics and Linguistic Theory, 26:705– 731.
- Nora Kassner and Hinrich Schütze. 2020. Negated and misprimed probes for pretrained language models: Birds can talk, but cannot fly. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7811–7818, Online. Association for Computational Linguistics.
- Paul Kiparsky and Carol Kiparsky. 1970. *FACT*, pages 143–173. De Gruyter Mouton, Berlin, Boston.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2023. Large language models are zero-shot reasoners.
- Hanmeng Liu, Ruoxi Ning, Zhiyang Teng, Jian Liu, Qiji Zhou, and Yue Zhang. 2023. Evaluating the logical reasoning ability of chatgpt and gpt-4.
- Isabelle Lorge and Janet B. Pierrehumbert. 2023. Not wacky vs. definitely wacky: A study of scalar adverbs in pretrained language models. *ArXiv*, abs/2305.16426.
- Bill MacCartney and Christopher D. Manning. 2008. Modeling semantic containment and exclusion in natural language inference. In *Proceedings of the 22nd International Conference on Computational Linguistics (Coling 2008)*, pages 521–528, Manchester, UK. Coling 2008 Organizing Committee.
- Tom McCoy, Ellie Pavlick, and Tal Linzen. 2019. Right for the wrong reasons: Diagnosing syntactic heuristics in natural language inference. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3428–3448, Florence, Italy. Association for Computational Linguistics.

William Merrill, Alex Warstadt, and Tal Linzen. 2022. Entailment semantics can be extracted from an ideal language model. In *Proceedings of the 26th Conference on Computational Natural Language Learning (CoNLL)*, pages 176–193, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Alicia Parrish, Sebastian Schuster, Alex Warstadt, Omar Agha, Soo-Hwan Lee, Zhuoye Zhao, Samuel R. Bowman, and Tal Linzen. 2021. NOPE: A corpus of naturally-occurring presuppositions in English. In *Proceedings of the 25th Conference on Computational Natural Language Learning*, pages 349–366, Online. Association for Computational Linguistics.

Vagelis Plevris, George Papazafeiropoulos, and Alejandro Jiménez Rios. 2023. Chatbots put to the test in math and logic problems: A preliminary comparison and assessment of chatgpt-3.5, chatgpt-4, and google bard.

Chengwei Qin, Aston Zhang, Zhuosheng Zhang, Jiaao Chen, Michihiro Yasunaga, and Diyi Yang. 2023. Is chatgpt a general-purpose natural language processing task solver?

Abhilasha Ravichander, Matt Gardner, and Ana Marasović. 2022. Condaqa: A contrastive reading comprehension dataset for reasoning about negation.

Marco Tulio Ribeiro, Tongshuang Wu, Carlos Guestrin, and Sameer Singh. 2020. Beyond accuracy: Behavioral testing of NLP models with CheckList. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4902–4912, Online. Association for Computational Linguistics.

Anna Rogers, Olga Kovaleva, and Anna Rumshisky. 2020. A primer in bertology: What we know about how bert works.

Alexis Ross and Ellie Pavlick. 2019. How well do NLI models capture verb veridicality? In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2230–2240, Hong Kong, China. Association for Computational Linguistics.

Xinyue Shen, Zeyuan Chen, Michael Backes, and Yang Zhang. 2023. In chatgpt we trust? measuring and characterizing the reliability of chatgpt.

Elias Stengel-Eskin and Benjamin Van Durme. 2022. The curious case of control. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 11065–11076, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton

Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, 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, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, 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, 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.

Thinh Truong, Timothy Baldwin, Trevor Cohn, and Karin Verspoor. 2022. Improving negation detection with negation-focused pre-training. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4188–4193, Seattle, United States. Association for Computational Linguistics.

Thinh Hung Truong, Timothy Baldwin, Karin Verspoor, and Trevor Cohn. 2023. Language models are not naysayers: an analysis of language models on negation benchmarks. In *Proceedings of the 12th Joint Conference on Lexical and Computational Semantics* (\*SEM 2023), pages 101–114, Toronto, Canada. Association for Computational Linguistics.

Jindong Wang, Xixu Hu, Wenxin Hou, Hao Chen, Runkai Zheng, Yidong Wang, Linyi Yang, Haojun Huang, Wei Ye, Xiubo Geng, Binxin Jiao, Yue Zhang, and Xing Xie. 2023. On the robustness of chatgpt: An adversarial and out-of-distribution perspective.

Deirdre Wilson and Dan Sperber. 1979. *Ordered Entailments: An Alternative to Presuppositional Theories*, volume 11, pages 299–323.

Hitomi Yanaka, Koji Mineshima, Daisuke Bekki, and Kentaro Inui. 2020. Do neural models learn systematicity of monotonicity inference in natural language? In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6105–6117, Online. Association for Computational Linguistics.

Hitomi Yanaka, Koji Mineshima, Daisuke Bekki, Kentaro Inui, Satoshi Sekine, Lasha Abzianidze, and Johan Bos. 2019a. Can neural networks understand monotonicity reasoning? In *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 31–40, Florence, Italy. Association for Computational Linguistics.

Hitomi Yanaka, Koji Mineshima, Daisuke Bekki, Kentaro Inui, Satoshi Sekine, Lasha Abzianidze, and

Johan Bos. 2019b. HELP: A dataset for identifying shortcomings of neural models in monotonicity reasoning. In *Proceedings of the Eighth Joint Conference on Lexical and Computational Semantics* (\*SEM 2019), pages 250–255, Minneapolis, Minnesota. Association for Computational Linguistics.

- Hitomi Yanaka, Koji Mineshima, and Kentaro Inui. 2021. Exploring transitivity in neural nli models through veridicality.
- Mengyu Ye, Tatsuki Kuribayashi, Jun Suzuki, Goro Kobayashi, and Hiroaki Funayama. 2023. Assessing step-by-step reasoning against lexical negation: A case study on syllogism. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 14753–14773, Singapore. Association for Computational Linguistics.
- Yuhan Zhang, Edward Gibson, and Forrest Davis. 2023. Can language models be tricked by language illusions? easier with syntax, harder with semantics. In *Conference on Computational Natural Language Learning*.
- Qihuang Zhong, Liang Ding, Juhua Liu, Bo Du, and Dacheng Tao. 2023. Can chatgpt understand too? a comparative study on chatgpt and fine-tuned bert.
- Guido Zuccon and Bevan Koopman. 2023. Dr chatgpt, tell me what i want to hear: How prompt knowledge impacts health answer correctness.

## A Human annotation details

In order to construct a human baseline, we sampled 1100 examples (100 per inference type) and had 3 human annotators label each example, asking if based on Text 1, Text 2 is true, false or neutral. We then use majority vote to determine the final human label for each example.

For human annotation, we used the Amazon Mechanical Turk platform. We did not use any qualification tests and filters for annotator selection, but instead employed 5 annotators who showed themselves reliable in handling annotation tasks assigned by us in the past. We paid \$0.10 per example. Labeling each example normally (excluding outliers) took about 21.7 seconds which resulted in a payment of about \$16.5 per hour.

We obtained an overall accuracy of 97.55%. In 91% of the cases all three annotators confirmed the expected labels unanimously. The instructions received by the annotators are shown in Figure 2.

# **B** Main entailment experiments prompts

The prompt below is used throughout the experiments described in Section 3 for all the OpenAI models (see Section 3).

You are given a pair of texts. Say about this pair: given Text 1, is Text 2 true, false or neutral (you can't tell if it's true or false)? Reply in one word.

Text 1: "text1"

Text 2: "text2"

For LLaMA 2 we use the following prompt:

You have two texts, and your task is to determine the truthfulness of Text 2 based on Text 1. Provide a one-word response indicating whether Text 2 is true, false, or neutral (indeterminable). Here are the texts:

Text 1: "text1"

Text 2: "text2"

The model outputs one of three possible labels: "true" (corresponding to "entailment"), "false" (corresponding to "contradiction") or "neutral". 18

# C Complete Results Table

All the results from all LLMs including both Chat-GPT and GPT-4 versions are shown in Table 3. Accuracy across all tested LLMs obtained in the experiments in Section 3 is compared in Table 4.

# D Prompt Paraphrasing

We experimented with the following prompt variations: (Original prompt)

You are given a pair of texts. Say about this pair: given Text 1, is Text 2 true, false or neutral (you can't tell if it's true or false)? Reply in one word.

*Text 1: "text1" Text 2: "text2"* 

(Variation 1)

You have two texts, and your task is to determine the truthfulness of Text 2 based on Text 1. Provide a one-word response indicating whether Text 2 is true, false, or neutral (indeterminable). Here are the texts:

Text 1: "text1"

Text 2: "text2"

(Variation 2)

Assess the veracity of Text 2 based on Text 1: Is Text 2 true, false, or indeterminable? Provide a one-word response.

Text 1: "text1"

Text 2: "text2"

The comparison between the three prompts is shown in Table 5.

# E Chain-of-thought prompting

## **E.1** Chain-of-thought experiments

For the chain-of-thought experiments, we added the phrase "Let's think step by step" (Kojima et al., 2023) to the original prompt as follows:

You are given a pair of texts. Say about this pair: given Text 1, is Text 2 true, false or neutral (you can't tell if it's true or false)?

Text 1: "text1"

Text 2: "text2"

Let's think step by step.

<sup>&</sup>lt;sup>18</sup>In the rare cases when the model outputs a different label, we normalize it to one of the three expected forms. E.g. "truthful" is normalized to "true".

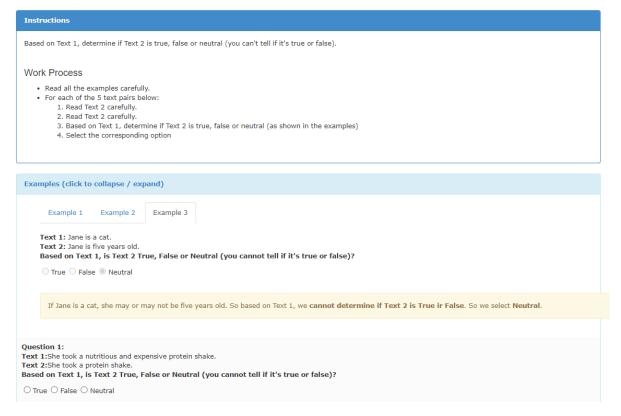


Figure 2: Instructions for human annotation. Note that in order to not implicitly train annotators towards the linguistic inferences we consider in this work, we provided examples that demonstrate the meaning of the neutral/entailing/contradiction labels, but on other inference types, not covered in this work.

		human		n text-davinci-003		gpt-3.5-turbo-0301		gpt-3.5-turbo-0613		gpt-4-0314			gpt-4-0613			LLaMA 2 70b						
		e	n	c	e	n	c	e	n	c	e	n	c	e	n	c	e	n	c	e	n	c
	P	100	0	0	39	56	5	53	34	13	9	53	38	84	16	0	100	0	0	31	61	8
SA	MP	89	11	0	25	46	29	43	51	6	21	35	44	39	61	0	<u>70</u>	30	0	36	59	5
SA	MN	5	94	1	22	28	50	37	42	21	12	37	51	6	88	6	33	56	11	46	35	19
	UA	0	100	0	77.56	4.67	17.78	79.89	9.11	11.00	34	19	47	30	<u>70</u>	0	61	38	1	89.89	9	1.11
	P	100	0	0	65.35	27.83	6.82	72.30	23.35	4.35	6	42	52	90	10	0	100	0	0	32.00	58.13	9.87
PT	MP	97	3	0	38.04	25.13	36.83	55.65	37.22	7.13	17	40	43	39	61	0	<u>71</u>	29	0	36.87	44.87	18.26
FI	MN	5	95	0	39.78	19.04	41.18	52.83	29.39	17.78	21	24	55	9	<u>87</u>	4	36	53	11	32.43	38.52	29.05
	UA	1	99	0	87.78	2.65	9.57	86.91	7.52	5.57	28	25	47	37	<u>63</u>	0	70	30	0	58.39	29.35	12.26
	P	0	100	0	51.00	41.00	8.00	56.65	40.04	3.31	4	49	47	55	45	0	97	3	0	29.74	64.74	5.52
NF	MP	1	99	1	28.83	27.26	43.91	60.17	30.17	9.66	13	37	50	16	<u>84</u>	0	44	56	0	34.35	52.35	13.30
	MN	0	100	0	30.30	17.43	52.26	52.91	22.35	24.74	15	23	62	4	<u>94</u>	2	23	69	8	29.30	41.96	28.74

Table 3: Experiment results across all models. Background indicates the expected label (green - "entailment", blue - "neutral"). The green and red font colors indicate increase or decrease compared to standalone premises. The numbers in **bold** also indicate *accuracy* (e.g. for uncertainty adverbs accuracy coincides with the percentage of neutrals). The **underlined bold** numbers indicate *the best accuracy for a specific inference type* across all models (human baseline excluded). **SA** - standalone premises; **PT** - premises under presupposition triggers; **NF** - premises under non-factives; **P** - pronouns (grammatically-specified-entailments); **MP** - monotonicity entailment (positive); **MN** - monotonicity entailment(negative), **UA** - uncertainty adverbs; **e** - entailment (%); **n** - neutral (%); **c** - contradiction (%).

As can be seen, we 1) removed the requirement to return a one-word answer; 2) added the words "Let's think step by step" at the end.

After the model outputs a chain of thought, an additional step is needed to obtain a final one-word answer. For this *answer extraction* step we use an additional prompt:

Therefore, the one-word answer (True,

False or Neutral) is

For the CoT experiments with standalone premises we use the same 100-example test sets as for the original experiments (see Section 3 for details). For experiments with embeddings we sample 100 sentence pairs out of each 2300-example test set.

The results of the original experiments (see Sec-

	Inference type	human	LLaMA 2 70b	text-davinci-003	gpt-3.5-turbo-0301	gpt-3.5-turbo-0613	gpt-4-0314	gpt-4-0613
	pronouns	100.00	31.00	39.00	53.00	9.00	84.00	100.00
Standalone	monotonicity positive	89.00	36.00	25.00	43.00	21.00	39.00	70.00
Standarone	monotonicity negative	94.00	35.00	28.00	42.00	37.00	88.00	56.00
	uncertainty adverbs	100.00	9.00	4.67	9.11	19.00	70.00	38.00
Under	pronouns	100.00	32.00	65.35	72.30	6.00	90.00	100.00
presupposition	monotonicity positive	97.00	36.87	38.04	55.65	17.00	39.00	71.00
triggers	monotonicity negative	95.00	38.52	19.04	29.39	24.00	87.00	53.00
triggers	uncertainty adverbs	99.00	29.35	2.65	7.52	25.00	63.00	30.00
Under	pronouns	100.00	64.74	41.00	40.04	49.00	45.00	3.00
non-factives	monotonicity positive	99.00	52.35	27.26	30.17	37.00	84.00	56.00
	monotonicity negative	100.00	41.96	17.43	22.35	23.00	94.00	69.00
	overall accuracy	97.55	36.98	27.95	36.78	24.27	71.18	58.73

Table 4: Human accuracy and accuracy (%) across LLMs compared. LLaMA 2, GPT-3.5 (text-davinci-003) and gpt-3.5-turbo-0301 were tested on full data. Human performance, gpt-3.5-turbo-0613 and both GPT-4 versions were tested on 100 examples per inference type.

		original promp	ot	paraphrased p	rompt 1	paraphrased prompt 2		
		accuracy (%)	entailment (%)	accuracy (%)	entailment (%)	accuracy (%)	entailment (%)	
	pronouns	53.00	53.00	67.00	67.00	62.00	62.00	
Standalone	monotonicity positives	43.00	43.00	68.00	68.00	65.00	65.00	
Standarone	monotonicity negatives	42.00	37.00	8.00	53.00	37.00	42.00	
	uncertainty adverbs	9.1	79.9	5.67	80.33	8.44	78.67	
	pronouns	72.30	72.30	86.70	86.70	76.61	76.61	
Under presupposition	monotonicity positives	55.65	55.65	78.30	78.30	70.39	70.39	
triggers	monotonicity negatives	29.39	52.83	6.17	59.35	29.70	50.70	
	uncertainty adverbs (sample)	7.52	86.91	1.39	88.65	9.61	82.65	
Under	pronouns	40.04	56.65	20.09	72.87	40.52	56.65	
non-factives	monotonicity positives	30.17	60.17	9.39	72.83	29.57	63.61	
non-ractives	monotonicity negatives	22.35	52.91	6.00	48.87	40.00	38.39	

Table 5: The experiment results for the original prompt and its two paraphrases suggested by ChatGPT itself. The background color indicates the expected label: green - "ENTAILMENT", blue - "NEUTRAL". The pink cells indicate the results that do not fit the pattern exhibited by the original prompt (see Appendix B): predicting more entailment under all embedding types. The bold figures indicate the highest accuracy for a specific inference type across all 3 prompts.

tion 3) and the CoT experiments are compared in Table 6.

# E.2 Manual analysis of ChatGPT's reasoning chains

This expands on the results presented in Section 4, "Analysis of CoT Results". We sampled 55 examples of reasoning chains produced by the ChatGPT at the time of the CoT experiments (5 from each test set) and analyzed them manually.

Some chains of thought (about 14.5%) are clear, correct, lead to a correct answer and demonstrate understanding of the task and the input premise-hypothesis pair. For example:

**Premise:** *It's worth noting that all philosophers are mortal.* 

**Hypothesis**: All hegelians are mortal.

Chain of thought: Step 1: Examine Text 1. It states that "All philosophers are mortal." This means that every philosopher who ever existed or will exist must eventually die. This statement is categorical and absolute.

Step 2: Compare Text 2 to Text 1. Text 2 states that "All hegelians are mortal." Hegelians are a specific subset of philosophers who follow the philosophy of Georg Wilhelm Friedrich Hegel. Since Text 1 states that all philosophers are mortal, and Hegelians fall under the category of philosophers, it follows that Text 2 is true.

Answer: True.

**Prediction**: Entailment.

In some other cases (about 9.1%) the prediction is correct, but the reasoning chain is too vague, the logical transitions between the steps are unclear, so that it's impossible to tell if the model reasoned correctly or not. For example:

**Premise:** Sam violated the school's jewelry policy.

**Hypothesis**: Sam violated something.

Chain of thought: 1. Start by analyzing Text 1. It states that "Sam violated the school's jewelry policy." This implies

		original promp	ot	chain-of-though prompt				
		accuracy (%)	entailment (%)	accuracy (%)	entailment (%)	neutral (%)		
	pronouns	53.00	53.00	7.00	7.00	90.00		
Standalone	monotonicity positives	43.00	43.00	44.00	44.00	53.00		
Standarone	monotonicity negatives	42.00	37.00	53.00	39.00	53.00		
	uncertainty adverbs	9.1	79.9	46.56	46.44	46.56		
	pronouns	72.30	72.30	8.00	8.00	91.00		
Under presupposition	monotonicity positives	55.65	55.65	26.00	26.00	70.00		
triggers	monotonicity negatives	29.39	52.83	56.00	31.00	56.00		
	uncertainty adverbs (sample)	7.52	86.91	52.00	43.00	52.00		
Under	pronouns	40.04	56.65	99.00	1.00	99.00		
non-factives	monotonicity positives	30.17	60.17	78.00	21.00	78.00		
non-ractives	monotonicity negatives	22.35	52.91	58.00	30.00	58.00		

Table 6: The experiment results for the original prompt and the CoT prompt. The background color indicates the expected labels: green – "ENTAILMENT"; blue – "NEUTRAL". The bold figures indicate which prompt scored higher on a specific inference type.

that there is a school policy regarding wearing jewelry and Sam broke this rule.

- 2. Next, move on to Text 2. It states that "Sam violated something." It is not specified what Sam violated, but it is implied that Sam did break a rule or a policy of some sort.
- 3. Based on Text 1 and Text 2, it is safe to say that Text 2 is true. The fact that Text 1 explicitly mentions a specific policy that Sam violated supports the idea that a violation did occur. Text 2 confirms this by stating that Sam violated something, thus aligning with the information provided in Text 1.

Therefore, the answer is: Text 2 is true.

**Prediction**: *Entailment*.

Yet in other cases (about 23.6%) the model is "right for the wrong reasons" ariving at a correct prediction through a wrong reasoning chain. For example:

**Premise:** It was reported that every alien used some quarters.

Hypothesis: Every alien used coins.

Chain of thought: Step 1: Let's analyze Text 1 to extract the key information. The statement says "Every alien used some quarters." This means that all aliens used quarters, but it is possible they also used other types of coins or currency.

Step 2: Now, let's examine Text 2 in light of Text 1. The statement in Text 2 says "Every alien used coins." This

statement seems to suggest that all aliens used coins, but it does not specify which type of coins they used.

Step 3: Since Text 2 does not specify which types of coins alien used, we cannot say for sure whether it is true or false. However, we can say that the statement in Text 2 is neutral since it does not contradict Text 1, but also does not provide us with any new information that supports Text 1.

**Prediction**: Neutral.

The correct label here is indeed "neutral", however the reason cited in the chain of thought ("Since Text 2 does not specify which types of coins alien used, we cannot say for sure whether it is true or false") is wrong. In reality, the relation is neutral because of the non-factive expression "It was reported".

Remarkably, in one case, the model does not arrive at the gold label, but seems to be - at least pragmatically - correct in its reasoning, and the CoT clearly makes sense:

**Premise:** it's worth noting that I don't eat chicken skin.

**Hypothesis**: I don't eat chicken skin at any cost.

Chain of thought: Step 1: Analyze Text 1 The first text states that the person speaking does not eat chicken skin. This seems pretty clear and straightforward.

Step 2: Analyze Text 2 The second text states that the person does not eat chicken skin at any cost. This seems to be an exaggerated statement compared

to the first text. The phrase "at any cost" implies that the person is willing to go to great lengths or make great sacrifices to avoid eating chicken skin. This level of emphasis is not present in Text 1, so it's possible that Text 2 is an exaggerated version of the truth. Step 3: Decide if Text 2 is true, false, or neutral Based on the information we have, we can't definitively say whether Text 2 is true or false. It's possible that the person speaking really does avoid eating chicken skin at all costs, but it's also possible that they are exaggerating 

Model's prediction: Neutral.

fore, the answer is neutral.

Gold label: Entailment.

Finally, in the remaining cases both the chain of thought and the resulting prediction are wrong. For example:

and simply avoid it when possible. There-

**Premise:** The school's deputy headmaster, Howard Jones, told her to remove the crucifix around her neck.

**Hypothesis**: The school's deputy headmaster, Howard Jones, told her to do something.

Chain of thought: 1. Text 1 states that Howard Jones, the deputy headmaster of the school, asked a student to remove her crucifix.

- 2. Text 2 only confirms that Howard Jones did tell the student to do something, but it doesn't specify what that something was.
- 3. Therefore, Text 2 cannot be used to determine whether or not Text 1 is true or false.
- 4. However, it is important to note that if Text 1 is true, and Howard Jones did ask a student to remove their crucifix, this would be a violation of the student's religious freedom and may constitute discrimination.

**Model's prediction**: *Neutral*.

**Correct prediction**: *Entailment*.

Here the model not only demonstrates misunderstanding of the task trying to determine the veracity of Text 1 based on Text 2, while clearly instructed to do the opposite, but also makes a moral judgement (see step 4) which is totally out of place in this task. Below we list some more observations regarding the analyzed reasoning chains:

- 1. The embedding context (presupposition triggers or non-factives) are sometimes mentioned in the chain of thought, but are never used as a basis for the final decision. (One consequence of this is that for the test sets with non-factives all the correct answers result from wrong reasoning chains.)
- 2. For premises with uncertainty adverbs, in 8 out of the 10 analyzed cases the adverb is mentioned and its meaning explained, but only in 3 out of the 10 cases the adverb serves as a basis for the final answer.
- 3. The model constantly misinterprets indefinite pronouns as referring to a specific entity. Hence the incorrect "neutral" labels for most cases of grammatically-specified entailment. For example ChatGPT decides that "Mary lent him money" does not entail "Someone lent him money" because "someone" in Text 2 "could be referring to someone other than Mary". (The correct answer is, of course, "entailment" because "someone" is a generic term encompassing any individual including Mary.)
- 4. The model often gets confused about the monotonicity entailment directions (upward vs. downward), stating, for example, that "No alien ate pork" entails "No alien ate meat" since "pork is a type of meat", but "Every alien used some quarters" does not entail "Every alien used some coins" because Text 2 "does not specify which type of coins they used".
- 5. More generally, the model usually predicts entailment when Text 2 contains a more specific mention than Text 1 (which is, in fact, only correct for cases of downward entailment), and vice versa.
- 6. The reasoning chains are mostly vague, excessively wordy, with unclear logical relations

correct CoT/correct label	23.6%
wrong CoT/correct label	23.6%
wrong CoT/wrong label	50.9%
correct CoT /wrong label	1.82%
CoT coherent and clear	16.4%
underlying LP mentioned in CoT	49.1%
correct understanding of the underlying LP reflected in CoT	23.6%
underlying LP explicitly used in prediction	14.5%
CoT demonstrates correct understanding of the task	81.8%
CoT reflects correct understanding of the input sentences	80.0%

Table 7: Manual CoT analysis results. LP stands for "linguistic phenomena". Some numbers are approximate, since not all the cases are clear-cut, and some reasoning chains are unclear and difficult to analyze.

between steps, which makes them hard to understand and analyse, and often contain obvious logical errors (e.g. "Text 2 is likely true, as it directly contradicts the assumption made in Text 1").

- 7. The CoT can sometimes misrepresent the contents of the input sentences. For example the model claims that the text "*I love something outside the city*" does not mention "love".
- 8. Different chains of thought exhibit contradictory logics. For example, one CoT says "There is no contradiction between the two texts... Therefore, Text 2 can be determined as true", while another reasoning chain states: "Text 2 does not contradict Text 1, so it is neutral."

Quantitatively, the results of this analysis are represented in Table 7.

The analysis shows that zero-shot CoT prompting fails to improve ChatGPT's performance on the task because of various flaws in the generated reasoning chains.