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

Anonymous ACL submission

#### Abstract

We evaluate LLMs' language understanding 001 capacities on simple inference tasks that most humans find trivial. Specifically, we target 004 (i) grammatically-specified entailments, (ii) premises with evidential adverbs of uncertainty, and (iii) monotonicity entailments. We de-007 sign evaluation sets for these tasks and conduct experiments in both zero-shot and chainof-thought setups, and with multiple prompts. The models exhibit moderate to low perfor-011 mance on these evaluation sets in all settings. Subsequent experiments show that embedding 013 the premise under presupposition triggers or non-factives, which should exhibit opposite lin-015 guistic behavior, causes ChatGPT to predict entailment more frequently in the zero-shot and 017 less frequently in the chain-of-thought setup, and in each case regardless of the correct label. Similar experiments with LLaMA 2 exhibit dif-019 ferent yet equally flawed tendencies. Overall these results suggest that, despite LLMs' celebrated language understanding capacity, they have blindspots with respect to certain types of entailments, and that certain informationpackaging structures act as "blinds" overshadowing the semantics of the embedded premise.

#### 1 Introduction

027

033

037

041

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 on 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. Therefore LLMs' errors on such tasks, while making them unreliable for certain applications, do not necessarily mean that they haven't reached *humanlevel* performance on reasoning and understanding. 042

043

044

047

048

053

054

056

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

078

079

081

082

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 but nonetheless pose a challenge to LLMs. Using NLI we reveal some of the models' blind spots that indicate that they are far from a genuine *humanlevel* understanding. Furthermore, we show that some information-packaging structures may act as "blinds" that hinder the semantics of embedded premises, again in contrast to human-like behavior.

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 uncertainty* (*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

		Standalone		Embedded within a Pressuposition Trigger		Embedded within a Non-factive Clause	
Condition	Section	Expected	Model	Expected	Model	Expected	Model
Pronouns	3.1	<b>100</b> % Entail	53% Entail	<b>100</b> % Entail	72.3% Entail	<b>100</b> % Neutral	56.7% Entail
Monotonicity Positive	3.2.3	<b>100</b> % Entail	43% Entail	<b>100</b> % Entail	55.7% Entail	<b>100</b> % Neutral	60.2% Entail
Monotonicity Negative	3.2.4	<b>100</b> % Neutral	37% Entail	100% Neutral	52.8% Entail	<b>100</b> % Neutral	52.9% Entail
Uncertainty Adverbs	3.3	<b>100</b> % Neutral	79.9% Entail	<b>100</b> % Neutral	86.9% Entail		
Pronouns (grammatically-specified entailments) Sue was hungry → someone was hungry Monotonicity (positive) John saw a dog → John saw an animal				Embedded Johr	within a presu n realized s	pposition trigger Sue was hun ne was hungr	igry Ƴ
Monotonicity (negative) John saw an animal AJohn saw a dog Uncertainty adverbs John allegedly saw a dog AJohn saw a dog				Embedded within a non-factive clause John <i>thought</i> Sue was hungry someone was hungry			ıgry ry

Figure 1: High-level summary of the experiments and results.

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.

084

880

100

102

103

104

105

106

107

108

109

We next check how embedding of the premise in a larger grammatical context affects the prediction. Such embedding can take several forms. In particular, 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-factive verbs (e.g. I feel that [...], They thought that [...], He imagined that [...]) may cancel it. We experiment with both types of contexts and show that ChatGPT<sup>1</sup> has a hard time distinguishing the two cases, incorrectly treating both as hints towards entailment (for regular prompting) or against it (for chain-of-thought prompting). These or similar trends are consistent across different prompts and models (such as GPT-3.5 and LLaMA 2), and in both regular and chainof-thought prompting showing that these inference failures are robust, and merit further investigation. Figure 1 summarizes our main results.

## 2 Linguistic background

First, we briefly explain the linguistic phenomena 111 that we use in the experiments described below. 112

110

113

114

115

116

117

118

119

120

126

127

128

129

130

131

132

133

**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.	121
You've eaten all of something.	122
You've eaten something.	123
You've done something.	124
Someone's eaten all my apples.	125

**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 dogs134need water.135

<sup>&</sup>lt;sup>1</sup>https://openai.com/chatgpt

Liu et al. (2023) show that ChatGPT's perfor-136 mance on monotonicity entailment datasets (HELP 137 (Yanaka et al., 2019b) and MED (Yanaka et al., 138 2019a)) improves on previous models but is still far 139 from perfect (42.13% and 55.02% respectively). 140

141

147

151

152

153

154

155

156

157

158

159

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

182

183

Evidential Adverbs "express degrees of certitude with respect to the speaker's subjective per-142 ception of the truth of a proposition" (Haumann, 143 2007). We run experiments that test LLMs' abil-144 ity to understand evidential adverbs expressing un-145 certainty (allegedly, purportedly, supposedly etc.). 146 Introducing such adverbs into a clause cancels the entailment of the remaining part of the clause. For 148 example, Mike allegedly worked all night does not 149 entail Mike indeed worked all night. The relation 150 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 doesn't 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, (We detail them in Appendix A).

Non-factive Verbs and Expressions (Kiparsky 178 and Kiparsky, 1970), such as believe, claim, feel, 179 hope, suspect, think, do not entail either truth or 180 falsity of their complements. For example, given: 181

- **a**. Jane thinks that Bill bought bread.
- **b**. Bill bought bread.

c. Bill didn't buy bread.

186

187

188

189

190

191

192

198

200

201

202

203

204

206

207

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

Sentence a doesn't entail either b or c. The relation between **a** and **b** is neutral, and so is the relation between **a** and **c**.

Presupposition Triggers, Non-Factives and NLI It is important to note that embedding a premise under a presupposition trigger doesn't affect the relations between the premise and hypothesis. By contrast, if we embed the premise under a nonfactive, the relation becomes neutral. For example:

<b>a</b> . A balloon hit a light post.	194
$\models$ <b>b</b> . Something hit a light post	195
Premise <b>a</b> above entails hypothesis <b>b</b> . If we embed	196
premise <b>b</b> under a presupposition trigger as in <b>a</b> ':	197

premise **b** under a presupposition trigger as in **a**': a'. She realized that a balloon hit a light post,

 $\models$  **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" doesn't entail b.

#### 3 **Main Experiments**

**Setup.** In this section our experiments use the gpt-3.5-turbo-0301 ChatGPT version (while Section 4 tests more models). We access ChatGPT through OpenAI's API<sup>2</sup> with the default settings, and use 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). This prompt yields a 71%accuracy<sup>3</sup> on a 300-instance SNLI sample, consistent with previous SNLI literature (Qin et al., 2023; Wang et al., 2023; Jang and Lukasiewicz, 2023).

Probing ChatGPT for Phenomena Explaining. Before each experiment we probe ChatGPT's ability to explain the relevant linguistic phenomenon. In humans, being able to explain a concept typically implies being able to apply it. However in LLMs, explaining does not necessarily imply knowing. While not our study's main focus, this probing exemplifies a contrast between humans' and LLMs' language comprehension attested abilities.

<sup>&</sup>lt;sup>2</sup>https://openai.com/product

<sup>&</sup>lt;sup>3</sup>All the results reported in this paper are from a single run.

Figure 2: ChatGPT correctly explains grammatically specified entailment when prompted for it directly.

#### 3.1 Grammatically-Specified Entailments

# **3.1.1** Probing ChatGPT for the phenomenon explanation

When prompted explicitly with the concept, Chat-GPT correctly explains the meaning and usage of indefinite pronouns (see Figure 2).

#### 3.1.2 Testing ChatGPT in an NLI setting

231

235

239

240

242

243

245

247

248

249

251

**Data:** we manually curated a dataset of 100 pairs with grammatically specified entailments by selecting naturally occurring sentences and, in each one, replacing a noun-phrase with an indefinite pronoun:

**Premise:** Crown Princess Mary of Denmark has given birth to a healthy baby boy.

**Hypothesis**: Someone has given birth to a healthy baby boy.

This is a seemingly very easy dataset, in which all the gold labels are "ENTAILMENT", trivial for any human to solve.

**Results:** The model's accuracy on this (seemingly simple) dataset is 53%, failing in almost half (47%) of the cases.

#### 3.1.3 Embedding the Premises under Presupposition Triggers

**Data:** Next we modify this dataset by embedding the original premise (but not the hypotheses) under a presupposition trigger. For instance:

**Original premise:** *Crown Princess Mary of Denmark has given birth to a healthy baby boy.* 

Modified premise: We are happy that Crown Princess Mary of Denmark has given birth to a healthy baby boy.

#### The hypotheses remain unmodified:

262Someone has given birth to a healthy263baby boy.

Since presupposition triggers don't cancel entailment, the gold label in all the cases also remains unchanged: "ENTAILMENT".

264

265

266

267

269

270

271

272

273

274

275

276

277

278

280

281

303

304

305

306

307

308

310

We used 23 presupposition 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. <u>After he came, I cooked</u> *dinner*. Embedding the 100 original premises under 23 triggers resulted in 2300 examples.

**Results:** In the same setting as before, on this data we achieve a significant increase in accuracy: 72.30% for embedded premises vs. 53% for original, "non-embedded", ones. The model predicted entailment more frequently than it originally did.

# 3.1.4 Embedding the Premises under Non-Factive Verbs

**Data:** We modify the dataset, embedding the premises under non-factive predicates as follows:

**Original premise:** Crown Princess Mary of Denmark has given birth to a healthy baby boy. Modified premise: I hope Crown Princess Mary of Denmark has given birth to a healthy baby boy. 287 The hypotheses remain unmodified: 288 Someone has given birth to a healthy baby boy. 290 Since non-factives cancel entailment, the gold label 291 in all the cases changes to "NEUTRAL". 292 We embedded each of the 100 original premises 293 under 23 different non-factives, such as *feel*, hope, 294 believe etc., obtaining 2300 examples. **Results:** We run the same NLI experiment. Sur-296 prisingly, the model still predicts "ENTAILMENT" 297 in 56.65% of the cases - more often than on the orig-298 inal dataset with non-embedded premises (53%). 3.1.5 Bottomline 300 Our results suggest that ChatGPT struggles with 301 grammatically-specified entailments. It tends to 302

predict "ENTAILMENT" more often when the

premise is embedded under a presupposition trig-

ger, which creates a seeming accuracy increase.

But embedding the premise under non-factives also

causes the LLM to predict "ENTAILMENT" more,

even though in this case it is wrong. So, the mere

syntactic embedding affects NLI consistently, re-

gardless of the semantics of the embedding clause.



Figure 3: ChatGPT correctly explains set-membership relations.

#### 3.2 Monotonicity Entailment

311

312

313

314

315

316

318

319

321

322

323

324

325

326

327

328

332

333

335

336

337

#### Probing ChatGPT for the phenomenon 3.2.1 explanation

ChatGPT correctly explains the set-membership relations (which are typically the basis of monotonicity entailments) when prompted for it directly (see Figure 3). Since the example in the figure only covers a specific set-membership pair and thus might not be representative enough, we run a similar test with more pairs extracted from our monotonicity dataset and get similarly accurate explanations in all the cases (see Appendix C for more details).

#### 3.2.2 Testing ChatGPT in an NLI setting

We sample two subsets of the Monotonicity Entailment Dataset (MED) (Yanaka et al., 2019a):

- 100 positive examples (the gold label is "EN-TAILMENT").
- 100 negative examples (the gold label is "NEUTRAL").
- Example from the positive subset:

**Premise:** She planted blue and purple pansies in the flower bed. Hypothesis: She planted pansies in the flower bed. Label: Entailment Example from the negative subset: Premise: Susan made a dress for Jill.

- Hypothesis: Susan made a long dress for Jill. 339
  - Label: Neutral

We first describe our experiments with the posi-341 tive subset, and then with the negative one. 342

#### 3.2.3 Positives: Standalone Premises

Our usual NLI experiment on the positive part of the monotonicity data, yields an accuracy of 43%. 345

Embedding the premises under presupposition	346
triggers As above, we embed each premise under	347
23 types of triggers, resulting in 2300 examples.	348
We modify the premises accordingly:	349
Original premise: She planted blue and	350
nurnle nansies in the flower bed	351
	001
<b>New premise</b> : After she planted blue	352
and purple pansies in the flower bea,	353
she shuhed planting other jiowers.	354
The hypothesis remains unchanged:	355
She planted pansies in the flower bed.	356
Since presupposition triggers don't change the rela-	357
tions between the premise and the hypothesis, the	358
gold labels remain unmodified: "ENTAILMENT".	359
<b>Desulta:</b> After this modification the accuracy for	000
positive examples improves significantly: 55.65%	300
again the model predicts entailment more often	260
- again, the model predicts channelit more often.	302
Embedding the premises under non-factive	363
verbs We use the same 23 non-factives as for the	364
grammatically-specified entailment experiments	365
above. We apply each non-factive to each original	366
premise, obtaining 2300 examples. For example:	367
<b>Original premise:</b> She planted blue and	368
purple pansies in the flower bed.	369
<b>New premise</b> : I think she planted blue	370
and purple pansies in the flower bed.	371
The hypothesis remains unmodified:	372
She planted pansies in the flower bed.	373
Since non-factives cancel entailment, the gold label	37/
now becomes "NEUTRAL" for all the pairs	375
now becomes Theorem for an me pairs.	010
<b>Results:</b> The model predicts entailment even	376
more often than for premises under presupposition	377
	070
triggers, 60.17% of the time, even though in this	3/0
triggers, 60.17% of the time, even though in this case it is wrong.	379
<ul><li>triggers, 60.17% of the time, even though in this case it is wrong.</li><li>Bottomline: As in the case of grammatically-</li></ul>	379 380
<ul><li>triggers, 60.17% of the time, even though in this case it is wrong.</li><li>Bottomline: As in the case of grammatically-specified entailments, the mere embedding of the</li></ul>	<ul><li>378</li><li>379</li><li>380</li><li>381</li></ul>
<ul><li>triggers, 60.17% of the time, even though in this case it is wrong.</li><li>Bottomline: As in the case of grammatically-specified entailments, the mere embedding of the premise, whether under presupposition triggers or</li></ul>	<ul><li>378</li><li>379</li><li>380</li><li>381</li><li>382</li></ul>
<ul><li>triggers, 60.17% of the time, even though in this case it is wrong.</li><li>Bottomline: As in the case of grammatically-specified entailments, the mere embedding of the premise, whether under presupposition triggers or non-factives, causes the LLM to predict entailment</li></ul>	<ul> <li>376</li> <li>379</li> <li>380</li> <li>381</li> <li>382</li> <li>383</li> </ul>
<ul><li>triggers, 60.17% of the time, even though in this case it is wrong.</li><li>Bottomline: As in the case of grammatically-specified entailments, the mere embedding of the premise, whether under presupposition triggers or non-factives, causes the LLM to predict entailment much more, regardless of the correct inference.</li></ul>	<ul> <li>370</li> <li>379</li> <li>380</li> <li>381</li> <li>382</li> <li>383</li> <li>384</li> </ul>
<ul> <li>triggers, 60.17% of the time, even though in this case it is wrong.</li> <li>Bottomline: As in the case of grammatically-specified entailments, the mere embedding of the premise, whether under presupposition triggers or non-factives, causes the LLM to predict entailment much more, regardless of the correct inference.</li> <li>3.2.4 Negatives: Standalone Premises</li> </ul>	<ul> <li>376</li> <li>379</li> <li>380</li> <li>381</li> <li>382</li> <li>383</li> <li>384</li> <li>385</li> </ul>
<ul> <li>triggers, 60.17% of the time, even though in this case it is wrong.</li> <li>Bottomline: As in the case of grammatically-specified entailments, the mere embedding of the premise, whether under presupposition triggers or non-factives, causes the LLM to predict entailment much more, regardless of the correct inference.</li> <li>3.2.4 Negatives: Standalone Premises</li> <li>In the NLI experiment on the negative part of the</li> </ul>	<ul> <li>376</li> <li>379</li> <li>380</li> <li>381</li> <li>382</li> <li>383</li> <li>384</li> <li>385</li> <li>386</li> </ul>
<ul> <li>triggers, 60.17% of the time, even though in this case it is wrong.</li> <li>Bottomline: As in the case of grammatically-specified entailments, the mere embedding of the premise, whether under presupposition triggers or non-factives, causes the LLM to predict entailment much more, regardless of the correct inference.</li> <li>3.2.4 Negatives: Standalone Premises</li> <li>In the NLI experiment on the negative part of the monotonicity data, the model predicts entailment</li> </ul>	<ul> <li>376</li> <li>379</li> <li>380</li> <li>381</li> <li>382</li> <li>383</li> <li>384</li> <li>385</li> <li>386</li> <li>387</li> </ul>
<ul> <li>triggers, 60.17% of the time, even though in this case it is wrong.</li> <li>Bottomline: As in the case of grammatically-specified entailments, the mere embedding of the premise, whether under presupposition triggers or non-factives, causes the LLM to predict entailment much more, regardless of the correct inference.</li> <li>3.2.4 Negatives: Standalone Premises</li> <li>In the NLI experiment on the negative part of the monotonicity data, the model predicts entailment in 37% of the cases (which is wrong). It predicts</li> </ul>	<ul> <li>376</li> <li>379</li> <li>380</li> <li>381</li> <li>382</li> <li>383</li> <li>384</li> <li>385</li> <li>386</li> <li>387</li> <li>388</li> </ul>
<ul> <li>triggers, 60.17% of the time, even though in this case it is wrong.</li> <li>Bottomline: As in the case of grammatically-specified entailments, the mere embedding of the premise, whether under presupposition triggers or non-factives, causes the LLM to predict entailment much more, regardless of the correct inference.</li> <li>3.2.4 Negatives: Standalone Premises</li> <li>In the NLI experiment on the negative part of the monotonicity data, the model predicts entailment in 37% of the cases (which is wrong). It predicts the correct label "NEUTRAL" for 42% of the pairs</li> </ul>	<ul> <li>376</li> <li>379</li> <li>380</li> <li>381</li> <li>382</li> <li>383</li> <li>384</li> <li>385</li> <li>386</li> <li>387</li> <li>388</li> <li>389</li> </ul>

5

and "CONTRADICTION" for the remaining 21%.

403

404

405

406

407

408

409

412 413

414

415

416

417

418

419

420

**Embedding the premises under presupposition triggers** We embed the premises under presuppo-

sition triggers, as before:

**New premise**: *They are aware that Susan made a dress for Jill.* 

8 The hypothesis remains unchanged:

Susan made a long dress for Jill.

Since presupposition triggers don't change the relations between the premise and hypothesis, the gold
labels remain unmodified: "NEUTRAL".

Applying 23 trigger types to 100 premises results in 2300 new examples.

**Results:** After this modification, the model (incorrectly) predicts "ENTAILMENT" even more: in 52.83% of the cases. The correct label, "NEUTRAL", is predicted for 29.39% of the cases, and the remaining 17.78% are "CONTRADICTION".

410 Embedding the premises under non-factive411 verbs We modify the premises accordingly:

**Original premise:** *Susan made a dress for Jill.* 

**New premise**: *They believe Susan made a dress for Jill.* 

The hypothesis remains unmodified:

Susan made a long dress for Jill

The gold label remains "NEUTRAL" for all pairs. To each premise we apply the same set of 23 non-factives as above, obtaining 2300 pairs.

421**Results:** The model, again, predicts "ENTAIL-422MENT" even slightly more often than for premises423under presuppositon triggers: in 52.91% of the424cases. The correct label "NEUTRAL" appears in42522.35% of the cases, and the remaining 24.74% of426the labels are "CONTRADICTION".

**Bottomline** Thus, both on the positive and negative monotonicity entailment test sets the model demonstrates the same pattern we observed before: embedding the premise inside a clause, under presupposition triggers or non-factives, causes the model to predict entailment more often, regardless of the embedding clause's semantic content.



Figure 4: ChatGPT correctly explains the meaning of evidential adverbs when prompted for it directly.

### 3.3 Adverbs of Uncertainty

# **3.3.1 Probing ChatGPT for the phenomenon** explanation

434

435

436

437

438

439

440

441

442

443

444

ChatGPT correctly explains the meaning and usage of evidential adverbs of uncertainty when prompted for it directly (see Figure 4).

## 3.3.2 Testing ChatGPT in an NLI setting

We manually create a dataset of 100 sentence pairs where the only difference between the premise and the hypothesis is that the premise contains an adverb of uncertainty, while the hypothesis omits it:

<b>Premise:</b> These persons were allegedly	445
inhabiting the home.	446
Hypothesis: These persons were inhab-	447
iting the home.	448
We apply 9 uncertainty adverbs ( <i>allegedly, hope-</i>	449
fully, possibly, presumably, probably, purportedly,	450
reportedly, seemingly, supposedly) to each of the	451
100 pairs, obtaining 900 examples, 100 per adverb.	452
The gold label for all the pairs is "NEUTRAL".	453
<b>Results:</b> The model predicted "ENTAILMENT"	454
in 79.9% of the cases, which is wrong. It pre-	455
dicts "NEUTRAL" in 9.1% of the cases (which is	456
correct) and "CONTRADICTION" in 11% of the	457
cases. The results are mostly consistent across all	458
the individual adverbs (see Appendix D for details).	459
3.3.3 Embedding the premises under	460
presupposition triggers	461
We randomly sample 100 examples from our	462
dataset of 900 sentence pairs, and apply each of the	463
23 triggers to each sampled example:	464
Original premise: These persons were	465
allegedly inhabiting the home.	466
New premise: <u>The owner was aware</u>	467
that these persons were allegedly inhab-	468
iting the home.	469
The hypothesis remains the same.	/70
	470

475

476

477

478

479

481

482

483

484

485

486

487

488

489

490

491

492

493

494

497

498

499

503

504

505

507

510

511

512

513

514

#### These persons were inhabiting the home.

472 Since presupposition triggers don't affect the origi473 nal relation, the gold label also remains the same:
474 "NEUTRAL".

**Results:** The model predicts "ENTAILMENT" even more often: in 86.91% of the cases (which is wrong). It predicts "NEUTRAL" (the correct label) only for 7.52% of the pairs, and "CONTRADIC-TION" - for the remaining 5.57%.

# 3.3.4 Embedding the premises under non-factive verbs

We omit this experiment for adverbs of uncertainty 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, which creates a tautology.

#### 3.3.5 Bottomline

These experiments exhibit the same pattern we observed earlier: when the premise is embedded under a presupposition trigger, the model tends to predict entailment. As a result, in presupposed content ChatGPT appears to overlook evidential adverbs even more than in unembedded premises.

See Figure 1 for a summary of the experimental results described in this section.

#### 4 Model and Prompt Variations

To assess further the LLMs' performance on these phenomena, we use model and prompt variations.

**Prompt variation.** We ask ChatGPT to rephrase our prompt template (see Appendix B), obtaining two templates (see Appendix G) that we verify to be semantically equivalent w.r.t. the task. The first prompt shows the same pattern as our original prompt; the second one yields higher accuracy, and even shows a more reasonable trend of predicting entailment *less* often in non-factive contexts (which indeed cancel entailment). But even this, better, prompt is very far from solving the task.

These differences in the behavior of the three prompts again stress ChatGPT's inconsistency: according to ChatGPT itself they all have the same meaning and describe the same task; however, they produce different outputs. See Appendix G for full experiment details, results and their comparison.

515 Model variation. We repeated the experiments516 from Section 3 with additional models.

1. GPT-3.5. We run the same experiments over the text-davinci-003 model,<sup>4</sup> with a temperature of 0. The results are overall lower than those of ChatGPT, while exhibiting the same trends. At the same time ChatGPT tends to overpredict entailment more than text-davinci-003. See Appendix E for full GPT-3.5 results and their comparison with ChatGPT results. 2. LLaMA 2. To verify that the problem is not model-specific, we run the same set of experiments using the LLaMA-2 Chat model with 70 billion parameters<sup>5</sup> (Touvron et al., 2023).<sup>6</sup> We first assess its performance on the NLI task using the same SNLI sample as for ChatGPT (see Section 3). With our original prompt LLaMA 2 only achieves a 39% accuracy on this test. Therefore we search for another prompt and find that prompt (1) (see Appendix G) from our prompt variation experiments (see above) yields a 61% accuracy. Using this prompt, the temperature of 0.01 (lowest possible) and top-k=1, we run the set of experiments described in Section 3.

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

LLaMA 2 also exhibits moderate to low accuracy (mostly far below its own SNLI results) on all the test sets. For unembedded premises, the LLM tends to predict the opposite of the expected label ("entailment" for uncertainty adverbs where "neutral" is expected; "neutral" for grammatically-specified entailments where "entailment" is expected, etc.)

Like for GPT models, embedding premises in larger grammatical contexts affects the predictions, but the pattern is different. For all inference types (grammatically-specified entailments, monotonicity entailments, uncertainty adverbs), under presupposition triggers and, even more, under nonfactives the accuracy increases relative to the standalone premises. Looking just at the numbers, it seems that such embeddings in larger grammatical contexts make the model more sensitive to simple linguistic inferences, that is, to some extent more accurate.

However, a closer examination of the results suggests that the increase in accuracy might be due to the wrong reasons. For example, for grammaticallyspecified and positive monotonicity entailments, the "neutral" predictions prevail with or without the embedding clauses, regardless of the correct label. Hence the results are seemningly better under non-factives, where the more-frequent neutral

<sup>&</sup>lt;sup>4</sup>https://platform.openai.com/docs/models/
gpt-3-5

<sup>&</sup>lt;sup>5</sup>The largest currently available LLaMA 2 version.

<sup>&</sup>lt;sup>6</sup>We access LLaMA 2 through the Replicate API: https: //replicate.com

label happens to be correct. See Appendix F forfull results.

**Chain-of-Thought Prompting (CoT).** Using the 566 gpt-3.5-turbo-0301 model as in our main experi-567 ments (see Section 3), we investigate CoT prompting (Kojima et al., 2023), and find that the results exhibit a different pattern that is equally incorrect. CoT reverses ChatGPT's trends for embedded con-571 texts observed in Section 3: the number of entailment predictions *decreases* when the premises are embedded under presupposition triggers or non-574 factives, while the number of neutral predictions increases - again, regardless of the correct label. 576 CoT prompting improves the accuracy only in a one-sided way: scorring much higher on all the 578 "neutral" test sets, but much lower on almost all the "entailment" ones.

Analysis of CoT Results. The CoT technique 581 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 584 585 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 understand-591 ing of the task expressed in the prompt. In half of the cases (49.1%) the CoT mentioned the undelying linguistic phenomena explicitly, but only 593 in half of those (23.6% of the total) it reflected a correct understanding of the phenomena and only 595 in 14.5% of the cases used them as a basis for the 596 final prediction.

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

#### 5 Conclusions and discussion

598

599

606

607

610

612

Recent studies have highlighted the need for improvement in LLMs' inferencing and reasoning abilities, since they impact the LLM's consistency, reliability, practical applicability and performance on downstream tasks (e.g., Liu et al., 2023; Jang and Lukasiewicz, 2023; Shen et al., 2023; Zheng et al., 2023; Gao et al., 2023; Jin et al., 2023; Plevris et al., 2023; Luo et al., 2023; Lin and Zhang, 2023).

Flawed heuristics and biases in earlier NLI systems have been extensively studied, forming an important research direction (e.g., Poliak et al., 2018; Nie et al., 2018; Glockner et al., 2018; Sanchez et al., 2018; Gururangan et al., 2018; McCoy et al., 2019; Ross and Pavlick, 2019; Zhou and Bansal, 2020; Asael et al., 2022; Gubelmann and Handschuh, 2022). State-of-the-art LLMs' limitations in NLI indicate the presence of fallible tendencies and blind spots in recent foundation models as well, necessitating further investigation. Our study is an attempt to pinpoint some of these weaknesses. 613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

Our experiments focus on simple inference tasks that are typically trivial for humans. We conduct experiments testing LLMs' ability to handle grammatically specified entailments, evidential adverbs of uncertainty, and monotonicity entailment. The findings suggest that LLMs struggle with these types of inferences, exhibiting a gap between their performance and human-level text understanding.

We observe a trend in ChatGPT to predict entailment more often (regardless of the correct label) when premises are embedded in main clauses, containing presupposition triggers or non-factives. Our experiments with text-davinci-003 suggest that this erroneous trend was inherited by ChatGPT from the earlier model despite the upgrading process. Interestingly, in the CoT setting the trend is reversed: clause-embedding contexts cause the LLM to predict entailment less - again, regardless of the correct label. Using paraphrases of the same prompt suggested by ChatGPT itself yields other (but similarly inaccurate) results, stressing again the LLM's lack of self-consistency and oversensitivity to prompt phrasing. The LLaMA 2 experiments exhibit low accuracy and other, but also flawed, trends, confirming that the issue is not limited to GPT models, but is likely to affect up-to-date LLMs in general.

Our work shows 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 only for three specific models, we do expect them to hold more generally.

We share the evaluation sets and methodology we created, facilitating study with other models. These results also call for further research to uncover other fallible tendencies and blind spots that might hinder LLMs' ability to make accurate textual inferences, which critically affect their reasoning abilities and their potential for real-world applications.

664

667

671

673

679

6 Limitations

This work has some limitations which we list herein.

First, as a consequence of LLMs' sensitivity to prompts, there may exist prompts that can potentially modify the reported results. Moreover, we only tested the LLMs in a regular zero-shot and zero-shot chain-of-thought setting and did not try other approaches like in-context learning or fewshot chain-of-thought prompting. 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, not all possible types of non-factive verbs have been considered. For example, we did not try adding negation to the non-factives or using nonfactive 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, we didn't try other types of clauseembedding predicates (e.g., implicative verbs).

Finally, since ChatGPT undergoes continuous updates, the test results presented here may vary over time. The same is true for LLaMA and future versions of newer models that may become available. That said, our data and methodology for benchmarking these capabilities is model-agnostic and remains intact.

#### References

- Dimion Asael, Zachary Ziegler, and Yonatan Belinkov. 2022. A generative approach for mitigating structural biases in natural language inference. In *Proceedings* of the 11th Joint Conference on Lexical and Computational Semantics, pages 186–199, Seattle, Washington. Association for Computational Linguistics.
- 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.
- Ido Dagan, Oren Glickman, and Bernardo Magnini. 2005. The pascal recognising textual entailment challenge. In *Machine Learning Challenges Workshop*, pages 177–190. Springer.
- Kai Fintel. 2008. What is presupposition accommodation, again? *Philosophical Perspectives*, 22:137– 170.
- Jinglong Gao, Xiao Ding, Bing Qin, and Ting Liu. 2023. Is chatgpt a good causal reasoner? a comprehensive evaluation.
- Max Glockner, Vered Shwartz, and Yoav Goldberg. 2018. Breaking NLI systems with sentences that require simple lexical inferences. In *Proceedings* of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), pages 650–655, Melbourne, Australia. Association for Computational Linguistics.
- Reto Gubelmann and Siegfried Handschuh. 2022. Uncovering more shallow heuristics: Probing the natural language inference capacities of transformer-based pre-trained language models using syllogistic patterns.
- Suchin Gururangan, Swabha Swayamdipta, Omer Levy, Roy Schwartz, Samuel Bowman, and Noah A. Smith. 2018. Annotation artifacts in natural language inference data. In Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers), pages 107–112, New Orleans, Louisiana. Association for Computational Linguistics.
- Dagmar Haumann. 2007. Adverb Licensing and Clause Structure in English.
- 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.
- 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.
- Zhijing Jin, Jiarui Liu, Zhiheng Lyu, Spencer Poff, Mrinmaya Sachan, Rada Mihalcea, Mona Diab, and Bernhard Schölkopf. 2023. Can large language models infer causation from correlation?

755

756

757

758

759

760

761

762

763

764

765

766

69 69 7(

702

703

710

767

- 810 811 812
- 813 814 815
- 816 817 818
- 819 820

- Lauri Karttunen. 2016. Presupposition: What went wrong? Semantics and Linguistic Theory, 26:705-731.
- 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.
- Guo Lin and Yongfeng Zhang. 2023. Sparks of artificial general recommender (agr): Early experiments with chatgpt.
- Hanmeng Liu, Ruoxi Ning, Zhiyang Teng, Jian Liu, Qiji Zhou, and Yue Zhang. 2023. Evaluating the logical reasoning ability of chatgpt and gpt-4.
- Zheheng Luo, Qianqian Xie, and Sophia Ananiadou. 2023. Chatgpt as a factual inconsistency evaluator for text summarization.
- 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.
- Yixin Nie, Yicheng Wang, and Mohit Bansal. 2018. Analyzing compositionality-sensitivity of nli models.
- 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.
- Adam Poliak, Jason Naradowsky, Aparajita Haldar, Rachel Rudinger, and Benjamin Van Durme. 2018. Hypothesis only baselines in natural language inference. In Proceedings of the Seventh Joint Conference on Lexical and Computational Semantics, pages 180-191, New Orleans, Louisiana. Association for Computational Linguistics.
- Chengwei Qin, Aston Zhang, Zhuosheng Zhang, Jiaao Chen, Michihiro Yasunaga, and Diyi Yang. 2023. Is chatgpt a general-purpose natural language processing task solver?

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.

821

822

823

824

825

827

828

829

830

831

832

833

834

835

836

837

838

839

840

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

866

868

869

870

871

872

873

874

875

876

877

878

- Ivan Sanchez, Jeff Mitchell, and Sebastian Riedel. 2018. Behavior analysis of NLI models: Uncovering the influence of three factors on robustness. In Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers), pages 1975–1985, New Orleans, Louisiana. 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.
- 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.
- 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, 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

880	Johan Bos. 2019b. HELP: A dataset for identify-
881	ing shortcomings of neural models in monotonic-
882	ity reasoning. In Proceedings of the Eighth Joint
883	Conference on Lexical and Computational Semantics
884	(*SEM 2019), pages 250-255, Minneapolis, Min-
885	nesota. Association for Computational Linguistics.

- Shen Zheng, Jie Huang, and Kevin Chen-Chuan Chang. 2023. Why does chatgpt fall short in providing truthful answers?
- 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.
- Xiang Zhou and Mohit Bansal. 2020. Towards robustifying nli models against lexical dataset biases.
  - Guido Zuccon and Bevan Koopman. 2023. Dr chatgpt, tell me what i want to hear: How prompt knowledge impacts health answer correctness.

#### A Presupposition triggers

897

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

920

921

922

924

926

929

931

933

934

937

939

940

Presupposition triggers are constructions or lexical items that give rise to *presuppositions* (see Section 2). Presupposition trigger types used in this work are factive predicates, such as *know*, *realize*, *regret*, which presuppose the truth of their complement clause (*I regret leaving the party* presupposes that *I left the party*), temporal and other adverbial clauses (*Paul worked/didn't work as a driver before he started his company* presupposes that *Paul started his company*), and embedded *wh*questions (*This explains/doesn't explain why Jane likes him* presupposes that *Jane likes him*).

There are many other types of presupposition triggers, for example, definite descriptions (*The king of France is bold* presupposes that a king of France exists), aspectual/change of state predicates (*Kate quit smoking* presupposes that *Kate smoked before*), implicative predicates like manage, dare (Joe managed to book a table presupposes that Joe tried to book a table), cleft sentences (It was/wasn't my brother who brought wine presupposes that Someone brought wine), counterfactual conditionals (If Mike were a politician, I would vote for him presupposes that Mike isn't a politician) and others.

#### **B** Main entailment experiments prompt

The prompt below is used throughout the experiments described in Section 3, as well as in the experiments with the text-davinci-003 model (see Section 4 and Appendix E).

> 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"

The model outputs one of three possible labels: "*true*" (corresponding to "*entailment*"), "*false*" (corresponding to "*contradiction*") or "*neutral*".<sup>7</sup>

#### C Probing ChatGPT for explanations of set-membership relations

To probe ChatGPT for explanations of setmembership relations (which are typically the basis of monotonicity entailments) more systematically, we ran a test on 30 set-membership pairs manually extracted from our monotonicity data, asking the model "Is X a kind of Y or is Y a kind of Y?" In each case the model returned a reasonable answer, correctly explaining the concept, similar to the one shown in Figure 3. 941

942

943

944

945

946

947

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

970

971

972

973

Since monotonicity entailments can also be based on pairs of less specific and more specific verbal phrases (e.g. *drinking coffee - drinking coffee at night*) we also extracted 30 VP pairs of this type from our monotonicity entailment data snd asked the model about each pair: "Which description is more specific: X or Y?" Similarly to setmembership pairs, in each case the model provided an accurate explanation, e.g. ""Drinking coffee at night" is more specific because it specifies the time of day when the activity is taking place."

These results show that the model is able to correctly explain relations that form the basis of monotonicity entailments.

#### **D** Adverbs: detailed experiment results

In Table 1 we present the results of the experiments with unembedded premises with uncertainty adverbs using the original prompt (see Appendix B) and the gpt-3.5-turbo model (see Subsection 3.3.2 for details). The table shows the accuracy and the percentage of entailment predictions by individual adverb, as well as the overall results.

adverb	accuracy (%)	entailment(%)
allegedly	12	74
hopefully	10	81
possibly	10	62
presumably	10	82
probably	5	82
purportedly	7	85
reportedly	15	81
seemingly	4	92
supposedly	9	80
overall	9.1	79.9

Table 1: The results of the experiment with unembedded premises with uncertainty adverbs, using the original prompt and the gpt-3.5-turbo model.

#### **E GPT-3.5** experiment results

The results of our experiments with ChatGPT (see Section 3) and GPT-3.5 (see Section 4) are compared in Table 2.

<sup>&</sup>lt;sup>7</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".

		ChatGPT		GPT-3 5	
		accuracy(%)	entailment (%)	accuracy(%)	entailment (%)
Standalone	pronouns	53.00	53.00	39.00	39.00
	monotonicity positives	43.00	43.00	25.00	25.00
	monotonicity negatives	42.00	37.00	28.00	22.00
	uncertainty adverbs	9.1	79.9	4.67	77.56
Under presupposition triggers	pronouns	72.30	72.30	65.35	65.35
	monotonicity positives	55.65	55.65	38.04	38.04
	monotonicity negatives	29.39	52.83	19.04	39.78
	uncertainty adverbs (sample)	7.52	86.91	2.65	87.78
Under non-factives	pronouns	40.04	56.65	41.00	51.00
	monotonicity positives	30.17	60.17	27.26	28.83
	monotonicity negatives	22.35	52.91	17.43	30.30

Table 2: Comparison between ChatGPT's and GPT3.5's experiment results

Looking only at the accuracy columns, we see that ChatGPT constantly beats its predecessor. However, the entailment columns show that it also more readily predicts entailment in almost all the experiments regardless of whether or not it is correct. It is also clear that both models share the same erroneous tendencies. Therefore the higher accuracy scores as such do not necessarily mean that the newer model is more consistent and reliable.

#### F LLaMA 2 experiment results

The results of our experiments with the LLaMA 2 model with 70 bilion parameters are detailed in Table 3.

#### G Consistency across prompts

It is well known that GPT models are very sensitive to the phrasing of prompts: even slight modifications of the prompt can produce considerably different results. Researchers point out that ChatGPT "even produces inconsistent outputs for paraphrased inputs generated by itself" (Jang and Lukasiewicz, 2023). To find out if this is the case also for the inference types explored in this work, we ask ChatGPT twice to generate a paraphrase of our original prompt (see Appendix B). We obtain the following two paraphrases.

(1)

974

975

976

977

978

979

981

982

985

987

989

990

991

993

995

997

1000

1002

1003

1004

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:

1006 Text 1: "text1"

1007 Text 2: "text2"

(2) 10	0
--------	---

Assess the veracity of Text 2 based on	1009
Text 1: Is Text 2 true, false, or indeter-	101
minable? Provide a one-word response.	101

Text 1: "text1" 1012

1014

1015

1016

1017

1018

1019

1020

1021

1023

1024

1025

1026

1028

1029

1030

1031

1032

1033

1034

1035

1036

1037

1038

1039

1040

For both prompts we run the same set of experiments as described in Section 3, with the same data and setup (changing only the prompt).

The first paraphrased prompt demonstrates the same pattern as the original prompt (see Section 3): embedding the premises under either a presupposition trigger or a non-factive predicate makes the model more likely to predict entailment (with only one exception - see Table 4).

Interestingly, the second paraphrased prompt behaves differently: while presupposition triggers again increase entailment predictions, non-factives have the opposite effect: they cause the model to predict entailment *less* often.

The comparison between the three prompts is shown in Table 4.

While the behavior of the second paraphrased prompt might seem somewhat more logical (it predicts less entailment under non-factives which indeed cancel entailment), it is not always more accurate: the original prompt has the best accuracy for the "neutral" test sets with standalone prompts; the first paraphrased prompt yields the best accuracy for the "entailing" test sets, the second paraphrased prompt - for the "neutral" test sets with embeddings, and in general they all struggle with the inference types discussed in this work.

		accuracy(%)	entailment (%)	neutral (%)
Standalone	pronouns	31.00	31.00	61.00
	monotonicity positives	36.00	36.00	59.00
	monotonicity negatives	35.00	46.00	35.00
	uncertainty adverbs	9.00	89.89	9.00
Under presupposition triggers	pronouns	32.00	32.00	58.13
(the initial labels aren't expected to change	monotonicity positives	36.87	36.87	44.87
	monotonicity negatives	38.52	32.43	38.52
	uncertainty adverbs (sample)	29.35	58.39	29.35
Under non-factives	pronouns	64.74	29.74	64.74
(all the labels should change to "NEUTRAL")	monotonicity positives	52.35	34.35	52.35
	monotonicity negatives	41.96	29.30	41.96

Table 3: LLaMA 2 70b results. The green background color means that the expected label is "ENTAILMENT"; the blue background color means that the expected label is "NEUTRAL".

		original prompt		paraphrased prompt 1		paraphrased prompt 2	
		accuracy (%)	entailment (%)	accuracy (%)	entailment (%)	accuracy (%)	entailment (%)
Standalone	pronouns	53.00	53.00	67.00	67.00	62.00	62.00
	monotonicity positives	43.00	43.00	68.00	68.00	65.00	65.00
	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
Under presupposition	pronouns	72.30	72.30	86.70	86.70	76.61	76.61
triggers	monotonicity positives	55.65	55.65	78.30	78.30	70.39	70.39
	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
	monotonicity negatives	22.35	52.91	6.00	48.87	40.00	38.39

Table 4: The experiment results for the original prompt and its two paraphrases suggested by ChatGPT itself. The green rows indicate the datasets where the expected label is always "ENTAILMENT". The blue rows indicate the datasets where the expected lable is always "NEUTRAL". The pink cells indicate the results that don't fit the pattern exhibited by the original prompt (see Appendix B). The bold figures indicate the highest accuracy for a specific test set across all 3 prompts.

#### H Chain-of-thought prompting

1041

1042

1043

1045

1046

1047

1049

1050

#### H.1 Chain-of-thought experiments

As part of our research, we ran the set of experiments described in Section 3, using zero-shot chainof-thought (CoT) prompting, which, basically, consists in adding the phrase "*Let's think step by step*" to the end of the prompt (Kojima et al., 2023). Here we present the experiment details.

We slightly modified our original prompt (see Appendix B) as follows:

1051You are given a pair of texts. Say about1052this pair: given Text 1, is Text 2 true,1053false or neutral (you can't tell if it's true1054or false)?1055Text 1: "text1"1056Text 2: "text2"1057Let's think step by step.

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 comparison between the original experiments described in section 3 and the CoT experiments is shown in Table 5.

As can be seen in the table, the CoT approach hardly proves helpful for the inferences discussed

1058

1059

1060

1061

1062

1063

1064

1065

1067

1068

1069

1070

1071

1072

1073

1074

1075

1076

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

Table 5: The experiment results for the original prompt and the CoT prompt. The green rows indicate the datasets where the expected label is always "ENTAILMENT". The blue rows indicate the datasets where the expected label is always "NEUTRAL". The bold figures indicate which prompt scored higher on a specific test set.

in this work. The CoT method scores much higher on all the "neutral" test sets, but much lower (with one exception) on all the "entailment" test sets.

1078

1079 1080

1081

1082

1083 1084

1085

1086

1087

1088 1089

1090

1091

1092

1093

1094

1095

1096

1099

1100

1101

1102

1103

1104

1105

1106

1107

1111

Another important observation is that the results of the CoT experiments do not share the pattern observed in the previous experiments (more entailment predictions for embedded premises). Rather surprisingly - it exhibits quite the opposite pattern! Namely, the number of entailment predictions decreases when we embed the premises under presupposition triggers or non-factives. At the same time, the number of neutral predictions rises when the premises are embedded (to show this, Table 5 includes an additional column reflecting the percentage of neutral predictions in each CoT experiment). In other words, the pattern we observed earlier is, in a sense, preserved, but now it affects the "neutral", rather than the "entailment" label (again regardless of the correct inference). As for the entailment predictions, now under embeddings they become less likely, rather than more likely as we observed before. In our experiments the regular prompt and the CoT prompt expose mirror-opposite tendencies, but both of these trends are equally flawed.

#### Manual analysis of ChatGPT's reasoning **H.2** chains

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.

Here we present the results of this analysis. 1108 Some chains of thought (about 14.5%) are clear, 1109 correct, lead to a correct answer and demonstrate 1110 understanding of the task and the input premisehypothesis pair. For example:

<b>Premise:</b> It's worth noting that all	1113	
philosophers are mortal.	1114	
Hypothesis: All hegelians are mortal.	1115	
Chain of thought: Step 1: Examine Text	1116	
1. It states that "All philosophers are	1117	
mortal." This means that every philoso-	1118	
pher who ever existed or will exist must	1119	
eventually die. This statement is categor-	1120	
ical and absolute.	1121	
Step 2: Compare Text 2 to Text 1. Text	1122	
2 states that "All hegelians are mor-	1123	
tal." Hegelians are a specific subset of	1124	
philosophers who follow the philosophy	1125	
of Georg Wilhelm Friedrich Hegel. Since	1126	
Text 1 states that all philosophers are	1127	
mortal, and Hegelians fall under the cat-	1128	
egory of philosophers, it follows that Text	1129	
2 is true.	1130	
Answer: True.	1131	
Prediction: Entailment.	1132	
In some other cases (about 9.1%) the prediction	1133	
s correct, but the reasoning chain is too vague, the		
ogical transitions between the steps are unclear,		
so that it's impossible to tell if the model reasoned	1136	
correctly or not. For example:	1137	

Premise: Sam violated the school's jew-	1138
elry policy.	1139
Hypothesis: Sam violated something.	1140
Chain of thought: 1. Start by analyz-	1141
ing Text 1. It states that "Sam violated	1142

1143the school's jewelry policy." This implies1144that there is a school policy regarding1145wearing jewelry and Sam broke this rule.

11462. Next, move on to Text 2. It states that1147"Sam violated something." It is not spec-1148ified what Sam violated, but it is implied1149that Sam did break a rule or a policy of1150some sort.

3. Based on Text 1 and Text 2, it is safe to 1151 say that Text 2 is true. The fact that Text 1 1152 explicitly mentions a specific policy that 1153 Sam violated supports the idea that a vi-1154 olation did occur. Text 2 confirms this by 1155 1156 stating that Sam violated something, thus 1157 aligning with the information provided in Text 1. 1158

1159 Therefore, the answer is: Text 2 is true.

**Prediction**: *Entailment*.

1161

1162

1163

1164

1168

1169

1170

1171

1172

1173

1174

1175

1176

1177

1178

1179

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 alienused 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 1180 which types of coins alien used, we can-1181 not say for sure whether it is true or 1182 false. However, we can say that the state-1183 ment in Text 2 is neutral since it does 1184 not contradict Text 1, but also does not 1185 provide us with any new information that 1186 supports Text 1. 1187

1188**Prediction**: Neutral.

The correct label here is indeed "neutral", how-1189ever the reason cited in the chain of thought ("Since1190Text 2 does not specify which types of coins alien1191used, we cannot say for sure whether it is true or1192false") is wrong. In reality, the relation is neutral1193because of the non-factive expression "It was re-1194ported".1195

1196

1197

1198

1232

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:

1199 **Premise:** *it's worth noting that I don't* 1200 eat chicken skin. Hypothesis: I don't eat chicken skin at 1202 any cost. 1203 Chain of thought: Step 1: Analyze Text 1204 1 The first text states that the person speaking does not eat chicken skin. This 1206 seems pretty clear and straightforward. 1207 Step 2: Analyze Text 2 The second 1208 text states that the person does not eat 1209 chicken skin at any cost. This seems to 1210 be an exaggerated statement compared 1211 to the first text. The phrase "at any cost" 1212 implies that the person is willing to go 1213 to great lengths or make great sacrifices 1214 to avoid eating chicken skin. This level 1215 of emphasis is not present in Text 1, so 1216 it's possible that Text 2 is an exaggerated 1217 version of the truth. 1218 Step 3: Decide if Text 2 is true, false, 1219 or neutral Based on the information we 1220 have, we can't definitively say whether 1221 Text 2 is true or false. It's possible that 1222 the person speaking really does avoid 1223 eating chicken skin at all costs, but it's 1224 also possible that they are exaggerating 1225 and simply avoid it when possible. There-1226 fore, the answer is neutral. 1227 Model's prediction: Neutral. 1228 Gold label: Entailment. 1229 Finally, in the remaining cases both the chain of 1230 thought and the resulting prediction are wrong. For 1231

Premise: The school's deputy headmas-<br/>ter, Howard Jones, told her to remove the<br/>crucifix around her neck.1233

example:

1236Hypothesis: The school's deputy head-1237master, Howard Jones, told her to do1238something.

1239

1240

1241

1242

1243

1244

1245

1246

1247

1248

1249

1250

1251

1252

1253

1254

1255

1256

1257

1258

1259

1260

1261

1263

1264

1265

1266

1267

1268

1269

1270

1272

1273

1274

1275

1276

1277

**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.)
- 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.
- 12783. The model constantly misinterprets indefinite1279pronouns as referring to a specific entity (even1280though it "knows" that indefinite pronouns are1281"generic or underspecified" terms encompass-1282ing any entity or individual see Figure 2).

Hence the incorrect "neutral" labels for most 1283 cases of grammatically-specified entailment. 1284 For example ChatGPT decides that "Mary lent 1285 him money" does not entail "Someone lent him 1286 money" because "someone" in Text 2 "could 1287 be referring to someone other than Mary". 1288 (The correct answer is, of course, "entailment" 1289 because "someone" is a generic term encom-1290 passing any individual including Mary.) 1291

1292

1293

1294

1295

1296

1297

1298

1299

1300

1301

1302

1303

1304

1305

1306

1307

1308

1309

1310

1311

1312

1313

1314

1315

1316

1317

1318

1319

1320

1321

1322

1323

1324

- 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 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"*).
- The CoT can sometimes misrepresent the contents of the input sentences. For example the model claims that the text "I <u>love</u> something outside the city" doesn't mention "love".
- 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 <u>true</u>", while another reasoning chain states: "Text 2 does not contradict Text 1, so it is <u>neutral</u>."

Quantitatively, the results of this analysis are represented in Table 6.

The analysis shows that zero-shot CoT prompt-1325ing fails to improve ChatGPT's performance on1326the task because of various flaws in the generated1327reasoning chains.1328

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		
correct understanding of the underlying LP reflected in CoT		
underlying LP explicitly used in prediction		
CoT demonstrates correct understanding of the task		
CoT reflects correct understanding of the input sentences	80.0%	

Table 6: 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.