

Does ChatGPT Resemble Humans in Processing Implicatures?

Zhuang Qiu, Xufeng Duan, Zhenguang Cai

Department of Linguistics and Modern Languages
The Chinese University of Hong Kong
{zhuangqiu, zhenguangcai}@cuhk.edu.hk
xufeng.duan@link.cuhk.edu.hk

Abstract

Recent advances in large language models (LLMs) and LLM-driven chatbots, such as ChatGPT, have sparked interest in the extent to which these artificial systems possess human-like linguistic abilities. In this study, we assessed ChatGPT's pragmatic capabilities by conducting three preregistered experiments focused on its ability to compute pragmatic implicatures. The first experiment tested whether ChatGPT inhibits the computation of generalized conversational implicatures (GCIs) when explicitly required to process the text's truth-conditional meaning. The second and third experiments examined whether the communicative context affects ChatGPT's ability to compute scalar implicatures (SIs). Our results showed that ChatGPT did not demonstrate human-like flexibility in switching between pragmatic and semantic processing. Additionally, ChatGPT's judgments did not exhibit the well-established effect of communicative context on SI rates.

1 Introduction

In recent years, large language models (LLMs) have achieved unprecedented success in various linguistic tasks, such as disambiguation (Ortega-Martín, 2023), question answering (Brown et al., 2020) and translation (Jiao et al., 2023). However, there is still ongoing debate among researchers about whether these LLMs truly approximate human cognition and language use. On the pessimistic side, Chomsky et al. (2023) argued that “[LLMs] differ profoundly from how humans’ reason and use language. These differences place significant limitations on what these programs can

do, encoding them with ineradicable defects”. In contrast, others have taken a more optimistic view. Piantadosi (2023) argued that recent LLMs should be considered as cognitive models of how people represent and use language.

To address this ongoing debate, researchers have taken an empirical approach by subjecting LLMs to various psychological experiments. Binz and Schulz (2023) subjected GPT-3 to psychological experiments originally designed to study aspects of human cognition such as decision-making, information search and causal reasoning. They found that GPT-3 exhibited human-like or even better-than-human performance in tasks like gamble decisions and multiarmed bandit tasks, with signs of model-based reinforcement learning. Kosinski (2023) tested several language models using the false-belief tasks commonly used to test theory of mind (ToM) in humans. They found that recent GPT models, including GPT-4, GPT-3.5, and GPT-3, provided ToM-like responses similar to those of school children. However, more recent research suggests that ChatGPT's deployment of ToM was not as reliable as that of humans (Brunet-Gouet, Vidal, and Roux, 2023).

Cai et al. (2023) investigated whether ChatGPT resembles humans in language comprehension and production by conducting 12 experiments on psycholinguistic effects at different linguistic levels. They found that ChatGPT exhibited human-like patterns of language use in 10 out of the 12 experiments. For instance, in speech perception, it demonstrated sound-shape (Westbury, 2005) and sound-gender association (Cassidy, Kelly & Sharoni, 1999); in lexical processing, it updated meanings of ambiguous word according to recent input (Rodd et al., 2013); in syntactic processing, it reused recently-encountered syntactic structures (Bock, 1986); in semantic processing, it inferred

78 the likelihood that a sentence is implausible as a
79 result of noise corruption (Gibson et al., 2013) and
80 glossed over errors; at the discourse level, it drew
81 inferences and attributed causality of events
82 according to verb meanings; it was also sensitive to
83 the interlocutor in meaning access and word
84 choice. These results demonstrate that ChatGPT is
85 profoundly similar to humans in its language use.
86 However, it's worth noting that ChatGPT also
87 failed to replicate human patterns in two of the
88 experiments. In one, while humans tend to use
89 shorter words to express less information (e.g.,
90 Mahowald et al., 2013), ChatGPT did not display
91 this tendency. In another, ChatGPT did not make
92 use of context to disambiguate syntactic
93 ambiguities (Altmann and Steedman, 1988).

94 As we delve deeper into LLM-human
95 similarities, it is vital to scrutinize the degree to
96 which ChatGPT's language use aligns with that of
97 humans and to reflect on the implications of such
98 similarities for the evolution of artificial
99 intelligence. Thus, it is important that LLMs are
100 comprehensively tested in order to evaluate how
101 human-like their language use is. So far, one aspect
102 of language use that has not been examined is
103 pragmatics. A hallmark of human language is the
104 ability to convey meanings beyond the literal
105 meaning of the words, through the use of pragmatic
106 implicatures (Grice, 1975; 1978). Experimental
107 pragmatics research has shown that humans can
108 distinguish implicatures from the literal meaning of
109 utterances, and that the computation of
110 implicatures is influenced by the communicative
111 context (Doran et al., 2012; Zondervan, 2010;
112 Bonnefon, Feeney and Villejoubert, 2009). In this
113 project, we assessed the pragmatic capabilities of
114 LLMs by subjecting ChatGPT to three pre-
115 registered experiments that focused on the
116 computation of pragmatic implicatures. The first
117 experiment aimed to determine whether ChatGPT
118 is able to inhibit the computation of generalized
119 conversational implicatures (GCIs) when explicitly
120 required to process the literal meaning of the text.
121 The second and third experiments tested whether
122 the communicative contexts affect how ChatGPT
123 computes scalar implicatures (SIs).

124 2 Experiment 1

125 In this experiment, we tested whether ChatGPT can
126 distinguish “what is said” from “what is
127 implicated” as human beings do. According to
128 standard linguistic accounts, “what is said” refers

129 to the truth-conditional meaning of an utterance,
130 while “what is implicated” refers to the pragmatic
131 implicature, which is an additional level of
132 meaning that is enriched during the conversation
133 (Grice, 1975; 1978). For instance, consider the
134 sentence “Bill caused the car to stop” (Levinson,
135 2000, p. 39). While this sentence is semantically
136 compatible with the scenario in which Bill
137 slammed on the brakes, its implicature suggests
138 that Bill stopped the car in an unconventional way,
139 thus excluding the possibility that he stopped it
140 with the foot pedal.

141 The computation of such implicature is believed
142 to follow general principles of conversation and
143 involve reasoning about the possible alternatives
144 that the speaker could have used (Grice, 1975). For
145 example, interlocutors are expected to be truthful
146 while also making their utterances clear and
147 understandable. If Bill stopped the car in a typical
148 way, the speaker would have said something like
149 “Bill slammed on the brakes.” The fact that the
150 speaker didn't use this typical expression implies
151 that Bill didn't use the brakes to stop the car and
152 might have stopped it in an unconventional way.
153 This pragmatic implicature is enriched based on the
154 literal meaning of the utterance. We are so used to
155 interpreting utterances pragmatically that we often
156 bypass their literal meaning, unless the implicature
157 is explicitly canceled, as in “Bill caused the car to
158 stop, I mean he slammed on the brakes.”

159 A critical question in the study of pragmatic
160 implicatures is whether non-experts can
161 differentiate between “what is said” and “what is
162 implicated.” To address this issue, Doran, Ward,
163 Larson, McNabb, and Baker (2012) measured the
164 rate at which people compute a variety of
165 generalized conversational implicatures (GCIs) in
166 different experimental manipulations. These GCIs
167 are implicatures that can be inferred without
168 reference to the context (Grice, 1975). The study
169 found that, by default, participants were able to
170 derive the implicature of an utterance around half
171 the time. However, the computation of GCIs
172 decreased if participants were explicitly instructed
173 to focus only on the literal meaning of the
174 utterance. This suggests that non-experts without
175 training in linguistics can still distinguish
176 pragmatic implicature from the literal meaning. We
177 adopted the experimental design of Doran et al.
178 (2012) to investigate whether ChatGPT exhibits
179 similar patterns to human participants when
180 processing GCIs.

181 2.1 Design and stimuli

182 The design of this experiment was based on that of
183 Doran et al. (2012). As shown in (1), ChatGPT was
184 presented a mini dialogue, where Irene asked a
185 question and Sam responded to the question. The
186 mini dialogue was followed by a statement of the
187 fact. ChatGPT was then asked to decide, given the
188 factual statement, whether Sam’s response was true
189 or false.

190 1.Q-based GCI:

191 Irene: How much cake did Gus eat at his
192 sister’s birthday party?

193 Sam: He ate most of the cake.

194 FACT: By himself, Gus ate his sister’s entire
195 birthday cake.

196 In (1), the GCI in question belongs to what is called
197 a “Q-based” implicature (Levinson, 2000), where a
198 weaker quantifier (i.e., “most”) in the scale of
199 informativeness implicates the negation of a
200 stronger quantifier (i.e., “all”, as expressed by the
201 word “entire” in the factual statement). That is,
202 quantifiers “some-most-all (entire)” form a scale of
203 increasing informativeness in that if “all of X”
204 holds, then “most of X” holds, and “some of X”
205 must hold, but not vice versa. Given the scale, the
206 utterance “some of X” implicates the negation of
207 “most of X” and “all of/ entire X”; similarly, the
208 utterance of “most of X” implicates the negation of
209 “all of/ entire X”. Thus, based on the factual
210 statement, Sam’s response is logically true but
211 pragmatically infelicitous. Judging Sam’s response
212 as false indicates successful GCI computation and
213 judging it as true indicates the computation of the
214 literal meaning but not of GCI.

215 Apart from Q-based GCIs, Doran et al. (2012)
216 also investigated two other types of GCIs: “I-
217 based” implicatures and “M-based” implicatures.
218 The former refers to cases where the speaker says
219 as little as necessary while the listener needs to
220 “amplify the informational content of the speaker’s
221 utterance by finding the most specific
222 interpretation” (Levinson, 2000). For example, the
223 utterance “She walked into the bathroom. The
224 window was open.” has the implicature that the
225 window is in the bathroom, while the truth-
226 conditional meaning of the utterance allows for the

227 possibility that the window is located elsewhere.
228 “M-based” implicatures refer to cases where the
229 speaker uses a marked way in the description of a
230 common state of affairs, implicating that the
231 unmarked form of the state of affairs does not hold.
232 For instance, the phrase “waited and waited”
233 implies an extended duration of waiting, despite its
234 literal meaning being agnostic to the length of the
235 waiting period. The three types of GCIs each have
236 their own subcategories, as detailed in Appendix A.
237 Each subcategory consisted of four experimental
238 items, resulting in a total of 44 experimental items.
239 Additionally, 16 filler items were included (taken
240 from Doran et al., 2012), which did not require the
241 computation of GCIs.

242 The experiment had two conditions: pragmatic
243 and literal. In the pragmatic condition, ChatGPT
244 was instructed to evaluate the truth of Sam’s
245 response based on the factual statement. After each
246 dialogue and the factual statement, we prompted
247 ChatGPT with “Please judge whether what Sam
248 says is true or false based on the fact.” In the literal
249 condition, ChatGPT was instructed to interpret
250 Sam’s response literally. We prompted ChatGPT
251 with “Please judge whether what Sam says is
252 literally true or false based on the fact.” Doran et
253 al. (2012) found that, compared to the literal
254 condition, the pragmatic condition led human
255 participants to compute more GCIs (i.e., to evaluate
256 Sam’s responses more often as false). We aimed to
257 investigate whether ChatGPT exhibits similar
258 sensitivity to the instructions in drawing GCIs.

259 2.2 Procedure

260 We followed the data collection procedure
261 preregistered with the Open Science Framework
262 (<https://osf.io/cp29j>), eliciting responses
263 from ChatGPT (Feb 13 version)¹. In each run, we
264 used a Python script to simulate a human
265 interlocutor having a conversation with ChatGPT.
266 We first presented a training example (in the
267 pragmatic or literal condition), followed by actual
268 experimental stimulus (see Appendix A). ChatGPT
269 was instructed to respond by saying only “true” or
270 “false” without other words or explanations, and
271 we recorded the responses. In total, this study had
272 400 runs, with 200 runs for each condition.

¹ The original study of Doran et al. (2012) included a third condition known as the “literal Lucy” condition, which was also included in our preregistration. We specified that we would only collect data for this

condition if ChatGPT could pass a sanity check test. Our testing revealed that ChatGPT consistently failed the sanity check. As per our preregistration plan, we did not collect data for this condition.

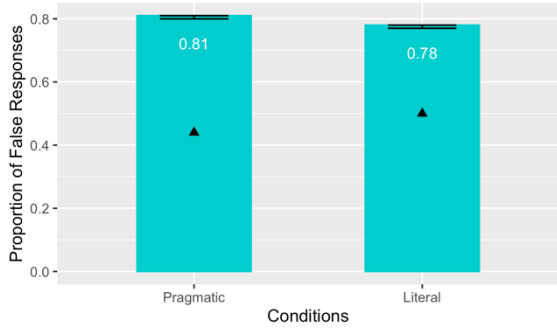


Figure 1: Proportion of false responses (i.e., GCIs) in the pragmatic and literal condition in Exp1. Note, the error bars represent confidence interval (computed using bootstrapping). The triangles represent conditional means from human participants in Doran et al. (2012).

2.3 Results and Discussion

Doran et al. (2012) found that human participants in the pragmatic condition were more likely to evaluate Sam’s response as false (50%) than those in the literal condition (44%), and such a difference was statistically significant. Given that in all the experimental items, Sam’s response was pragmatically infelicitous but logically compatible with the fact, the “false” judgements reflected the computation of GCIs. In this study, we found much higher rates of “false” judgements for the experimental items in both the pragmatic condition (81%) and the literal condition (78%) (see Figure 1). Following the preregistered analytical plan, we applied a Bayesian generalized linear model to trial-level responses (true or false, using true as the reference level), using condition (pragmatic vs. literal) as the predictor. The random effects structure consisted of by-item intercepts and slopes, which was the maximal random effects structure for a between-subjects design. Though there was a slight decrease of false responses in the literal compared to the pragmatic condition, this difference was not statistically significant ($\beta = -0.15$, $CI = [-0.9, 0.63]$). As an exploratory analysis, we investigated the possibility that the effect of the condition was modulated by the category of the GCIs. Another Bayesian generalized linear model was constructed using the condition (pragmatic vs literal, dummy-coded with the pragmatic condition being the reference level), the category of the GCIs (I-based, M-based, and Q-based, dummy-coded with the Q-based GCIs being the reference level), and their interactions to predict the probability of giving a false response (i.e., GCI). The results

showed that none of the effects in the model were statistically meaningful (see Table 1). Instead of showing human-like flexibility switching between pragmatic and semantic interpretation, ChatGPT was unable to inhibit the computation of GCIs even when it was instructed to do so.

3 Experiment 2

In this experiment, we aimed to further investigate ChatGPT’s ability to draw pragmatic inferences, specifically in relation to a type of Q-based GCIs known as scalar implicatures (SIs). SIs are a well-studied phenomenon where the presence of a lower scalar item implies the negation of the higher scalar items (Horn, 1972). For instance, the sentence

	Estimate	Est.Error	l-95% CI	u-95% CI
Intercept	3.89	1.22	1.63	6.40
Literal	-0.66	0.52	-1.69	0.36
M-Based GCIs	0.07	2.04	-3.93	4.08
I-Based GCIs	0.55	1.81	-3.00	4.12
Literal:M-Based GCIs	0.94	0.96	-0.87	2.92
Literal:I-Based GCIs	0.76	0.77	-0.71	2.34

Table 2: The effect of condition, the category of the GCIs and their interactions in Exp1. Note, an estimate is statistically meaningful when zero is not included within the 95% credible interval.

“Sam had a hot dog or a hamburger for lunch” implies that Sam did not have both a hot dog and a hamburger for lunch, even though the sentence’s literal meaning allows for this possibility.

Zondervan (2010) argued that an important contextual factor that influences the interpretation of scalar items is the information structure—whether the scalar item concerns the information focus or information background. For example, the sentence “Julie had found a crab or a starfish”, can be the answer to two different questions as follows:

- 2a. What had Julie found?
- 2b. Who had found a crab or a starfish?

Depending on the question, the same sentence “Julie had found a crab or a starfish” has different information structure. When it is the answer to question 2a, the second half of the sentence including the scalar item “or” is the information focus (new information), while the first half of the sentence including the subject and main verb is the information background (given information). On the other hand, if the same sentence is the answer to question 2b, the subject “Julie” becomes the information focus while the scalar item retreats to the information background. Zondervan conducted

347 a series of experiments, showing that readers are 397
348 more likely to derive the SI of “or” when it is part 398
349 of the information focus compared with the cases 399
350 in which the scalar item is part of the information 400
351 background. We wonder if ChatGPT resembles 401
352 human beings showing similar sensitivity to 402
353 conversational context when processing scalar item 403
354 “or”. If ChatGPT has acquired the pragmatic 404
355 knowledge similar to that of the humans, it should 405
356 be more likely to interpret the expression “A or B” 406
357 as “A or B but not both A and B” when it is part of 407
358 the information focus compared with the case in 408
359 which the expression “A or B” is part of the 409
360 information background. To further explore the 410
361 way ChatGPT processes scalar items, we replicated 411
362 the second experiment in Zondervan (2010) using 412
363 ChatGPT as the participant. 413

364 3.1 Design and stimuli 414

365 The experimental items of the study consisted of 415
366 six short story pairs, each followed by a true-or- 416
367 false question. All the stories ended with a 417
368 conversation between two characters, in which one 418
369 character used the scalar item “or” in his/her reply 419
370 to another character’s question (see 3 and 4). Each 420
371 story in a pair differed in terms of the context where 421
372 the scalar item occurred- whether the scalar item 422
373 being part of the information focus or the 423
374 information background. In the scalar-implicature- 424
375 relevant (SI-relevant) condition (see 3), the 425
376 question was about the object (“what” question), 426
377 and the scalar item “or” was part of the information 427
378 focus. In this case, the interpretation of the scalar 428
379 item as either “A or B but not both A and B” or “A 429
380 or B and possibly both A and B” had particular 430
381 relevance to the conversation. In the scalar- 431
382 implicature-irrelevant (SI-irrelevant) condition 432
383 (see 4), the question is about the subject (“who” 433
384 question), and the scalar item was part of the 434
385 information background. Thus, the interpretation 435
386 of the scalar item was not the major concern of the 436
387 conversation. Crucially, based on the information 437
388 provided in the story, the using of the scalar item 438
389 “or” was logically sound but pragmatically 439
390 infelicitous, and at the end of the story, ChatGPT 440
391 was asked to judge if the character’s answer was 441
392 true or false. If the SI of “or” was computed, 442
393 ChatGPT would respond with “false” to the 443
394 question; or conversely, if the SI was not computed, 444
395 a “yes” judgement would be given.

396 3. SI-relevant:

Julie and Karin were searching for marine animals on the beach. After some searching Julie found a crab. Not much later she also found a starfish. Unfortunately, Karin didn’t find anything. When Karin returned, her mother asked what kind of marine animals Julie had found. Karin answered that Julie had found a crab or a starfish.

Is Karin’s answer true or false?

4. SI-irrelevant:

Julie and Karin were searching for marine animals on the beach. After some searching Julie found a crab. Not much later she also found a starfish. Unfortunately, Karin didn’t find anything. When they returned, their mother asked who had found a crab or a starfish. Karin answered that Julie had found a crab or a starfish.

Is Karin’s answer true or false?

In Zondervan's original study (2010), the experimental items comprised six pairs of stories similar to (3) and (4) but written in Dutch. For the present study, we utilized the English versions of these stories as the experimental items. Additionally, we created 14 filler items that mirrored the length and structure of the experimental items. Each filler item contained a dialogue in which one character answered the question posed by the other character. Half of the filler items were designed to elicit a “true” response, while the other half were designed to elicit a “false” response. To balance the experimental conditions and the order of stimuli, we employed four pseudo-randomized lists of items, following Zondervan's original study.

432 3.2 Procedure 433

434 We followed the data collection procedure 435
436 preregistered with the Open Science Framework 437
438 (<https://osf.io/egm7v>), eliciting responses 439
440 from ChatGPT (Feb 13 version). In each run of the 441
442 experiment, we used a Python script to simulate a 443
444 human interlocutor having a conversation with 444
ChatGPT. At the start, the human interlocutor instructed ChatGPT to make truth-value judgements based on the content of the stories. Two practice trials were given to ChatGPT, the correct answer of which was “true” and “false” respectively. After the practice trial, ChatGPT was

randomly assigned to one list of items, which were presented sequentially. For each item, ChatGPT was instructed to respond by saying only “true” or “false” without other words or explanations, and we recorded the responses from ChatGPT. In total, this study had 200 runs of the script, with 50 runs for each list of items.

3.3 Results and Discussion

In Zondervan (2010), the rate of “false” judgements (i.e., SIs) was 67% in the SI-relevant condition and 41% in the SI-irrelevant condition. In our experiment, ChatGPT responded with “true” for more than 99% of the experimental items, regardless of whether the item was in the SI-relevant or SI-irrelevant condition. The “true” judgement meant that ChatGPT judged the pragmatic infelicitous usage of “or” as “true”, which suggested a lack of pragmatic interpretation. Only one trial in the SI-relevant condition and two

	“False”	“True”
Experimental items		
SI-relevant	1	599
SI-irrelevant	2	598
Filler items		
Correct Answer: False	1394	6
Correct Answer: True	96	1304

Table 2: A summary of judgements from ChatGPT for experimental items and filler items across different conditions in Exp2. Note, the column labels indicate the judgements provided by ChatGPT.

trials in the SI-irrelevant condition received a “false” judgement, which was typically interpreted as the computation of SIs (see Table 2). Given the large number of trials in the experiment, the difference between SI-relevant and SI-irrelevant condition regarding the rate of SI computation was not statistically meaningful (beta = -1.31, CI = [-10.81, 4.78]).

Our analysis of the filler items revealed that ChatGPT demonstrated sensitivity to the truth conditions of the statements (see Table 2). When the character in the story provided an untruthful response, and thus the correct answer to the question should have been “false”, ChatGPT provided more “false” judgments than “true” judgments (1394 vs. 6). Conversely, when the correct answer to the filler item was “true”,

ChatGPT provided more “true” judgments than “false” judgments (1304 vs. 96). To further explore the impact of the correct answer on ChatGPT’s judgments, we modeled the probability of ChatGPT providing a “false” judgment as a function of whether the correct answer to the filler item was “true” or “false” (both dummy coded with the “false” answer being the reference level). Maximal random effects structures were constructed including subject and item intercepts and slopes. We found that when the correct answer of the filler item was “true”, the “false” judgements from ChatGPT decreased at a statistically meaningful rate (beta = -19.64, CI = [-33.92, -11.66]). In total, the accuracy rate of ChatGPT in answering the filler items was above 85 percent.

In this experiment, we investigated whether ChatGPT exhibited human-like patterns of scalar implicature computation by responding to the information structure of the communicative context. Previous research on human participants has shown that when the scalar item “or” was in the information focus, they were more likely to derive the upper bounded reading (“A or B but not both A and B”) compared to when the scalar item was in the information background. Our findings suggest that ChatGPT consistently provided “true” responses when asked if “A or B” is true when both A and B occur, indicating that it interpreted the scalar item “or” as lower bounded (“A or B and possibly both A and B”) for over 99% of the trials, regardless of whether it appeared in the information focus or background. Furthermore, ChatGPT did not always provide “true” responses. For filler items where the correct answer was “false”, ChatGPT provided significantly more “false” responses than “true” responses, and its accuracy rate was high. Therefore, the reason why ChatGPT almost always provided a “true” response for experimental items was that it always endorsed the pure logical interpretation rather than the pragmatic interpretation of the scalar item “or”. The lack of scalar implicature computation for this scalar item and the insensitivity to the information structure of the communicative context differentiate ChatGPT from human participants.

4 Experiment 3

For human participants, the computation of SI is modulated by the conversational context, and the result of Experiment 2 suggested that ChatGPT lacked the sensitivity to the manipulation of

information structure, an important aspect of the conversational context. This experiment aimed to investigate whether conversational context affects how ChatGPT processes scalar implicature (SI) using a different contextual aspect and a different scalar item. Bonnefon, Feeney, and Villejoubert (2009) found that the rate of endorsing SIs for the scalar item “some” decreased when the lower bounded interpretation (“some and possibly all”) threatened the face of the listener, compared to when it boosted the listener’s face. In this experiment, we aimed to test whether ChatGPT shows similar sensitivity to conversational context. We adopted the same design as the first study in Bonnefon, Feeney, and Villejoubert (2009), comparing the rate of SI computation across two within-participants conditions. Unlike the original study, we did not recruit human participants but tested whether ChatGPT exhibits similar performance as human participants. Specifically, we examined whether ChatGPT is more likely to interpret the scalar item “some” as “some but not all” in the face-boosting context, but not so much when the scalar item “some” appears in the face-threatening context.

4.1 Design and stimuli

In this experiment, ChatGPT read two scenarios which were either face-threatening or face-boosting, and the scalar item “some” appeared in the description of the scenario. After reading each scenario, ChatGPT was required to answer a yes-no question. Specifically, we asked ChatGPT whether it would endorse the lower-bounded interpretation of some (which is “some and possibly all”). An example of the experimental item in the face-threatening and face-boosting context was shown in (5) and (6):

5. Face-threatening context:

Imagine that you have joined a poetry club, which consists of five members in addition to you. Each week, one member writes a poem, and the five other members discuss the poem in the absence of its author. This week, it is your turn to write a poem and to let others discuss it. After the discussion, one fellow member confides to you that “Some people hated your poem.”

Yes/No question: From what this fellow member told you, do you think it is possible that everyone hated your poem?

6. Face-boosting context:

Imagine that you have joined a poetry club, which consists of five members in addition to you. Each week, one member writes a poem, and the five other members discuss the poem in the absence of its author. This week, it is your turn to write a poem and to let others discuss it. After the discussion, one fellow member confides to you that “Some people loved your poem.”

Yes/No question: From what this fellow member told you, do you think it is possible that everyone loved your poem?

We included two scenarios like 5 and 6, creating two lists of items using the Latin Squared Design. All items in the experiment were directly adopted from Bonnefon, Feeney and Villejoubert (2009).

4.2 Procedure

We followed the data collection procedure preregistered with the Open Science Framework (<https://osf.io/3v9gn>), eliciting responses from ChatGPT (Feb 13 version). In each run of the experiment, we used a Python script to simulate a human interlocutor having a conversation with ChatGPT. At the start, the human interlocutor instructed ChatGPT to answer yes-no questions based on the description of scenarios. Two practice trials were given to ChatGPT, the correct answer of which was “yes” and “no” respectively. After that, ChatGPT was randomly assigned to one list of items, which were presented to ChatGPT in a random order. For each item, ChatGPT was instructed to respond by saying only “yes” or “no” without other words or explanations, and we recorded the responses from ChatGPT. In total, this study had 200 runs of the script, with 100 runs for each list of items.

	“No”	“Yes”
Face-boosting	198	0
Face-threatening	198	0

Table 3: A summary of judgements from ChatGPT for experimental items across different conditions in Exp3.

619 4.3 Results and Discussion

620 According to our preregistered data exclusion
621 criteria, we excluded data from two runs of the
622 experiment because ChatGPT answered the second
623 practice trial incorrectly, indicating that it may not
624 provide reliable judgments in that run of the
625 experiment. Therefore, we analyzed the data from
626 198 runs of the experiment. In Bonnefon, Feeney
627 and Villejoubert's (2009) study, 83% of human
628 participants responded with "no" when asked if the
629 lower bounded interpretation of "some" was
630 possible in the face-boosting context, while a
631 significantly lower 58% responded "no" in the
632 face-threatening context. In contrast, our study
633 found that ChatGPT always responded "no" to all
634 of the trials, regardless of whether the context was
635 face-boosting or face-threatening (see Table 3).

636 Though the exact mechanism is still unclear
637 regarding why human participants were more
638 likely to interpret the construction "some verb-ed
639 X" as "some and possibly all verb-ed X" in the face
640 threatening context than in the face boosting
641 context, Bonnefon, Feeney and Villejoubert (2009)
642 suggested that the listener may take into account
643 the intension of the speaker to use the word "some"
644 in an underinformative way in order to protect the
645 face of the listener. Although, the SI rate of "some"
646 decreased in the face threatening condition, in
647 general, human participants preferred the
648 pragmatic interpretation of "some" as "some but
649 not all", and that is why even in the face-
650 threatening condition, the majority of the human
651 participants (58%) provided a "no" judgement to
652 the question "Do you think it is possible that
653 everyone hated..." In our experiment with
654 ChatGPT, we clearly saw a stronger preference for
655 the pragmatic interpretation of "some" over the
656 truth-conditional interpretation. In fact, ChatGPT
657 exhibited zero variance in its judgements- for all
658 the trials that contained the scalar item "some",
659 ChatGPT always interpreted them as "some but not
660 all", and thus said "no" to the question, regardless
661 of whether the implicature was face threatening or
662 face boosting to the listener.

663 5 General Discussion and Conclusion

664 In three experiments, we investigated whether
665 LLMs like ChatGPT exhibit human-like
666 performance when processing pragmatic
667 implicatures. Previous research has shown that
668 humans distinguish implicatures from the truth-

669 conditional meaning of the utterance, and several
670 factors have been identified that modulate the
671 probability of implicature computation. While
672 pragmatic enrichment is an essential component of
673 successful communication, whether an implicature
674 is computed by a specific listener in a specific
675 communicative context is probabilistic in nature. In
676 contrast, our findings revealed that ChatGPT
677 lacked human-like flexibility in switching between
678 pragmatic and semantic interpretation, as it was
679 unable to inhibit the computation of GCIs even
680 when instructed to do so. Notably, the processing
681 of scalar items in ChatGPT exhibited a
682 deterministic pattern: whereas "some" always
683 received an upper bounded interpretation as "some
684 but not all", the expression "A or B" almost always
685 received a lower bounded interpretation as "A or B
686 and possibly both A and B".

687 Given ChatGPT's impressive human-like
688 performance across a range of language tasks (Cai
689 et al., 2023), one might question why humans and
690 LLMs differ in their computation of GCIs. Our
691 argument is that this difference can be explained by
692 the acquisition of GCIs and the computational
693 resources available to humans and machines.
694 Developmental research indicates that scalar items
695 are acquired with a lower bounded interpretation
696 before pragmatic enrichments (Noveck, 2001).
697 Consequently, adults have access to both the literal
698 and pragmatic interpretations of a scalar item,
699 whereas LLMs are exposed to language data that
700 are mainly pragmatically driven. This explains why
701 ChatGPT, in general, is more prone to pragmatic
702 interpretation compared with human participants.
703 However, it is still unclear why some specific word
704 like "or" almost always evokes a literal rather than
705 pragmatic interpretation. Furthermore, humans
706 possess limited computational resources compared
707 to machines. The principle of economy suggests
708 that the human mind enriches the truth-conditional
709 meaning only when the context necessitates it
710 (Noveck & Sperber, 2007). This echoes the fact
711 that the effect of contextual manipulation has only
712 been observed among human participants rather
713 than LLMs. It is consistent with the observation
714 that humans tend to use shorter forms of words
715 (e.g., math instead of mathematics) when the
716 meaning is predictable, while ChatGPT does not
717 (Cai et al., 2023). Overall, our experiments
718 demonstrate that although LLM-based chatbots
719 such as ChatGPT excel in many language tasks,

720 they do not mimic humans in their computation of 770
721 GCIs. 771 *phonology. Journal of Experimental Psychology: General* 128, no. 3 (1999): 362.

722 **Limitations**

723 The scope of our research is limited to uncovering 772
724 the distinction between humans and LLMs in a 773
725 specific aspect of pragmatic processing: the 774
726 computation of GCIs. While we offer tentative 775
727 explanations for the patterns we observed, our 776
728 study does not directly provide solutions for 777
729 improving the performance of LLMs. In this study, 778
730 we use ChatGPT as an example of LLMs due to its 779
731 prominence in current research. However, it 780
732 remains uncertain whether other LLMs exhibit 781
733 comparable characteristics and tendencies as 782
734 observed in ChatGPT. Moreover, it is important to 783
735 note that our findings may not generalize to the 784
736 processing of other types of pragmatic 785
737 implicatures. 786
787

738 **References**

- 739 Gerry Altmann, and Mark Steedman.1988. *Interaction*
740 *with context during human sentence processing.*
741 *Cognition* 30, no. 3: 191-238.
- 742 Marcel Binz and Eric Schulz. 2023. *Using cognitive*
743 *psychology to understand GPT-3. Proceedings of*
744 *the National Academy of Sciences* 120, no. 6:
745 e2218523120.
- 746 J. Kathryn Bock. 1986. *Syntactic persistence in*
747 *language production. Cognitive psychology* 18, no.
748 3: 355-387.
- 749 Jean-François Bonnefon, Aidan Feeney, and Gaëlle
750 Villejoubert. 2009. *When some is actually all:*
751 *Scalar inferences in face-threatening contexts.*
752 *Cognition* 112, no. 2: 249-258.
- 753 Tom Brown, Benjamin Mann, Nick Ryder, Melanie
754 Subbiah, Jared D. Kaplan, Prafulla Dhariwal,
755 Arvind Neelakantan et al. 2020. *Language models*
756 *are few-shot learners. Advances in neural*
757 *information processing systems* 33: 1877-1901.
- 758 Eric Brunet-Gouet, Nathan Vidal, and Paul Roux.
759 2023. *Do conversational agents have a theory of*
760 *mind? A single case study of ChatGPT with the*
761 *Hinting, False Beliefs and False Photographs, and*
762 *Strange Stories paradigms. HAL Open Science.*
763 <https://hal.science/hal-03991530/>
- 764 Zhenguang G Cai, David A. Haslett, Xufeng Duan,
765 Shuqi Wang, and Martin J. Pickering. 2023. *Does*
766 *ChatGPT resemble humans in language use? arXiv*
767 *preprint arXiv:2303.08014.*
- 768 Kimberly Wright Cassidy, Michael H. Kelly, and Lee'at
769 J. Sharoni. 1999. *Inferring gender from name*
770 *phonology. Journal of Experimental Psychology:*
771 *General* 128, no. 3 (1999): 362.
- 772 Noam Chomsky, Ian Roberts, and Jeffrey Watumull.
773 2023. *The false promise of ChatGPT. The New York*
774 *Times.*[https://www.nytimes.com/2023/03/08/opinio](https://www.nytimes.com/2023/03/08/opinion/noam-chomsky-chatgpt-ai.html)
775 [n/noam-chomsky-chatgpt-ai.html](https://www.nytimes.com/2023/03/08/opinion/noam-chomsky-chatgpt-ai.html)
- 776 Ryan Doran, Gregory Ward, Meredith Larson, Yaron
777 McNabb, and Rachel E. Baker. 2012. *A novel*
778 *experimental paradigm for distinguishing between*
779 *what is said and what is implicated. Language:* 124-
780 154.
- 781 Edward Gibson, Leon Bergen, and Steven T.
782 Piantadosi. 2013. *Rational integration of noisy*
783 *evidence and prior semantic expectations in*
784 *sentence interpretation. Proceedings of the National*
785 *Academy of Sciences* 110, no. 20: 8051-8056.
- 786 Herbert Paul Grice. 1975. *Logic and conversation. In*
787 *Speech acts*, pp. 41-58. Brill.
- 788 Herbert Paul Grice. 1978. *Further notes on logic and*
789 *conversation. In Pragmatics*, pp. 113-127. Brill
- 790 Laurence Robert Horn. 1972. *On the semantic*
791 *properties of logical operators in English.*
792 University of California, Los Angeles.
- 793 Wenxiang Jiao, Wenxuan Wang, Jen-tse Huang, Xing
794 Wang, and Zhaopeng Tu. 2023. *Is ChatGPT a good*
795 *translator? A preliminary study. arXiv preprint*
796 *arXiv:2301.08745.*
- 797 Michal Kosinski. 2023. *Theory of mind may have*
798 *spontaneously emerged in large language models.*
799 *arXiv preprint arXiv:2302.02083.*
- 800 Stephen C Levinson. 2000. *Presumptive meanings:*
801 *The theory of generalized conversational*
802 *implicature.* MIT press.
- 803 Kyle Mahowald, Evelina Fedorenko, Steven T.
804 Piantadosi, and Edward Gibson. 2013.
805 *Info/information theory: Speakers choose shorter*
806 *words in predictive contexts. Cognition* 126, no. 2
807 (2013): 313-318.
- 808 Ira A Noveck. 2001. *When children are more logical*
809 *than adults: Experimental investigations of scalar*
810 *implicature. Cognition* 78, no. 2: 165-188.
- 811 Ira A Noveck and Dan Sperber. 2007. *The why and how*
812 *of experimental pragmatics: the case of 'scalar*
813 *inferences, in Advances in Pragmatics*, ed N.
814 Burton-Roberts (Basingstoke: Palgrave), 184–212.
- 815 Miguel Ortega-Martín, Óscar García-Sierra, Alfonso
816 Ardoiz, Jorge Álvarez, Juan Carlos Armenteros, and
817 Adrián Alonso. 2023. *Linguistic ambiguity analysis*
818 *in ChatGPT. arXiv preprint arXiv:2302.06426*
- 819 Steven T Piantadosia. 2023. *Modern language models*
820 *refute Chomsky's approach to language. Lingbuzz*
821 *Preprint, lingbuzz/007180*

822 Jennifer M. Rodd, Belen Lopez Cutrin, Hannah Kirsch,
823 Alessandra Millar, and Matthew H. Davis. 2013.
824 Long-term priming of the meanings of ambiguous
825 words. *Journal of Memory and Language* 68, no. 2
826 (2013): 180-198.

827 Chris Westbury. 2005. Implicit sound symbolism in
828 lexical access: Evidence from an interference task.
829 *Brain and language* 93, no. 1 (2005): 10-19.

830 Arjen Zondervan. 2010. *Scalar implicatures or focus:
831 an experimental approach*. Netherlands Graduate
832 School of Linguistics.

833 A Appendix

834 An example of experimental items containing
835 GCIs of different categories in Exp1.

836

837

Dialogue	Fact	First Level Category	Second Level Category
Irene: Hey, Sam. Do you know who wrote <i>Pride and Prejudice</i> ? Sam: A British woman wrote it, and her last name was Austen.	FACT: Jane Austen, a British woman, wrote <i>Pride and Prejudice</i> .	Training Example	Training Example
Irene: How much cake did Gus eat at his sister's birthday party? Sam: He ate most of the cake.	FACT: By himself, Gus ate his sister's entire birthday cake.	Q_Based_GCIs	Quantifiers_Modals
Irene: How many children does Lisa have? Sam: Lisa has three children.	FACT: Lisa has quadruplets	Q_Based_GCIs	Cardinals
Irene: How would you say you're doing financially? Sam: I'm comfortable.	FACT: Sam just bought four condos at Lake Point Tower, in downtown Chicago, where Oprah Winfrey lives.	Q_Based_GCIs	Gradable_Adjectives
Irene: What kind of milk does your diet allow for? Sam: It allows for 1%.	FACT: The only type of milk prohibited by Sam's diet is full-fat milk.	Q_Based_GCIs	Rankings
Irene: I heard something big happened in the art studio yesterday. Sam: In a fit of rage, Rachel picked up a hammer and broke a statue.	FACT: After grabbing a hammer, Rachel angrily kicked a statue, causing it to fall over and break.	I_Based_GCIs	Argument_Saturation
Irene: What happened when Sue came over? Sam: She walked into the bathroom. The window was open.	FACT: The open windows are in the kitchen, and there are no windows in the bathroom.	I_Based_GCIs	Bridging_Inferences
Irene: Can the guys come to the reception? Sam: George and Steve play squash at the gym until 6:00 every day.	FACT: George plays squash at the YMCA until 6:00 daily, and Steve plays squash at SPAC until 6:00 every day.	I_Based_GCIs	Coactivities
Irene: I understand that George has had a really rough year. Sam: Last month, he lost his job and started drinking.	FACT: George started drinking on the 15th of last month and lost his job on the 20th of last month.	I_Based_GCIs	Conjunction_Buttrressing
Irene: Why is Stephen so upset? Sam: He caused Bill to die.	FACT: Stephen intentionally murdered Bill.	M_Based_GCIs	Verbal_Periphrasis
Irene: What happened at Doctor Witherspoon's office? Sam: Sasha waited and waited for her appointment.	FACT: Sasha waited 5 minutes for her appointment at DoctorWitherspoon's office.	M_Based_GCIs	Repeated_Verb_Conjuncts
Irene: What did Joseph do after finishing the marathon? Sam: He drank bottles and bottles of water.	FACT: Joseph drank one 20 oz bottle and one 16 oz bottle of water after finishing the marathon.	M_Based_GCIs	Repeated_Noun_Conjuncts