# Perceptions of Linguistic Uncertainty by Language Models and Humans

**Catarina Belem**[*,1], **Markelle Kelly**[*,1], **Mark Steyvers**[1,2], **Sameer Singh**[1], **Padhraic Smyth**[1],

[1]Department of Computer Science, University of California Irvine
[2]Department of Cognitive Sciences, University of California Irvine

## Abstract

*Uncertainty expressions* such as "probably" or "highly unlikely" are pervasive in human language. While prior work has established that there is population-level agreement in terms of how humans quantitatively interpret these expressions, there has been little inquiry into the abilities of language models in the same context. In this paper, we investigate how language models map linguistic expressions of uncertainty to numerical responses. Our approach assesses whether language models can employ theory of mind in this setting: understanding the uncertainty of another agent about a particular statement, independently of the model's own certainty about that statement. We find that 7 out of 10 models are able to map uncertainty expressions to probabilistic responses in a human-like manner. However, we observe systematically different behavior depending on whether a statement is actually true or false. This sensitivity indicates that language models are substantially more susceptible to bias based on their prior knowledge (as compared to humans). These findings raise important questions and have broad implications for human-AI and AI-AI communication.

## 1 Introduction

The expression of uncertainty is ubiquitous in human communication — in relaying predictions ("it is likely to rain tomorrow"), conveying imperfect knowledge ("I think I have a copy in my desk"), and describing unknown information ("the artifact could be more than 500 years old"). Expressing uncertainty is particularly critical in fields such as medicine, law, and politics, where statements including *uncertainty expressions* (*e.g.*, "likely," "doubtful") are frequently used to support medical, judicial, and political decisions ([Karelitz and

Budescu, 2004](#)). Domain experts use these expressions to communicate uncertainty across a variety of situations, such as the likelihood of side-effects of a medical treatment ([Sawant and Sansgiry, 2018](#); [Patt and Dessai, 2005](#)), the chances of a not-guilty verdict in legal cases ([Fore, 2019](#)), the probability of environmental events resulting from climate change ([Patt and Dessai, 2005](#); **?**), or the likelihood of emergence of military conflicts ([Duke, 2023](#)). Prior work has found that, in general, humans are well-attuned to the use of such uncertainty expressions, exhibiting population-level agreement in mapping these expressions to corresponding probabilities ([Wallsten et al., 1986a](#); [Willems et al., 20](#); [Fagen-Ulmschneider, 2019](#)).

However, the topic of how large language models (LLMs) interpret linguistic uncertainty has received relatively little attention. In particular, given text where a speaker expresses uncertainty about a particular statement, this paper investigates whether LLMs can interpret the uncertainty not as a function of the model's internal beliefs, but by objectively assessing the speaker's uncertainty about the statement. Consider the motivating example in Figure 1: when writing a headline for a statement qualified by the word "probable," ChatGPT expresses substantially different uncertainty depending on its prior belief about the statement.[1] In this example, ChatGPT is conflating the speaker's uncertainty with its own uncertainty about the statement—in effect, a failure of "theory of mind."

In this work, we investigate the abilities of LLMs to provide quantitative interpretations of uncertainty expressions, focusing in particular on how the prior knowledge of an LLM affects this ability. To this end, we propose to evaluate models' capabilities as a function of their ability to map text

---

*Authors contributed equally to this work. Correspondence: cbelem@uci.edu.

[1]When prompted about its belief about these statements, ChatGPT agrees with the first and disagrees with the second; see Figure 8 in the Appendix A.
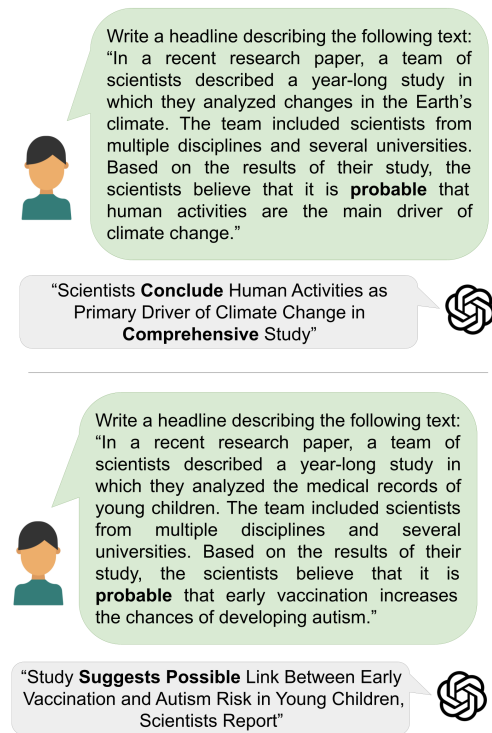
Figure 1: **Two interactions with `ChatGPT` (June 2024) concerning the generation of a headline for a short passage**. Both passages are structured identically and qualified with the word "probable," but the first is about climate change and the second about the link between vaccines and autism. For the first passage, `ChatGPT` generates a confident-sounding headline, using the words "conclude" and "comprehensive." The second headline is weaker, with words like "suggests" and "possible."

containing uncertainty expressions to numerical responses. We analyze the performance of both humans and 10 popular LLMs on this task, enabling direct comparison between humans and models.[2] We find that larger, newer models like `GPT-4` and `LLama3 (70B)` can consistently map uncertainty expressions to numerical responses that align with human population-level perceptions. However, we also show that the responses LLMs generate are susceptible to bias based on their prior knowledge—with much greater susceptibility than that of humans.

The models' sensitivity to their prior knowledge has concerning implications for the use of LLMs in tasks in which they must process and generate text containing uncertainty expressions, *e.g.*, summarizing scientific reports or writing news articles (Shao et al., 2024; Laban et al., 2024). When

---

[2]The code and data can be found at `https://github.com/UCIDataLab/llm-uncertainty-perceptions`.

an LLM's ability to quantify uncertainty can be "poisoned" by its beliefs, its downstream performance is dependent on its parametric or pretraining knowledge (which can be obsolete or wrong (Liang et al., 2022; Longpre et al., 2023)), rather than on critical contextual information (Longpre et al., 2021). Further, this means that the biases of a model (including the many well-documented potentially harmful biases of LLMs, *e.g.*, Wan et al. (2023); Kotek et al. (2023); Salewski et al. (2024); Scherrer et al. (2024); Motoki et al. (2024)) can subtly manifest in how it interprets and generates uncertainty language and, as a consequence, have broader implications for human-AI and AI-AI interactions.

## 2 Related Work

**Human Perceptions of Uncertainty Expressions.** In fields like medicine, finance, law, and politics, where it is impossible to make predictions with complete certainty, decisions are often informed by subjective probabilities (Karelitz and Budescu, 2004; Dhami and Wallsten, 2005; Fore, 2019). Subjective probabilities can be communicated quantitatively, *e.g.*, through numerical probabilities, odds, percentages, intervals, or qualitatively, through the use of uncertainty expressions or epistemological markers (*e.g.*, "I believe", "According to") (Dhami and Mandel, 2022). Although they are less precise than numerical values (Wallsten et al., 1986b; Brun and Teigen, 1988; **?**), humans generally prefer to use linguistic expressions to communicate uncertainty (Erev and Cohen, 1990; Wallsten et al., 1993).

Interested in the efficacy of how humans communicate uncertainty linguistically, researchers have examined how individuals map uncertainty expressions into numerical values across different fields and expertise levels (Windschitl and Wells (1996); Karelitz and Budescu (2004); Wallsten et al. (2008, 1986a); Fore (2019); *inter alia*). Although there can be considerable variation in responses at the individual level, these studies have revealed consistent and systematic patterns relating uncertainty expressions and numerical responses at the population level (Wallsten et al., 2008; Willems et al., 20; Fagen-Ulmschneider, 2019).

**Uncertainty Quantification in LLMs.** The need for more reliable LLMs has prompted researchers to investigate new methods for communicating the internal uncertainty of LLMs. Proposed methods

can be differentiated in terms of the information used to estimate the model's uncertainty: from token-level information (Jiang et al., 2021; Kuhn et al., 2023; Duan et al., 2024), to dissimilarities across multiple samples (Si et al., 2022; Chen and Mueller, 2023; Xiong et al., 2024; Hou et al., 2024; Lin et al., 2024; Aichberger et al., 2024), to training external classifiers using the inputs and/or LLMs' representations (Jiang et al., 2021; Mielke et al., 2022; Shrivastava et al., 2023), or even directly eliciting confidence estimates from LLMs as output tokens (Lin et al., 2022; Tian et al., 2023). Furthermore, several works have analyzed the impact of LLM-articulated uncertainty in human-AI interaction, finding that participants adjust their perception of LLMs' correctness when shown LLM outputs that include uncertainty expressions (Zhou et al., 2024; Kim et al., 2024; Steyvers et al., 2024). With the goal of calibrating human-AI interaction, Chaudhry et al. (2024) propose fine-tuning LLMs to convey uncertainty expressions that faithfully reflect their intrinsic uncertainty. While these works investigate how we can gauge LLMs' intrinsic uncertainty and how humans react to various uncertainty expressions in text, there has been far less work on the questions we focus on in this paper, *i.e.*, how LLMs interpret linguistic uncertainty and how closely these interpretations match those of humans.

**LLM Perceptions of Uncertainty Expressions.** A small body of recent work has begun to investigate the relationships between uncertainty expressions and model behavior. Recently, Yona et al. (2024) found low correlations between LLMs' intrinsic uncertainty and the use of uncertainty expressions. Sileo and Moens (2023) investigate whether LLMs are able to discriminate between two uncertainty expressions and reason in terms of compositions of expressions. However, the paper focuses on LLMs' binary rankings of expressions, as opposed to numerical interpretations, and does not compare LLM and human behavior. Most directly related to our work is that of Maloney et al. (2024) which compares numerical probability estimates from GPT-4 and humans using a small set of "context" prompts, and the work of Tang et al. (2024) who compare the numeric-textual mapping of uncertainty expressions across four different contexts in both Chinese and English. Our paper goes significantly beyond this work by assessing a broad range of LLMs using a more diverse and natural

set of contexts. We additionally conduct human experiments, assessing the performance of humans on the same task to facilitate a direct comparison between humans and LLMs. Further, our approach is designed to target "theory of mind"—the task requires humans and LLMs to quantify what an uncertainty expression reflects about the speaker's belief, rather than what the expression means to the human or LLM. Finally, our work is the first that we are aware of to investigate how LLMs can be biased by their prior knowledge in mapping uncertainty expressions to numerical responses.

## 3   Baseline Human Study

As a baseline for how people map uncertainty expressions to numerical probabilities, we first conducted an experiment in which 94 humans were shown uncertainty expressions and asked to provide corresponding numerical responses. We focused on a set of 14 uncertainty expressions (*e.g.*, "almost certain," "unlikely"—the full list is provided in Appendix C and is also shown on the y-axis in Figure 3), drawn from Wallsten et al. (1986a) and Wallsten et al. (2008). In this initial experiment, our goal is to assess how people perceive these uncertainty expressions "in the wild," putting them in the context of plausible real-world statements. In addition, we use types of statements that attempt to minimize the potential for people to conflate their own beliefs about these statements with their assessment of the confidence of the person making the statement.

To this end, we constructed a set of statements $(u, s, e)$ which include uncertainty expressions $u \in \mathcal{U}$ used by speakers $s \in \mathcal{S}$ to convey their degree of certainty about the truthfulness or falsehood of a statement or event $e \in \mathcal{E}$. By presenting statements as being made by a specific speaker $s$,[3] we are asking participants to use theory of mind to estimate how likely it is that the speaker believes that the statement is true. We then query participants about the speaker's degree of certainty, clearly distinguishing this notion from the participant's own beliefs. For instance, given the statement "Sonia believes it is unlikely it will rain today," we can ask participants to quantify with a numerical response how likely *Sonia* thinks it is that it will rain today, distinct from the participants' own beliefs about how likely it is to rain. We use the term *numerical*

---

[3] Names are selected arbitrarily from a pre-defined pool of names. Participants see each name once throughout the experiment. For more details, see Appendix C.

*response* throughout the paper to refer to the participant's numerical assessment of the likelihood that speaker $s$ believes statement $e$ to be true when using uncertainty expression $u$. We instruct participants that responses should be numbers between 0 and 100, where 0 implies that speaker $s$ believes there is a $0\%$ chance that statement $e$ is true while 100 implies the speaker believes there is a $100\%$ chance that the statement is true.

In our baseline experiment, we use *non-verifiable* statements to separate the meaning of the uncertainty expressions from uncertainty about the statements themselves. Non-verifiable statements are statements that are not sufficiently grounded with specific contextual information to allow an external observer to be confident in either the truth or falsity of the statement. For example, in the context of a prompt such as "Maria believes it is likely that [statement]," we consider statements such as *her boss has two pets* or *her flight will land around 6pm* non-verifiable—as there is not enough context for an observer to be able to assess the likelihood that the statement is true. In contrast, *verifiable* statements (which we discuss further in Section 4.1) can be verified as correct or incorrect in a context-free sense (e.g., *the capital city of Peru is Lima*); humans and LLMs will often have strong prior beliefs about the likelihood that such statements are true.

We manually constructed a set of 60 non-verifiable statements and systematically combined these with 14 uncertainty expressions. We randomly selected speaker names, generating sentences describing the belief of a hypothetical speaker in the form: "[Speaker] believes it is [uncertainty expression] that [statement]." For each sentence, participants were asked to quantify the speaker's belief about the statement. In particular, they were asked what the probability is *from the speaker's perspective* that the statement is correct. Participants then provided their response quantized to numerical bins $0, 5, 10, \ldots, 95, 100$ (see example in Figure 2). Each of the 94 participants provided annotations for 28 randomly chosen statements (and speaker names). As a result, every uncertainty expression was annotated twice by each participant.

The outcome of this experiment[4] is an empirical distribution over the $2 \times 94$ numerical responses that participants associated with each uncertainty
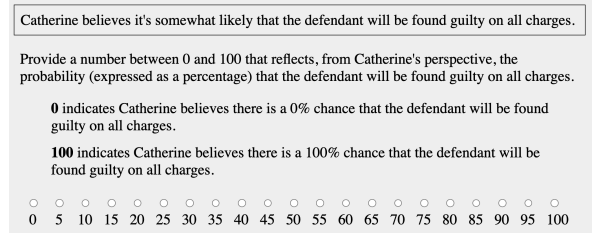
---

Figure 2: **Example of a non-verifiable statement provided to participants in the baseline experiment**. Each example uses a unique name and statement. Participants see one question at a time.
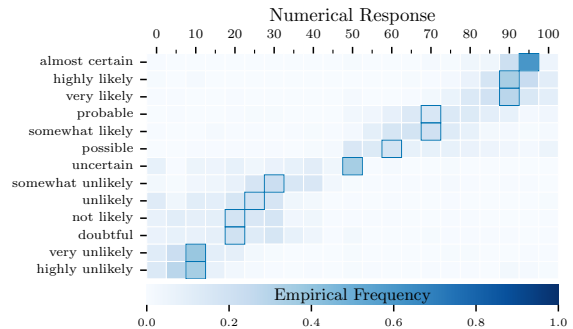


Figure 3: **Human empirical distributions of numerical responses per uncertainty expression in the non-verifiable setting**. Highlighted blue boxes represent the mode value for each expression. Overall, population-level perceptions increase monotonically with the use of more confident uncertainty expressions.

expression. For example, Figure 9 reflects the histogram of responses assigned to the phrases "very likely" and "very unlikely"; results for all 14 uncertainty expressions can be summarized in a heatmap as shown in Figure 3. We note that our results are in agreement with prior work on human perceptions of uncertainty expressions (Wallsten et al., 2008, 1986a; Willems et al., 20), including consistent ordering in aggregate population patterns, in terms of the mode of the empirical distributions.

## 4 Methodology

Our full set of experiments includes experiments with both humans and LLMs as participants and with both non-verifiable (NV) and verifiable (V) statements. The baseline experiment, as described in Section 3, consists of human participants and non-verifiable statements (denoted human+NV). In this section, we describe extensions to include verifiable statements and LLM participants, resulting in three additional sets of experiments: human+V, LLM+V, LLM+NV. To help us draw comparisons between the three additional settings and the baseline

experiment, we conclude this section with the description of two metrics.

## 4.1 Verifiable Statements

In addition to the non-verifiable statements described in Section 3, our dataset also includes *verifiable* statements, for the purpose of assessing the effects of prior knowledge on quantifying linguistic uncertainty. To increase the chances that both LLMs and humans are familiar with the verifiable statements, we focus on concise, general-knowledge statements based on widely recognized facts (*e.g.*, geography, history of art, science). Specifically, we create 60 verifiable statements based on a multiple-choice question-answering trivia dataset from The Question Company.[5] Starting with 30 of the dataset's "easy" questions and corresponding multiple-choice options, we write *true* statements that use the correct answer and *false* statements using one of the incorrect answers. Examples of verifiable statements and additional details about the dataset are included in Appendix B. We focus on results for both humans and LLMs with these 60 verifiable statements in the main paper but include a validation analysis using 400 additional statements in Section 5.3. These additional statements are extracted via a similar procedure from AI2-ARC (Clark et al., 2018), a grade-school level, multiple-choice science question answering dataset used to evaluate state-of-the-art LLMs' reasoning capabilities (Beeching et al., 2023; Jiang et al., 2023; Achiam et al., 2024).

## 4.2 Numerical responses from LLMs

To obtain uncertainty estimates from LLMs, we create prompts similar to the queries provided to humans (see Appendix C). Our goal is to estimate an empirical distribution, per uncertainty expression $u$, over each LLM's generated responses, in a manner similar to how empirical distributions for humans are generated (*e.g.*, see Figure 3). In the results in this paper we focus on greedy decoding, where we select the numerical response that has the highest probability in the next-token probability distribution generated by the LLM conditioned on the prompt, *i.e.*, decoding with temperature=0. Because this sampling approach requires no knowledge about the weights or next-token probabilities, it is applicable to any model, including those behind black-box APIs, such as Gemini (Anil et al.,

2024) and GPT-4 (Achiam et al., 2024). While focusing primarily on this greedy sampling approach allow us to efficiently compare the modal behavior of different LLM model families in equal terms (regardless of the available information), it does not provide insight into distributional behavior. Thus, in Section 5 we include results using probabilistic decoding with temperature=1 to evaluate the sensitivity of our conclusions to decoding method. Additional information on the extraction methodologies used can be found in Appendix D.

## 4.3 Metrics

We treat the empirical distribution obtained for the non-verifiable statements with human participants (described in Section 3) as our *reference distribution* for evaluation purposes, since it reflects human perceptions of uncertainty expressions in a setting that is designed to be free of prior information or biases about the corresponding statements. For every uncertainty expression $u \in \mathcal{U}$, we define a *reference conditional probability distribution* $P(k|u)$, $k = 0, 5, 10, \ldots, 95, 100$, where $P(k|u)$ is the empirical distribution from the baseline experiment. Given a response from any agent, human or LLM, in the context of a particular uncertainty expression $u$ we measure the quality of the response using the reference distribution $P(k|u)$.

The primary quality metric that we propose is **Proportional Agreement (PA)**. PA can be defined as follows: if an agent's response matches bin $k$ for uncertainty expression $u$, then the PA value for that response is defined as $P(k|u)$, where $P$ is the reference (population) distribution defined above. Intuitively, for an expression $u$, this PA score $P(k|u)$ represents the probability that the agent's response $k$ agrees with that of a randomly selected individual, and is upper bounded for any expression by $\arg\max_k P(k|u)$, *i.e.*, by the mode of the $P(k|u)$ values. The higher the PA value, the better the quality of the response in terms of agreement with the aggregate human population (as reflected by $P(k|u)$). To compute a single score for a particular LLM or individual human, we average the PA score over multiple responses and over the 14 uncertainty expressions.[6]

---

[6]Note that the PA metric is similar to the log-probability metric widely used to score probabilistic models in machine learning. However, it is not a likelihood in the sense that a likelihood corresponds to measuring the probability mass a model assigns to an observed outcome. Thus, in this non-likelihood context it is appropriate to average the PA scores directly (rather than taking products of probabilities as would

One drawback of using the PA metric is that it penalizes deviations from the reference distribution's mode. As a result, it may fail to capture the nuanced distributional differences that emerge when conveying uncertainty about events with varying severity or base rate (Wallsten et al., 1986a; Weber and Hilton, 1990; Willems et al., 20). However, in our experiments, we do not expect to observe systematic distributional differences, since the uncertainty expressions are used in neutral contexts that concern unknown people and extraneous events/facts. Nonetheless, we include an alternative to the PA metric that compares histograms of responses, *e.g.*, based on multiple responses from agents for a particular uncertainty expression $u$. In particular, we provide numerical results for histogram comparisons (using the Wasserstein distance between histograms) in Appendix E.

As an additional measure of alignment between the reference distribution and the agent's distribution, we also compute the **Mean Absolute Error (MAE)**, for each uncertainty expression $u$, defined as the absolute difference between (i) the mean of the responses across statements involving $u$ for an agent, and (ii) the mean of the reference distribution for $u$, $P(k|u)$. We then average across the 14 expressions $u$ to get a single score per agent.

## 5 Results

This section examines the ability of several well-known LLMs to interpret uncertainty expressions. We begin by assessing models' abilities to produce numerical responses that resemble human-like trends (*e.g.*, higher numerical responses assigned to higher-certainty expressions and vice-versa). We then study the effect of prior knowledge in the perception of uncertainty of both humans and models. We conclude with an assessment of the generalizability of our findings.

### 5.1 How well do LLMs perceive uncertainty?

As established in prior work (Wallsten et al., 1986b; Fagen-Ulmschneider, 2019; Willems et al., 20) and in our baseline experiment (Section 3), humans show population-level agreement in mapping uncertainty expressions to numerical responses. In this section, we assess whether LLMs possess a similar ability to ascribe numerical responses to uncertainty expressions. To this end, we prompt LLMs to provide responses for the same non-verifiable
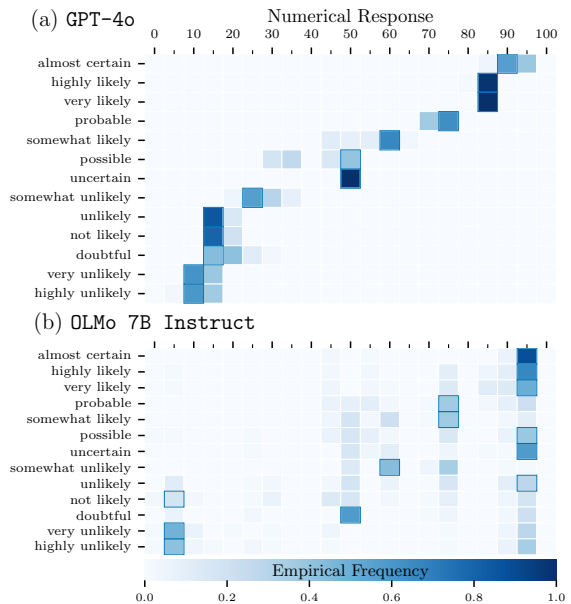
be done under an IID likelihood assumption).



Figure 4: **Model empirical distributions of numerical responses per uncertainty expression in the non-verifiable setting (LLM+NV).** Highlighted boxes represent the mode value for each expression. Even though we found no evidence of explicit instruction tuning datasets focusing on uncertainty estimation tasks, these results suggest that GPT-4o generally manifests human-like behavior, whereas OLMo (7B) does not.

(NV) statements as in the baseline experiment. Figure 4 shows the expression-wise histograms for these responses for two LLMs: GPT-4o and OLMo (7B), each of which can be directly compared to the histogram for humans in Figure 3.

Visually, we observe that GPT-4o matches the human distributions well, with smaller variance per distribution, while OLMo (7B) is less aligned. In Figure 14 in Appendix E, we show that most LLMs map uncertainty expressions to numerical responses in a manner consistent with human behavior, with higher values for expressions that are perceived by humans as higher-certainty (*e.g.*, "almost certain," "highly likely") and lower values for lower-certainty expressions (*e.g.*, "very unlikely"). Only two of the LLMs evaluated, OLMo (7B) and Gemma (2B), fail to reproduce this "increasing" pattern across expressions. One clear difference between humans and LLMs is that the conditional distributions of LLMs have lower entropy (or variance) relative to the human distributions, with the LLM distributions tending to be much more concentrated[7] on a small number of responses com-

---

[7]Appendix E.3 shows evidence that greedy decoding results in lower variability in estimated distributions compared to humans. Additionally, while probabilistic decod-

Table 1: **Human-LLM agreement for non-verifiable statements**. Average Proportional Agreement (PA), PA as a fraction of the *Human Mode* results (% PA), and absolute error between mean responses (MAE). *Human Mode* represents the mode of the human NV distribution, whereas *Human Individual* represents the PA score of individual human responses relative to the population.

| | PA | % PA | MAE |
|---|---|---|---|
| Human Mode | 27.6 | — | — |
| Human Individual | 17.6 | 63.8 | 8.91 |
| ChatGPT | 19.7 | 71.4 | 6.80 |
| GPT-4 | 24.4 | 88.4 | 4.64 |
| GPT-4o | 18.9 | 68.5 | 5.58 |
| Gemini | 25.4 | 92.0 | 4.09 |
| LLama3 (8B) | 17.8 | 64.5 | 11.99 |
| LLama3 (70B) | 23.6 | 85.5 | 5.56 |
| Mixtral 8x7B | 21.8 | 79.0 | 5.88 |
| Mixtral 8x22B | 21.8 | 79.0 | 7.20 |
| OLMo (7B) | 11.1 | 40.2 | 21.41 |
| Gemma (2B) | 8.1 | 29.3 | 20.17 |

Table 2: **Human-LLM agreement for verifiable statements**. Average Proportional Agreement (PA), absolute error between mean responses (MAE), and the difference between these scores and those from the non-verifiable statements (Table 1) ($\Delta$ PA and $\Delta$ MAE, respectively). Again, *Human Mode* represents the mode of the human NV distribution, whereas *Human Individual* represents the average behavior across individual humans on the verifiable setting.

| | PA | $\Delta$ PA | MAE | $\Delta$ MAE |
|---|---|---|---|---|
| Human Mode | 27.6 | — | — | — |
| Human Individual | 16.7 | -0.9 | 9.35 | 0.44 |
| ChatGPT | 15.3 | -4.4 | 8.57 | 1.77 |
| GPT-4 | 22.1 | -2.3 | 3.84 | -0.80 |
| GPT-4o | 15.2 | -3.7 | 7.05 | 1.47 |
| Gemini | 21.3 | -4.1 | 7.23 | 3.14 |
| LLama3 (8B) | 10.1 | -7.7 | 16.59 | 4.60 |
| LLama3 (70B) | 18.9 | -4.7 | 13.73 | 8.17 |
| Mixtral 8x7B | 15.2 | -6.6 | 12.23 | 6.35 |
| Mixtral 8x22B | 18.6 | -3.2 | 9.78 | 2.58 |
| OLMo (7B) | 7.6 | -3.5 | 33.66 | 12.25 |
| Gemma (2B) | 5.3 | -2.8 | 25.07 | 4.9 |

pared to the variance in responses from a population of humans.

These observations are reflected more precisely by the PA scores in Tables 1 and 9. We observe that larger and newer LLMs (in particular, GPT-4, LLama3 (70B), and Gemini) perform especially well on this task under the PA metric, being at 85% or above in terms of matching the modal scores that a human population assigns to each uncertainty expression. In fact, 7 out of the 10 LLMs evaluated are significantly better matched to population modal responses than are individual humans on average[8]. This aligns with the high-level findings of Maloney et al. (2024), in particular, that the difference between the numerical responses of GPT-4 and humans were similar to (or smaller than) inter-human differences. In the context of our experiments, these high scores reflect that LLMs tend to be more consistent than individual humans in terms of agreement with aggregate human responses.

The MAE scores in Table 1 (lower is better) are highly anti-correlated with the PA scores and tell a similar story in terms of which models perform better. To provide a sense of scale, the MAE numbers are lower-bounded by 0 and upper-bounded by 25 (the expected MAE for random responses).

## 5.2 Does knowledge affect uncertainty perceptions of LLMs?

In this section, we assess the extent to which LLMs, and humans, are biased by their prior knowledge or beliefs in mapping uncertainty expressions to numerical responses. To investigate this question we collect responses from humans and LLMs on our verifiable (V) dataset, which includes both true and false common-knowledge statements (based on correct or incorrect answers, respectively, to multiple-choice questions). We find that average PA scores for both humans and LLMs are systematically lower for verifiable statements compared to the non-verifiable responses (Tables 2 and 10). This suggests that prior knowledge about a statement[9] makes it more difficult to quantify the beliefs of someone else about that statement. While humans show a small drop in their PA score, this reduction in PA is particularly pronounced for LLMs: all 10 LLMs demonstrated a significant reduction in PA, averaging a 4.3 point drop in score (across all models), compared to a 0.9 point drop for humans.

To investigate these differences in more detail, we consider the mean response values produced by the 6 models exhibiting the highest PA score. These values differ systematically depending on whether the statement is true or false: across the 6 LLMs in Figure 5, the mean response is 7.0 points

---

ing (temperature=1) generally increases variability for most models, this effect is not observed for GPT-4.

[8]The average performance of individual humans is represented by the Human Individual row in Table 1.

[9]We validempirically ated our use of true correctness as a proxy for the LLMs' beliefs in Appendix B.3.
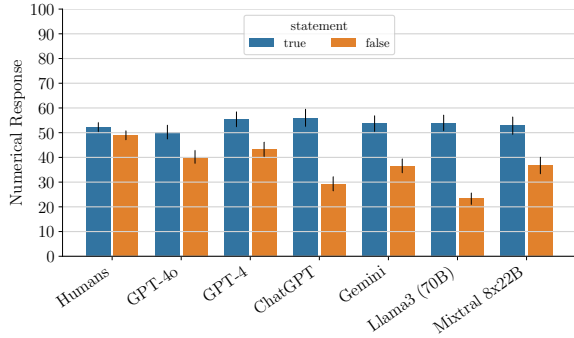
Figure 5: **Mean numerical response for the verifiable statements discriminated by truthfulness of statements**. The mean numerical responses produced by LLMs when evaluated in the context of true statements is significantly larger than when evaluated with the false statements. This difference is much larger in magnitude than the difference shown by a human population.
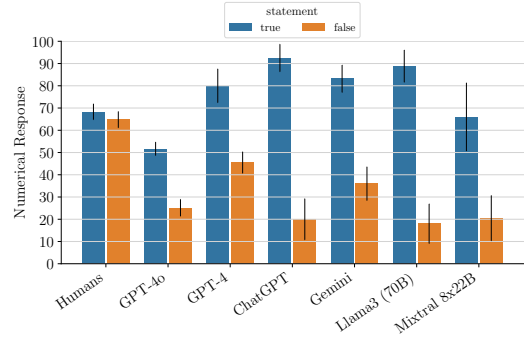
lower for false than true statements. This indicates that the models assign higher response values to the same uncertainty expression when they believe the associated statement is true than when they believe it to be false (providing a quantitative validation of the ChatGPT example in Section 1).

Results for a subset of specific uncertainty expressions are shown in Figure 6. We observe that the prior-knowledge bias is remarkably different depending on the uncertainty expression: the difference between true and false statements is much higher (49.5 percentage points) for the expression "possible" than for the expression "uncertain" (for this expression most models are close to the average 5.7 percentage point difference).
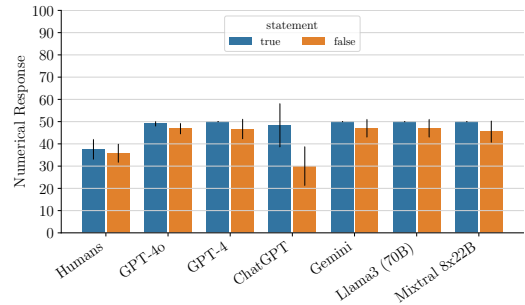
Overall, we find that all of the evaluated LLMs demonstrate significant biases based on their prior knowledge, well beyond those of humans. Our results indicate that when an LLM believes a statement is false, it tends to perceive the speaker's certainty as low, regardless of the actual uncertainty expressed by the speaker (and vice versa), lacking "theory of mind." Our findings in general are also consistent with recent results in the context of LLM reading comprehension, where the efficacy of LLM comprehension is sensitive to the degree of prior knowledge it has about the text it is analyzing (Basmov et al., 2024).

### 5.3 How generalizable are our findings?

In the previous sections, our analyses are conducted on a manually curated set of 120 statements, comprised of 60 NV statements and 60 V statements. To further validate our findings concerning LLMs'



(a) "possible"



(b) "uncertain"

Figure 6: **Mean response for verifiable statements discriminated by truthfulness of statements for three different uncertainty expressions**. We observe that the differences in mean numerical responses differ by uncertainty expression. Despite changes in magnitude, ChatGPT systematically exhibits the largest gap.

prior knowledge biases, we re-assess the impact of knowledge in LLMs' perceptual capabilities by obtaining their responses for 400 additional verifiable statements. Similarly to the original study, Figure 15 shows that, on average, all models except Gemma (2B) exhibit significant perceptual differences between true and false statements—between 5.87 (OLMo (7B)) and 17.26 (LLama3 (70B)) percentage points. While the magnitude of the difference is different across the two datasets (potentially due to semantic differences between the two QA datasets used to curate the verifiable statements), the directionality of the results with this larger dataset nonetheless corroborates our knowledge bias finding by showing that these perceptual differences persist in a different context. See Appendix F for a more detailed description of the experimental setup and additional results.

### 5.4 How does decoding impact our findings?

The previous analyses employ greedy decoding (i.e., temperature=0) when obtaining numerical responses from the LLMs. In this section, we investigate the impact of the decoding technique in the

Table 3: **Differences in average proportional agreement and mean responses from non-verifiable to verifiable settings when considering probabilistic decoding (`temperature=1`).** Even with a different decoding, we observe the same decrease in LLM perceptions when comparing non-verifiable with verifiable settings.

|         | $\Delta$ PA | $\Delta$ MAE |
|---------|-------------|--------------|
| ChatGPT | -3.6        | 0.4          |
| GPT-4   | -3.0        | 1.9          |
| GPT-4o  | -4.2        | -0.6         |

models' abilities to perceive linguistic uncertainty, by considering richer probability information (*i.e.*, `temperature=1`) when obtaining the response.[10]

Table 3 summarizes the change in agreement between LLM and human responses between the verifiable and non-verifiable settings (in terms of change in PA and MAE) when using probabilistic decoding. Validating the results reported in Section 5.2 with greedy decoding, we observe a clear difference in PA score between non-verifiable and verifiable statements when using probabilistic decoding. Further, comparing responses across true and false statements, we observe large mean response differences of 11.4, 11.7, and 27.9 percentage points for GPT-4o, GPT-4, and ChatGPT, respectively (see Figure 7 and Figure 18 in the Appendix for a breakdown across expressions). Although GPT-4o mean responses are considerably lower than in the greedy decoding setting (with a 20 percentage points drop), the gap between true and false statements persists. Ultimately, this analysis confirms the robustness of our previous findings to the decoding strategy.

## 6 Discussion

**Theory of mind.** A growing body of work aims to assess the *theory of mind* capabilities of LLMs in different contexts (*e.g.*, (Street et al., 2024; Verma et al., 2024; Sap et al., 2022; Zhou et al., 2023)). The task of mapping uncertainty expressions to numerical responses, from the perspective of some speaker, is one component of a general theory of mind ability. Our results indicate that LLMs have room for improvement in this area, in particular, that they are prone to confusing their own belief about a statement with the belief of someone else.



Figure 7: **Mean response on the verifiable statements discriminated by truthfulness of statements when decoding probabilistically `temperature=1`.** Despite being to a less extent, the prior-knowledge bias remains larger than that exhibited by human population.

**Connection to human behavior simulation using LLMs.** Our experiments reveal that, despite agreeing with population-level perceptions of linguistic uncertainty, models do not capture the full diversity of human behavior. While this lack of variability depends on the decoding algorithm used to extract the numerical response, we note that, for popular models like GPT-4, the disparity persists regardless of the chosen algorithm. Given the recent interest in using LLMs to simulate human participants (Aher et al., 2023; Gui and Toubia, 2023; Dillion et al., 2023; Park et al., 2023; Namikoshi et al., 2024), our work raises important questions about whose opinions and behaviors are being simulated (Santurkar et al., 2023; Motoki et al., 2023) and reveals a new dimension in which human and model diversity can differ.

## 7 Conclusions

We evaluate the abilities of LLMs to interpret uncertainty in language and evaluate numerous models in this context. Our results show that many LLMs can competently map uncertainty expressions to numerical responses in a way that aligns with population-level human perceptions, although the responses they choose can be much less diverse than those by humans. Additionally, we find that LLMs are more susceptible to conflating their own uncertainty about a statement with the statement speaker's uncertainty, resulting in output that is biased by the LLM's belief about the statement. By highlighting systematic inconsistencies related to the perceptions of linguistic uncertainty in the presence of knowledge, we shed light into overlooked model behaviors that are critical for understanding human-AI communication and downstream LLM performance.

---

[10]This analysis requires full next-token probability information from an LLM, which is prohibitively expensive to obtain empirically through sampling as it would require a large sample size (per $(u, s, e)$) to obtain accurately. As a result, we limit our analysis to OpenAI models for which the top 20 next-token probabilities are available. See Appendix G for additional discussion on this topic.
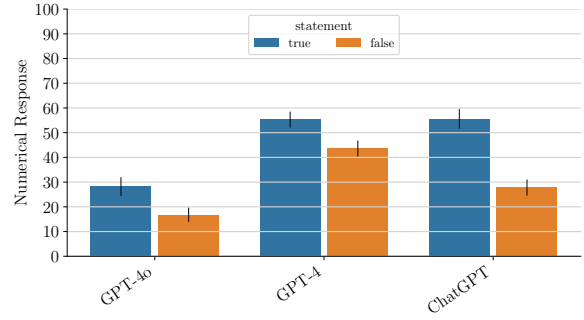
## Limitations

*Biases in Interpreting Uncertainty Expressions:*
Prior work has raised several concerns about the consistency of humans' interpretations of uncertainty expressions, demonstrating that they are subject to a number of biases and nuances. For example, people may conflate the speaker's confidence with the speaker's estimated uncertainty (Fleiner and Vennekens, 2024), statements worded in terms of confidence ("I am almost certain") or likelihood ("I believe it is almost certain") are interpreted as primarily communicating different types of uncertainty (epistemic and aleatoric, respectively) (Ülkümen et al., 2016), and these statements are directional, emphasizing either the occurrence or non-occurrence of an event (Teigen, 2023). Further, these interpretations are context-dependent, affected by factors including the individual's perception of the speaker and the severity of the event in question (Budescu and Wallsten, 1985; Collins and Hahn, 2018; **?**; Juanchich and Sirota, 2013; Brun and Teigen, 1988; Weber and Hilton, 1990). Thus, in interpreting uncertainty expressions, high variability can occur both across (Zhang et al., 2023; Collins and Hahn, 2018) and within (Clarke et al., 1992; Van Der Bles et al., 2019) individuals. In this paper we do not explore these dimensions of the interpretation of uncertainty expressions, which could limit the generality of our conclusions.

*US Centric View:* In this paper we focus on a small set of uncertainty expressions in English and our baseline is drawn from participants located in the United States. Investigating the role that cultural and language differences play in communicating uncertainty is important future work that will help better characterize the downstream abilities of LLMs for all users. The recent work by Tang et al. (2024) represents a first step in this direction, but it does not account for the cultural context that would inform the meaning of a speaker (Huang and Yang, 2023). In particular, the meaning of "probable" and "very likely" may differ significantly between English and Chinese. Such cultural difference may suffice to explain the observed differences in perceptions measured between US population and GPT-4 when prompted with Chinese contexts.

*Lack of Explanation:* Our results highlight the LLMs' abilities to interpret uncertainty phrases in a way that agrees with population-level human distribution in the non-verifiable and to less extent in the verifiable setting. It is not clear why models have acquired this general ability, given that there are not similarly framed tasks in available instruction-tuning and human feedback datasets (Wang et al., 2022; Bai et al., 2022). We hope that future work will explore the origins of this behavior.

## References

Josh Achiam, Steven Adler, Sandhini Agarwal, et al. 2024. GPT4 Technical Report. *Preprint*, arXiv:2303.08774.

Gati Aher, Rosa I. Arriaga, and Adam Tauman Kalai. 2023. Using large language models to simulate multiple humans and replicate human subject studies. In *Proceedings of the 40th International Conference on Machine Learning*, ICML'23. JMLR.org.

Lukas Aichberger, Kajetan Schweighofer, Mykyta Ielanskyi, and Sepp Hochreiter. 2024. Semantically diverse language generation for uncertainty estimation in language models. *Preprint*, arXiv:2406.04306.

Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, et al. 2024. Gemini: A family of highly capable multimodal models. *Preprint*, arXiv:2312.11805.

Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, Nicholas Joseph, Saurav Kadavath, Jackson Kernion, Tom Conerly, Sheer El-Showk, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Tristan Hume, Scott Johnston, Shauna Kravec, Liane Lovitt, Neel Nanda, Catherine Olsson, Dario Amodei, Tom Brown, Jack Clark, Sam McCandlish, Chris Olah, Ben Mann, and Jared Kaplan. 2022. Training a helpful and harmless assistant with reinforcement learning from human feedback. *Preprint*, arXiv:2204.05862.

Victoria Basmov, Yoav Goldberg, and Reut Tsarfaty. 2024. LLMs' reading comprehension is affected by

parametric knowledge and struggles with hypothetical statements. *ArXiv preprint*, abs/2404.06283.

Edward Beeching, Clémentine Fourrier, Nathan Habib, Sheon Han, Nathan Lambert, Nazneen Rajani, Omar Sanseviero, Lewis Tunstall, and Thomas Wolf. 2023. Open llm leaderboard. `https://huggingface.co/spaces/open-llm-leaderboard/open_llm_leaderboard`.

Wibecke Brun and Karl Halvor Teigen. 1988. Verbal probabilities: Ambiguous, context-dependent, or both? *Organizational Behavior and Human Decision Processes*, 41(3):390–404.

David V Budescu and Thomas S Wallsten. 1985. Consistency in interpretation of probabilistic phrases. *Organizational behavior and human decision processes*, 36(3):391–405.

Arslan Chaudhry, Sridhar Thiagarajan, and Dilan Gorur. 2024. Finetuning language models to emit linguistic expressions of uncertainty. *ArXiv preprint*, abs/2409.12180.

Jiuhai Chen and Jonas Mueller. 2023. Quantifying uncertainty in answers from any language model and enhancing their trustworthiness. *Preprint*, arXiv:2308.16175.

Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. Think you have solved question answering? try arc, the ai2 reasoning challenge. *ArXiv preprint*, abs/1803.05457.

Valerie A Clarke, Coral L Ruffin, David J Hill, and Arthur L Beamen. 1992. Ratings of orally presented verbal expressions of probability by a heterogeneous sample. *Journal of Applied Social Psychology*, 22(8):638–656.

Peter J. Collins and Ulrike Hahn. 2018. Chapter three - communicating and reasoning with verbal probability expressions. volume 69 of *Psychology of Learning and Motivation*, pages 67–105. Academic Press.

Mandeep K. Dhami and David R. Mandel. 2022. Communicating uncertainty using words and numbers. *Trends in Cognitive Sciences*, 26(6):514–526.

Mandeep K. Dhami and Thomas S. Wallsten. 2005. Interpersonal comparison of subjective probabilities: Toward translating linguistic probabilities. *Memory & Cognition*, 33(6):1057–1068.

Danica Dillion, Niket Tandon, Yuling Gu, and Kurt Gray. 2023. Can ai language models replace human participants? *Trends in Cognitive Sciences*, 27:597–600.

Jinhao Duan, Hao Cheng, Shiqi Wang, Alex Zavalny, Chenan Wang, Renjing Xu, Bhavya Kailkhura, and Kaidi Xu. 2024. Shifting attention to relevance: Towards the predictive uncertainty quantification of free-form large language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics*, Bangkok, Thailand. Association for Computational Linguistics.

Misty C. Duke. 2023. Probability and confidence: How to improve communication of uncertainty about uncertainty in intelligence analysis. *Journal of Behavioral Decision Making*, 37(1).

Ido Erev and Brent L Cohen. 1990. Verbal versus numerical probabilities: Efficiency, biases, and the preference paradox. *Organizational Behavior and Human Decision Processes*, 45(1):1–18.

Wade Fagen-Ulmschneider. 2019. Perception of probability words. Accessed: [June 12, 2024].

Christian Fleiner and Joost Vennekens. 2024. Towards effective management of verbal probability expressions using a co-learning approach. In *HHAI 2024: Hybrid Human AI Systems for the Social Good*, pages 124–133. IOS Press.

Joe Fore. 2019. "a court would likely (60-75%) find...." defining verbal probability expressions in predictive legal analysis. *Legal Comm. & Rhetoric: JAWLD*, 16:49.

Zorik Gekhman, Gal Yona, Roee Aharoni, Matan Eyal, Amir Feder, Roi Reichart, and Jonathan Herzig. 2024. Does fine-tuning llms on new knowledge encourage hallucinations? *ArXiv preprint*, abs/2405.05904.

George Gui and Olivier Toubia. 2023. The challenge of using llms to simulate human behavior: A causal inference perspective. *ArXiv preprint*, abs/2312.15524.

Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2020. The curious case of neural text degeneration. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.

Bairu Hou, Yujian Liu, Kaizhi Qian, Jacob Andreas, Shiyu Chang, and Yang Zhang. 2024. Decomposing uncertainty for large language models through input clarification ensembling.

Jing Huang and Diyi Yang. 2023. Culturally aware natural language inference. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 7591–7609, Singapore. Association for Computational Linguistics.

Maor Ivgi, Ori Yoran, Jonathan Berant, and Mor Geva. 2024. From loops to oops: Fallback behaviors of language models under uncertainty. *ArXiv preprint*, abs/2407.06071.

Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, et al. 2023. Mistral 7b. *Preprint*, arXiv:2310.06825.

Zhengbao Jiang, Jun Araki, Haibo Ding, and Graham Neubig. 2021. How can we know when language models know? on the calibration of language models for question answering. *Transactions of the Association for Computational Linguistics*, 9:962–977.

Marie Juanchich and Miroslav Sirota. 2013. Do people really say it is "likely" when they believe it is only "possible"? effect of politeness on risk communication. *Quarterly Journal of Experimental Psychology*, 66(7):1268–1275. PMID: 23782394.

Saurav Kadavath, Tom Conerly, Amanda Askell, Tom Henighan, Dawn Drain, Ethan Perez, Nicholas Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli Tran-Johnson, Scott Johnston, Sheer El-Showk, Andy Jones, Nelson Elhage, Tristan Hume, Anna Chen, Yuntao Bai, Sam Bowman, Stanislav Fort, Deep Ganguli, Danny Hernandez, Josh Jacobson, Jackson Kernion, Shauna Kravec, Liane Lovitt, Kamal Ndousse, Catherine Olsson, Sam Ringer, Dario Amodei, Tom Brown, Jack Clark, Nicholas Joseph, Ben Mann, Sam McCandlish, Chris Olah, and Jared Kaplan. 2022. Language models (mostly) know what they know. *ArXiv preprint*, abs/2207.05221.

Tzur M. Karelitz and David V. Budescu. 2004. You say "probable" and i say "likely": Improving interpersonal communication with verbal probability phrases. *Journal of Experimental Psychology: Applied*, 10(1):25–41.

Sunnie S. Y. Kim, Q. Vera Liao, Mihaela Vorvoreanu, Stephanie Ballard, and Jennifer Wortman Vaughan. 2024. "i'm not sure, but...": Examining the impact of large language models' uncertainty expression on user reliance and trust. In *Proceedings of the 2024 ACM Conference on Fairness, Accountability, and Transparency*, FAccT '24, page 822–835, New York, NY, USA. Association for Computing Machinery.

Hadas Kotek, Rikker Dockum, and David Sun. 2023. Gender bias and stereotypes in large language models. In *Proceedings of The ACM Collective Intelligence Conference*, pages 12–24.

Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. 2023. Semantic uncertainty: Linguistic invariances for uncertainty estimation in natural language generation. In *Proceedings of the 11th International Conference on Learning Representations*, ICLR'23.

Philippe Laban, Alexander R. Fabbri, Caiming Xiong, and Chien-Sheng Wu. 2024. Summary of a haystack: A challenge to long-context llms and rag systems. *ArXiv preprint*, abs/2407.01370.

Weixin Liang, Girmaw Abebe Tadesse, Daniel Ho, L. Fei-Fei, Matei Zaharia, Ce Zhang, and James Zou. 2022. Advances, challenges and opportunities in creating data for trustworthy AI. *Nature Machine Intelligence*, 4(8):669–677.

Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. Teaching models to express their uncertainty in words. *Preprint*, arXiv:2205.14334.

Zhen Lin, Shubhendu Trivedi, and Jimeng Sun. 2024. Generating with confidence: Uncertainty quantification for black-box large language models. *Preprint*, arXiv:2305.19187.

Shayne Longpre, Kartik Perisetla, Anthony Chen, Nikhil Ramesh, Chris DuBois, and Sameer Singh. 2021. Entity-based knowledge conflicts in question answering. *ArXiv preprint*, abs/2109.05052.

Shayne Longpre, Gregory Yauney, Emily Reif, Katherine Lee, Adam Roberts, Barret Zoph, Denny Zhou, Jason Wei, Kevin Robinson, David Mimno, et al. 2023. A pretrainer's guide to training data: Measuring the effects of data age, domain coverage, quality, & toxicity. *ArXiv preprint*, abs/2305.13169.

Laurence T. Maloney, Maria F. Dal Martello, Vivian Fei, and Valerie Ma. 2024. A comparison of human and gpt-4 use of probabilistic phrases in a coordination game. *Scientific Reports*, 14(1).

Potsawee Manakul, Adian Liusie, and Mark Gales. 2023. SelfCheckGPT: Zero-resource black-box hallucination detection for generative large language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 9004–9017, Singapore. Association for Computational Linguistics.

Sabrina J. Mielke, Arthur Szlam, Emily Dinan, and Y-Lan Boureau. 2022. Reducing conversational agents' overconfidence through linguistic calibration. *Transactions of the Association for Computational Linguistics*, 10:857–872.

Donna L. Mohr, William J. Wilson, and Rudolf J. Freund. 2022. Chapter 1 - data and statistics. In Donna L. Mohr, William J. Wilson, and Rudolf J. Freund, editors, *Statistical Methods (Fourth Edition)*, fourth edition edition, pages 1–64. Academic Press.

Fabio Motoki, Valdemar Pinho Neto, and Victor Rodrigues. 2023. More human than human: measuring chatgpt political bias. *Public Choice*, 198(1–2):3–23.

Fabio Motoki, Valdemar Pinho Neto, and Victor Rodrigues. 2024. More human than human: Measuring chatgpt political bias. *Public Choice*, 198(1):3–23.

Keiichi Namikoshi, Alexandre L. S. Filipowicz, David A. Shamma, Rumen Iliev, Candice Hogan, and Nikos Aréchiga. 2024. Using llms to model the beliefs and preferences of targeted populations. *ArXiv preprint*, abs/2403.20252.

Joon Sung Park, Joseph O'Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S. Bernstein. 2023. Generative agents: Interactive simulacra of human behavior. In *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology*, UIST '23, New York, NY, USA. Association for Computing Machinery.

Anthony Patt and Suraje Dessai. 2005. Communicating uncertainty: lessons learned and suggestions for climate change assessment. *Comptes rendus. Géoscience*, 337(4):425–441.

Leonard Salewski, Stephan Alaniz, Isabel Rio-Torto, Eric Schulz, and Zeynep Akata. 2024. In-context impersonation reveals large language models' strengths and biases. *Advances in Neural Information Processing Systems*, 36.

Shibani Santurkar, Esin Durmus, Faisal Ladhak, Cinoo Lee, Percy Liang, and Tatsunori Hashimoto. 2023. Whose opinions do language models reflect? *Preprint*, arXiv:2303.17548.

Maarten Sap, Ronan Le Bras, Daniel Fried, and Yejin Choi. 2022. Neural theory-of-mind? on the limits of social intelligence in large lms. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 3762–3780.

Ruta Sawant and Sujit Sansgiry. 2018. Communicating risk of medication side-effects: role of communication format on risk perception. *Pharmacy Practice*, 16(2):1174.

Nino Scherrer, Claudia Shi, Amir Feder, and David Blei. 2024. Evaluating the moral beliefs encoded in llms. *Advances in Neural Information Processing Systems*, 36.

Nikil Selvam, Sunipa Dev, Daniel Khashabi, Tushar Khot, and Kai-Wei Chang. 2023. The tail wagging the dog: Dataset construction biases of social bias benchmarks. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 1373–1386, Toronto, Canada. Association for Computational Linguistics.

Preethi Seshadri, Pouya Pezeshkpour, and Sameer Singh. 2022. Quantifying social biases using templates is unreliable. *Preprint*, arXiv:2210.04337.

Yijia Shao, Yucheng Jiang, Theodore A. Kanell, Peter Xu, Omar Khattab, and Monica S. Lam. 2024. Assisting in writing wikipedia-like articles from scratch with large language models. *ArXiv preprint*, abs/2402.14207.

Vaishnavi Shrivastava, Percy Liang, and Ananya Kumar. 2023. Llamas know what gpts don't show: Surrogate models for confidence estimation. *ArXiv preprint*, abs/2311.08877.

Chenglei Si, Chen Zhao, Sewon Min, and Jordan Boyd-Graber. 2022. Re-examining calibration: The case of question answering. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 2814–2829, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Damien Sileo and Marie-francine Moens. 2023. Probing neural language models for understanding of words of estimative probability. In *Proceedings of the 12th Joint Conference on Lexical and Computational Semantics (*SEM 2023)*, pages 469–476, Toronto, Canada. Association for Computational Linguistics.

Aaditya K. Singh and DJ Strouse. 2024. Tokenization counts: the impact of tokenization on arithmetic in frontier llms. *Preprint*, arXiv:2402.14903.

Mark Steyvers, Heliodoro Tejeda, Aakriti Kumar, Catarina Belem, Sheer Karny, Xinyue Hu, Lukas Mayer, and Padhraic Smyth. 2024. The calibration gap between model and human confidence in large language models. *Preprint*, arXiv:2401.13835.

Winnie Street, John Oliver Siy, Geoff Keeling, Adrien Baranes, Benjamin Barnett, Michael McKibben, Tatenda Kanyere, Alison Lentz, Robin IM Dunbar, et al. 2024. Llms achieve adult human performance on higher-order theory of mind tasks. *ArXiv preprint*, abs/2405.18870.

Zhisheng Tang, Ke Shen, and Mayank Kejriwal. 2024. An evaluation of estimative uncertainty in large language models. *ArXiv preprint*, abs/2405.15185.

Karl Halvor Teigen. 2023. Dimensions of uncertainty communication: What is conveyed by verbal terms and numeric ranges. *Current Psychology*, 42(33):29122–29137.

Katherine Tian, Eric Mitchell, Allan Zhou, Archit Sharma, Rafael Rafailov, Huaxiu Yao, Chelsea Finn, and Christopher Manning. 2023. Just ask for calibration: Strategies for eliciting calibrated confidence scores from language models fine-tuned with human feedback. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 5433–5442, Singapore. Association for Computational Linguistics.

Gülden Ülkümen, Craig R Fox, and Bertram F Malle. 2016. Two dimensions of subjective uncertainty: Clues from natural language. *Journal of experimental psychology: General*, 145(10):1280.

Anne Marthe Van Der Bles, Sander Van Der Linden, Alexandra LJ Freeman, James Mitchell, Ana B Galvao, Lisa Zaval, and David J Spiegelhalter. 2019. Communicating uncertainty about facts, numbers and science. *Royal Society open science*, 6(5):181870.

Mudit Verma, Siddhant Bhambri, and Subbarao Kambhampati. 2024. Theory of mind abilities of large language models in human-robot interaction: An illusion? In *Companion of the 2024 ACM/IEEE International Conference on Human-Robot Interaction*, HRI '24, page 36–45, New York, NY, USA. Association for Computing Machinery.

Thomas S Wallsten, David V Budescu, Amnon Rapoport, Rami Zwick, and Barbara Forsyth. 1986a. Measuring the vague meanings of probability terms. *Journal of Experimental Psychology: General*, 115(4):348.

Thomas S. Wallsten, David V. Budescu, Rami Zwick, and Steven M. Kemp. 1993. Preferences and reasons for communicating probabilistic information in verbal or numerical terms. *Bulletin of the Psychonomic Society*, 31(2):135–138.

Thomas S Wallsten, Samuel Fillenbaum, and James A Cox. 1986b. Base rate effects on the interpretations of probability and frequency expressions. *Journal of Memory and Language*, 25(5):571–587.

Thomas S. Wallsten, Yaron Shlomi, and Hisuchi Ting. 2008. Exploring intelligence analysts' selection and interpretation of probability terms: Final report for research contract 'expressing probability in intelligence analysis'.

Yixin Wan, George Pu, Jiao Sun, Aparna Garimella, Kai-Wei Chang, and Nanyun Peng. 2023. "kelly is a warm person, Joseph is a role model": Gender biases in llm-generated reference letters. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 3730–3748.

Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H. Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023a. Self-consistency improves chain of thought reasoning in language models. In *The Eleventh International Conference on Learning Representations*.

Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh Hajishirzi. 2023b. Self-instruct: Aligning language models with self-generated instructions. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 13484–13508, Toronto, Canada. Association for Computational Linguistics.

Yizhong Wang, Swaroop Mishra, Pegah Alipoormolabashi, et al. 2022. Super-naturalinstructions: Generalization via declarative instructions on 1600+ nlp tasks. *Preprint*, arXiv:2204.07705.

Elke Weber and Denis Hilton. 1990. Contextual effects in the interpretations of probability words: Perceived base rate and severity of events. *Journal of Experimental Psychology: Human Perception and Performance,*, 16.

S. J. W. Willems, Casper Johannes Albers, and Ionica Smeets. 20. Variability in the interpretation of dutch probability phrases - a risk for miscommunication. *ArXiv preprint*, abs/.

Paul D Windschitl and Gary L Wells. 1996. Measuring psychological uncertainty: Verbal versus numeric methods. *Journal of Experimental Psychology: Applied*, 2(4):343.

Miao Xiong, Zhiyuan Hu, Xinyang Lu, Yifei Li, Jie Fu, Junxian He, and Bryan Hooi. 2024. Can llms express their uncertainty? an empirical evaluation of confidence elicitation in llms. In *Proceedings of the 12th International Conference on Learning Representations*, ICLR'24.

Gal Yona, Roee Aharoni, and Mor Geva. 2024. Can large language models faithfully express their intrinsic uncertainty in words? In *Main Conference of the Association for Computational Linguistics: EMNLP 2024 (to appear)*.

Yuqian Zhang, Anura De Zoysa, and Corinne Cortese. 2023. Foreign language effect in accounting uncertainty expressions: Interpretation and probabilistic estimation. *Journal of International Accounting, Auditing and Taxation*, 50:100528.

Zihao Zhao, Eric Wallace, Shi Feng, Dan Klein, and Sameer Singh. 2021. Calibrate before use: Improving few-shot performance of language models. In *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 12697–12706. PMLR.

Kaitlyn Zhou, Jena D. Hwang, Xiang Ren, and Maarten Sap. 2024. Relying on the unreliable: The impact of language models' reluctance to express uncertainty. *ArXiv preprint*, abs/2401.06730.

Pei Zhou, Aman Madaan, Srividya Pranavi Potharaju, Aditya Gupta, Kevin R McKee, Ari Holtzman, Jay Pujara, Xiang Ren, Swaroop Mishra, Aida Nematzadeh, et al. 2023. How far are large language models from agents with theory-of-mind? *ArXiv preprint*, abs/2310.03051.

# A Motivating example

In the main paper, we discuss the divergence of model behavior in the use case of news headline generation when prompted with text containing uncertainty expressions. In particular, Figure 1 shows two different levels of conviction in the generated text: the first headline contains confident-sounding words like "conclude" and "comprehensive" whereas the second generated headline uses less confident language like "suggests" and "possible". To further understand the extent to which these differences could be explained by the models' own knowledge, we prompt the model a second time to assess its agreement or disagreement with the topic being discussed. Its answers are in Figure 8.

# B Datasets

To ensure that our results generalize and are not artifacts of the evaluation methodology and/or benchmarks used (Selvam et al., 2023; Seshadri et al., 2022), we intentionally include statements spanning multiple categories and different syntactic structures. Moreover, to simplify and avoid potential length-based artifacts, we opt for simple, short statements that communicate a single fact.

## B.1 Non-Verifiable Statements

As mentioned in Section 3, we manually craft the non-verifiable statements to be short ordinary statements that purposely lack grounding on specific
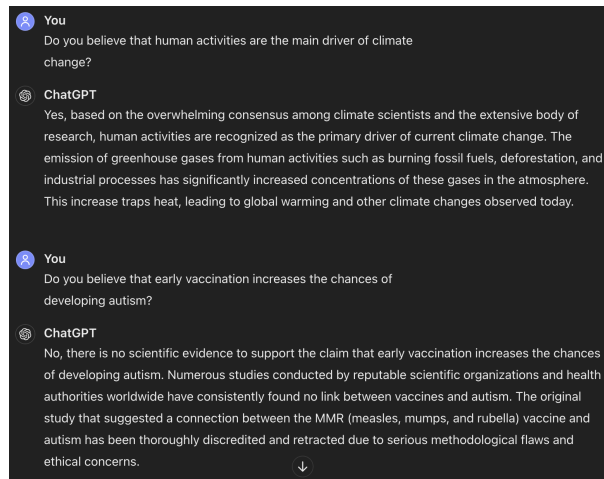
Figure 8: **Responses from `ChatGPT` when asked about its belief in the statements from Figure 1**. `ChatGPT` agrees that "human activities are the main driver of climate change," but disagrees with the statement that "early vaccination increases the chances of developing autism".

contextual information. In doing so, we hypothesize that an external observer will not have confidence in the truth or falsity of the statement. To ensure that our analysis covers a diverse set of plausible and realistic uses of uncertainty expressions, we create 15 statements for 4 different settings, including: (1) the forecasting of events; (2) the communication of imperfect knowledge; (3) to communicate possession; and (4) to communicate preferences. Below we list a random selection of 5 statements for each of the previous scenarios. Note that each of these statements are incorporated in the prompts listed in the main paper and the placeholders `[[they]]` and `[[their]]` are replaced by pronouns matching the gender of the statement speaker's name.

***Forecasting of future events***. Uncertainty expressions are often used to communicate uncertainty about future events.

- `[[they]]` will buy a new watch this Thanksgiving weekend.
- `[[they]]` will be offered a promotion this fall.
- the company will have another round of layoffs by mid July.
- there will be vegetarian options at the barbecue.
- `[[they]]` will visit New York over winter break.

***Imperfect knowledge***. Uncertainty expressions can also be used to communicate uncertainty imprecise information about events or outcomes.

- the restaurant near `[[their]]` apartment accepts reservations.
- the new museum is offering complimentary admission.
- there is a yoga studio within 2 miles of `[[their]]` workplace.
- there are more than eighty students in the auditorium right now.
- the temperature in the office is at least 72 degrees Fahrenheit.

***Possession***. Alternatively, uncertainty expressions can be used to convey uncertainty about the belongings of acquaintances.

- `[[their]]` boss owns a blue car.
- `[[their]]` friend has a leather jacket.
- `[[their]]` cousin has a vegetable garden.
- `[[their]]` classmate owns a guitar.
- `[[their]]` boss has a stereo amplifier.

***Preference***. Uncertainty expressions can be used to convey uncertainty about the preferences of acquaintances.

- `[[their]]` cousin prefers spinach over broccoli.
- `[[their]]` boss prefers coffee over tea.
- `[[their]]` friend prefers running over cycling.
- `[[their]]` neighbor prefers the beach over the mountains.

- `[[their]]` coworker prefers reading books over watching movies.

## B.2 Verifiable statements.

Unlike the non-verifiable statements, verifiable statements are created such that an external observer may have confidence in the truth or falsity of the statement. One way to ensure that the both humans and LLMs have high confidence in the truth or falsity of the statements is by focusing on popular general-knowledge facts, for example by focusing on trivia-like facts or student grade level science questions.

***TriviaQA dataset.*** For the experiments in the main paper, we propose to create short verifiable statements using the set of questions provided by The Question Company (as described in Section 4). In addition to covering multiple categories (*e.g.*, geography, entertainment, sports, astronomy), it also includes questions of varying difficulty (*i.e.*, easy to hard). Specifically, the "easy" subset of questions is designed to test knowledge about simple and straightforward common facts, without requiring specialized knowledge. For this reason, we choose to use the facts from this subset as a proxy for statements that an external observer is likely to know to be true or false. In particular, we use all 3 sets of 10 easy questions, spanning the topics "Cities/Geography," "Art/History," and "Science and Nature," and, for each question, we write one *true* statement using the correct answer choice and one *false* statement using the wrong choice (see examples in Table 4).

***AI2-ARC dataset.*** In addition to the 60 statements derived from the trivia dataset, we include an additional analysis using 400 verifiable statements (of which 200 are true and 200 are false). Using the same approach as previously described, we extract true and false verifiable statements from AI2-ARC (Clark et al., 2018), an elementary school-level multiple-choice question answering dataset focused on science facts. Despite its simplicity, AI2-ARC has been an important benchmark to measure progress in LLMs reasoning capabilities (Jiang et al., 2023; Achiam et al., 2024; Beeching et al., 2023), making it a relevant dataset to include in our analysis. Table 5 lists 6 true statements and 6 false statements derived from the AI2-arc datasets.

## B.3 Statement Correctness vs LLMs' beliefs

A key assumption underlying our selection of verifiable datasets is that, because our statements concern simple common facts, *LLMs must know whether the statements are true or false*. In this section, we conduct an empirical validation of this assumption by soliciting the LLMs' beliefs about each of the statements.

**Methodology.** Prompting is commonly used as a way to assess LLMs' knowledge. Given a question-answer pair, one common approach is to assume that a LLM ***knows*** the answer to $q$ is $a$ if it generates $a$ when prompted to answer $q$ (Gekhman et al., 2024; Kadavath et al., 2022; Manakul et al., 2023). Likewise, we can adopt this approach to determine an LLM's belief about the truth or falsity of a statement by asking the LLM about the veracity of a specific statement. To this end, we use the prompt:

"*Question: What is the veracity of the statement: "[[statement]]"?\nChoose from the following options: True, False, Unknown\nWrite only one of the answer choices and nothing else.*"

Using the prompt above, we empirically validate our assumption for four different LLMs: `Gemini`, `ChatGPT`, `GPT-4`, and `GPT-4o`. In particular, for every statement in non-verifiable and verifiable settings, we prompt each LLM to generate `n_samples=7` using `temperature=0.5` (Wang et al., 2023b). For each prompt, we compute the relative frequency of each of the possible responses 'true', 'false', and 'unknown' and use the mode response as the final prediction.

**Metrics.** Our goal is to provide supporting evidence that the models' belief of correctness differs for the three different settings evaluated in our experiments, in particular, we want to show that: (1) LLMs do not know the veracity of the non-verifiable statements; (2) LLMs knows which of the verifiable statements are true (dubbed "verifiable true") and which of them are false (dubbed "verifiable false"). To gauge model correctness, we report the average accuracy of the model. In particular, we consider a model to be correct (or accurate) if it generates "unknown" when prompted about the veracity of the non-verifiable statement; "true" when prompted about the veracity of a true verifiable statement; and "false" when prompted about the veracity of a false statement. To differentiate the accuracy of the models across the three different scenarios, we designate the models' accu-

Table 4: **Examples of true and false verifiable statements across different categories of the trivia dataset.** The true and false statements are created based on the correct and one of the incorrect choices of a trivia multiple-choice question answering dataset. The statements constitute short, simple, public knowledge facts.

| Category | True statement | False statement |
|---|---|---|
| Cities & Geography | "Great Britain directly borders *0* countries." <br> "New York is known as the **Big Apple**." <br> "the Colosseum, a famous landmark in Rome, was originally built as an **Amphitheatre**." | "Great Britain directly borders *2* countries." <br> New York is known as the **Big Orange**." <br> "the Colosseum, a famous landmark in Rome, was originally built as an **Cathedral**." |
| History & Art | "the Mona Lisa is a famous painting by **Leonardo da Vinci**." <br> "the Scream is the best known painting by **Edvard Munch**." <br> "**Andy Warhol** became a famous artist in the 1960s for painting soup cans and soap boxes." | "the Mona Lisa is a famous painting by **Tintoretto**." <br><br> "the Scream is the best known painting by **Jackson Pollock**." <br> "**Frida Kahlo** became a famous artist in the 1960s for painting soup cans and soap boxes." |
| Science & Nature | "**water**'s chemical formula is H2O." <br> "the nearest planet to the sun is **Mercury**." <br> "**oG** is a measure of the acidity or basicity of a substance." | "**carbon monoxide**'s chemical formula is H2O." <br> "the nearest planet to the sun is **Mars**." <br> "**pH** is a measure of the acidity or basicity of a substance." |

Table 5: **Examples of true and false verifiable statements derived from the AI2-Arc dataset** ([Clark et al., 2018](#)). The true and false statements are created based on the correct and one of the incorrect choices of the dataset. The statements constitute short school-level science facts.

| Category | True statement | False statement |
|---|---|---|
| Easy | "**Growing** and reproducing are two life processes that occur in both plants and humans." <br> "A light year refers to the **distance light travels in one year**." <br> "Carbon dioxide produced by cells is removed from the body primarily by the **respiratory system**." | "**Germinating** and reproducing are two life processes that occur in both plants and humans." <br> "A light year refers to the **time it takes light to travel from Earth to the Sun**." <br> "Carbon dioxide produced by cells is removed from the body primarily by the **immune system**." |
| Challenge | "**Trees** are a renewable natural resource that can be replenished over a period of time" <br> "When cold temperatures are produced in a chemical reaction, the reaction is known as **endothermic**." <br> "**Swimming fast** is an adaptive characteristic that helps dolphins survive life in the ocean." | "**Coal** is a renewable natural resource that can be replenished over a period of time." <br> "When cold temperatures are produced in a chemical reaction, the reaction is known as **exothermic**." <br> "**Traveling alone** is an adaptive characteristic that helps dolphins survive life in the ocean." |

racy in each setting as $\text{Acc}_{NV}$, $\text{Acc}_{VT}$, and $\text{Acc}_{VF}$, respectively.

**Results.** Table 6 shows the average accuracy results of applying the methodology described previously to four different LLMs. Overall, our results show that *LLMs are able to differentiate between verifiable and non-verifiable settings*. Specifically, GPT-4 and GPT-4o abstain from assigning a veracity judgment to the non-verifiable statements in about >99% of the examples, while correctly judging 100% of the true verifiable statements and >93% of the false verifiable statements.

Overall, we find that *models are able to correctly discriminate between true and false statements (>90% accuracy)* which corroborates the use of true correctness of the statements as a proxy for the LLMs' beliefs.

## C  Experiment Details

This section describes in greater detail various aspects of the experiments conducted in this paper, including the list of uncertainty expressions, the name selection strategy, the list of prompts, a list of statements, as well as additional details on the human experiments.

### C.1  Uncertainty Expressions

The uncertainty expressions are a subset of the expressions proposed in Wallsten et al.; Wallsten et al.; Willems et al.; Fore. The final list of uncertainty expressions used in this paper is listed below:

- almost certain, highly likely, very likely, likely, probable, somewhat likely, somewhat unlikely, uncertain, possible, unlikely, not likely, doubtful, very unlikely, highly unlikely

|                  | Avg Acc (%) | $\text{Acc}_{NV}$ (%) | $\text{Acc}_{VT}$ (%) | $\text{Acc}_{VF}$ (%) |
|------------------|-------------|-----------------------|-----------------------|-----------------------|
| Num examples     | 170         | 110                   | 30                    | 30                    |
| ChatGPT          | 90.00       | 90.00                 | 90.00                 | 90.00                 |
| GPT-4            | 98.82       | 100.00                | 100.00                | 93.33                 |
| GPT-4o           | 98.82       | 99.09                 | 100.00                | 96.67                 |
| Gemini           | 90.00       | 90.91                 | 93.33                 | 83.33                 |

Table 6: **Accuracy of LLMs' beliefs about the veracity of the non-verifiable (NV) and verifiable statements (V) in the main experiment.** Overall, models are very accurate in differentiating between non-verifiable and verifiable statements. Moreover, we observe that models are able are able to correctly discriminate between true and false statements (>90% accuracy), which corroborate the hypothesis that LLMs know the correctness of the statements they are tested with.

## C.2 Name Selection

All names used in our experiments were collected from a random name generator[11], which we ran iteratively until we obtained 32 unique names, half of each biological gender (as determined by the random generator).

- **Female names**: "Amanda", "Bonnie", "Camille", "Catherine", "Cheri", "Ethel", "Gabriela", "Jacquelyn", "Jessica", "Laura", "Olga", "Roxanne", "Silvia", "Tara", "Violet"

- **Male names**: "Brendan", "Bruce", "David", "Gary", "Isaac", "Jeffery", "Joey", "Johnnie", "Kenny", "Lance", "Marco", "Mike", "Nathan", "Nick", "Raul"

## C.3 Human Experiments

Human responses were collected using Prolific (https://www.prolific.com/). We recruited 100 participants for the non-verifiable experiment and 100 different participants for the second verifiable experiment. One of the 100 responses was not received due to a technical issue in both the first and second experiment, leaving a total of 99 responses for each. We recruited participants whose first language was English that were located in the United States. Participants were paid $2 for completing the study and the average completion time was 8 minutes and 48 sections; the average payment rate was $13.64/hour. The University of California, Irvine Institutional Review Board (IRB) approved the experimental protocol. Prior to the experiment, participants were given detailed instructions outlining the experimental procedure as well as how to understand and interact with the user interface.
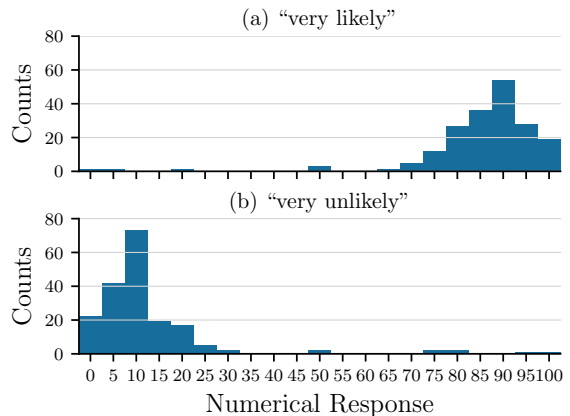


Figure 9: **Histogram of human participant responses for non-verifiable statements for two uncertainty expressions**. The modes of the two distributions sum to 100%, suggesting symmetry in the population-level interpretation of the two expressions.

Participants were asked to sign an integrity pledge after reading all of the instructions, stating that they would complete the experiment to the best of their abilities. After submitting their integrity pledge, participants were granted access to the experiment.

We filtered out low-quality responses with the following procedure. For each participant, we computed the Spearman correlation between the participant's responses and the overall ranking of uncertainty statements in the non-verifiable experiment. We removed participants with $\rho < 0.2$, a threshold chosen empirically to filter out only no-signal, spam-like responses. This filter removed 5 participants in the first experiment and 10 in the second experiment. In total, we remain with 94 participants in the non-verifiable experiment and 89 in the verifiable experiment.

---

[11] https://randomwordgenerator.com/name.php, last accessed on March 26th, 2024.

## C.4 Prompts

In our main paper, we conduct experiments using 2 demonstrations, since this better replicates the human setup. We use the prompt listed in Figure 10 to elicit the model's numerical response, and use the format described in Figure 11 to represent the demonstrations in the prompt.

During the course of our experiments, we carried experiments with varying assumptions: non-verifiable setup assessed models (and humans) perceptions in the absence of strong prior knowledge about the statements, whereas the verifiable setup focused on the evaluation of the same perceptions when knowledge was present. We used two different sets of exemplars (or demonstrations) in our experiments to reflect these differences. The set of exemplars used in the non-verifiable setting are defined in terms of the following:

- speaker: "Kathleen", uncertainty: "impossible", statement: "the cafe made a profit in the last 6 months".

- speaker: "Cedric", uncertainty: "certain", statement: "the new treatment will improve the patient's condition".

Conversely, for the set of exemplars used in the **verifiable** setting, we carried a preliminary experiment to determine whether the truthness/falsehood of the examples and their ordering had a significant impact in models' performance. We used the two pairs of statements ("only some metals can conduct electricity," "all metals can conduct electricity") and ("the Sun orbits around the planet Earth," "the planet Earth orbits around the Sun") in the preliminary experiments. We found negligible differences in the obtained distribution of numerical responses (as emphasized in Tables 7 and 8). As a result, we report the results using a false statement as the first example and a true statement as the second example in the prompt (denoted FT). These were associated with the following speaker names and uncertainty expressions:

- speaker: "Kathleen", uncertainty: "impossible", statement: "the Sun orbits around the planet Earth".

- speaker: "Cedric", uncertainty: "certain", statement: "all metals can conduct electricity".

Table 7: **Average distributional difference in LLMs conditional distributions when estimated based on different true/false configurations of the examples in the verifiable prompt**. We report the Wasserstein distance of the estimated conditional distributions averaged over the empirical distributions obtained using the other configurations and greedy decoding algorithm (see Table 8 for results using the probabilistic decoding). On average, we observe that changing the veracity of the examples in the prompt has minimal impact on the final distributions. Moreover, we explicitly avoid using the configurations FF or TT to avoid majority label bias (Zhao et al., 2021).

|     | ChatGPT | GPT-4 | GPT-4o |
| --- | ------- | ----- | ------ |
| FF  | 2.33    | 1.31  | 0.92   |
| FT  | 3.20    | 1.25  | 0.85   |
| TF  | 2.93    | 1.32  | 0.88   |
| TT  | 2.31    | 1.38  | 0.90   |

Table 8: **Average distributional difference in LLMs conditional distributions when estimated based on different true/false configurations of the examples in the verifiable prompt**. On average, we observe that changing the veracity of the examples has minimal impact on the final distributions.

|     | ChatGPT | GPT-4 | GPT-4o |
| --- | ------- | ----- | ------ |
| FF  | 1.91    | 1.53  | 1.74   |
| FT  | 3.01    | 1.36  | 1.36   |
| TF  | 2.77    | 1.33  | 1.81   |
| TT  | 1.91    | 1.48  | 1.39   |

## C.5 Language Models

Our evaluation concerns the study of instruction-tuned LLMs, some of which are accessible through black-box APIs and others through the use of the HuggingFace Python package. We use OpenAI to obtain the results for ChatGPT, GPT-4, and GPT-4o; Google's Vertex AI API to obtain results for Gemini, and TogetherAI[12] to run LLama3 (70B), Mixtral 8x7B, Mixtral 8x22B, and DBRX. We run LLama3 (8B) OLMo (7B) and Gemma (2B) locally on a single GPU 8 RTX A6000 (48 GB). All experiments were conducted from April through June 2024.

For simplicity, we have abbreviated the name of the evaluated models, removing information about the size and version of the model. For reproducibility, we list below the mapping from model name to

---
[12]https://www.together.ai/

In this experiment, you will be shown text reflecting an individual's beliefs about a statement. You will then be asked to judge, in your opinion, the strength of that individual's belief. To do so, you will provide a probability (expressed as a percentage) where:
- The number 0 indicates that the individual believes with certainty that the statement is false.
- The number 100 indicates that the individual believes with certainty that the statement is true.

You will be shown two example question and response pairs below to familiarize you with the experiment setup. After the examples, you will be asked to answer a new question.

## Example Question 1
[[demonstration1]]
Choose the answer from the following options:
0, 5, 10, 15, 20, 25, 30, 35, 40, 45, 50, 55, 60, 65, 70, 75, 80, 85, 90, 95, 100
Correct answer: 100

## Example Question 2
[[demonstration2]]
Choose the answer from the following options:
0, 5, 10, 15, 20, 25, 30, 35, 40, 45, 50, 55, 60, 65, 70, 75, 80, 85, 90, 95, 100
Correct answer: 100

## Question
Given the examples before, answer the following question by writing a single number as the answer.

[[test-example]]
Choose the answer from the following options:
0, 5, 10, 15, 20, 25, 30, 35, 40, 45, 50, 55, 60, 65, 70, 75, 80, 85, 90, 95, 100 Correct answer:

Figure 10: **Prompt used to conduct the experiments in the 2-shot setting.** [[demonstration1]], [[demonstration2]], and [[text-example]] are placeholders that will be composed of a statement (or event), question, and individual instructions. The prompt closely reflects the interface presented to the participants of the user studies. Similarly to the human interface, we encourage LLMs to select a multiple of 5 in the range [0, 100].

> Sentence: "[[speaker]] believes it's [[uncertainty]] that [[statement]]."
>
> Question: Provide a number between 0 and 100 that reflects, from [[speaker]]'s perspective, the probability (expressed as a percentage) that [[statement]].
> - 0 indicates [[speaker]] believes there is a 0% chance that [[statement]].
> - 100 indicates [[speaker]] believes there is a 100% chance that [[statement]].

Figure 11: **Structure of a demonstration (or exemplar), as defined in the 2-shot prompt**. The placeholders [[speaker]], [[uncertainty]], and [[statement]] are replaced by gendered names, uncertainty expressions, and statements, respectively.

exact version of the model used:

- ChatGPT: gpt-3.5-turbo-0125
- GPT-4: gpt-4-turbo-2024-04-09
- GPT-4o: gpt-4o-2024-05-13
- LLama3 (8B): meta-llama/Meta-Llama-3-8B-Instruct
- LLama3 (70B): meta-llama/Llama-3-70b-chat-hf
- Gemini: models/gemini-pro
- Mixtral 8x7B: mistralai/Mixtral-8x7B-Instruct-v0.1
- Mixtral 8x22B: mistralai/Mixtral-8x22B-Instruct-v0.1
- Gemma (2B): google/gemma-1.1-2b-it (we found Gemma (2B) to respond better empirically to the prompts than its 7B version, which tended to extrapolate the few-shot instructions with additional examples).
- OLMo (7B): allenai/OLMo-7B-Instruct

# D Extracting LLMs Numerical Responses

In this section, we outline the methodologies used to extract the numerical responses from auto-regressive LLMs. In an ideal world, given an uncertainty expression $u$, the model's conditional distribution $\hat{P}_{\text{model}}(k|u)$, would be fully observable. However, information about $\hat{P}_{\text{model}}$ is seldom available, for models are either (1) served through opaque closed-source APIs or (2) too large, often requiring large amounts of compute to estimate the full distribution. To circumvent such limitations, one idea is to empirically estimate such probabilities using sampling-based approaches (*e.g.*, self-consistency (Wang et al., 2023a)) or using greedy decoded sequences over multiple examples. In the ensuing sections, we describe the three different strategies considered in our work to estimate the LLMs' empirical distributions.

## D.1 Method 1: Full Probability Approach

The "full" methodology requires access to the next-token probability distribution of an auto-regressive LLM and, as a result, is currently applicable to open source models. Because we use the models' probability to compute the probability of producing any integer between 0 and 100, it is also our most time-consuming approach[13], requiring 101 model calls to obtain the *full* probability.

**Methodology.** An important aspect to account for when using LLMs to estimate the probability distribution over the set of integers ranging from 0 to 100 is the LLM's tokenizer. Specifically, models, such as Gemma (2B), LLama3 (70B), and OLMo (7B), use single-digit tokenization (Singh and Strouse, 2024), which implies that the textual representation of 100 (represented as "100") is tokenized into at least three individual tokens (*i.e.*, "1", "0", "0") and not in a single token (*e.g.*, "100"). In terms of probabilities, this is a challenge, since many of these LLMs generate digits in a left-to-right fashion and, by design, the probabilities of any number in $[10, 100]$ are always less than or equal to the probability of the first constituting digit (*e.g.*, probability of "10" subsumes the probability of "1" since in order to generate the string "10" the model needs to generate the string "1" first). Since we are interested in knowing the true independent probability assigned to every integer between 0 and 100, we need to adjust the LLMs' probability. Therefore, instead of reporting the probability that a number $i \in [0, 100]$ occurs, we compute the probability that $i$ occurs and is not followed by a number $j \in [0, 9]$:

---

[13]Running LLama3 (70B) across 5x GPU 8 RTX A6000 (48 GB) for 900 examples required approximately 9 full days. For that reason, we only apply this strategy to a few LLMs.

$$p(x_t = i | x_{<t}) - \sum_{j=0}^{9} p(x_t = i, x_{t+1} = j | x_{<t}),$$

where $p$ represents the LLM's next-token probability distribution. In sum, given a prompt parameterized with a speaker $s$, uncertainty expression $u$, statement $e$, and a prompt that combines these parameters, denoted prompt$(s, u, e)$ we use the expression above to obtain the full probability distribution for all $k \in 0, ..., 100$. We denote this corrected probability distribution as $p(k|\text{prompt}(s, u, e))$ and we refer to the probability mass that falls outside the set of strings $k \in 0, ..., 100$ as $\bar{p}$.

***Constructing the Greedy Histogram.*** Unlike sampling-based greedy decoding algorithms, we restrict the selection of the arg-max to the set of strings representing the numbers between $[0, 100]$, as opposed to sampling a series of tokens using temperature=0. In other words, we constrain the greedy decoding to be any of the sequences in {"0", "1", ..., "100"}, regardless of how little probability[14] is assigned to any of these numerical sequences. For every triplet $(u, s, e)$, we obtain the the arg-max prediction and then assign it to a bin $0, 5, 10, ..., 100$.

***Constructing the Probabilistic Histogram.*** The probabilistic histogram benefits from the probability information that was computed previously for a specific triplet $(u, s, e)$. In particular, for a specific example we accumulate all probability values in the corresponding bins $0, 5, ..., 100$. The remaining probability mass $\bar{p}$ is assigned to a default bin "-1". In other words, the bin "-1" will accumulate the probability of a number in $[0, 100]$ not following the specified prompt. While we could have normalized the probabilities of $0, 5, ..., 100$ to sum to 1, we decided to add a "-1" bin to allow for a fair comparison with the top-k approach, where only part of the probability distribution (*e.g.*, top-20 next probability distribution values are revealed), as is the case for the OpenAI models. After accumulating the probability over all $(u, s, e)$ triplets, we normalize the histogram by dividing by the number of triplets.

---

[14]We found this approach to be particular sensitive to the instruction format used. For example, in earlier iterations of this work we used Llama-2 and Gemma (7B) but found them to be particular sensitive to the whitespaces provides in the prompt.

## D.2 Method 2: Top-k approach

This method leverages information about the top-k values of the next-token probability information. Applicable to OpenAI models, namely ChatGPT, GPT-4, and GPT-4o, this approach allow us to obtain richer probability information with a single API call.

***Methodology.*** This method requires two properties to be satisfied: (1) numbers between 0 and 100 must be encoded with one single token each (*i.e.*, each integer is represented with a single token), and (2) exponentiating the log probabilities returned by the API must lead to a valid probability distribution (*i.e.*, numbers obtained for different prompts will be comparable to one another). We validate that the first requirement is satisfied by OpenAI models.
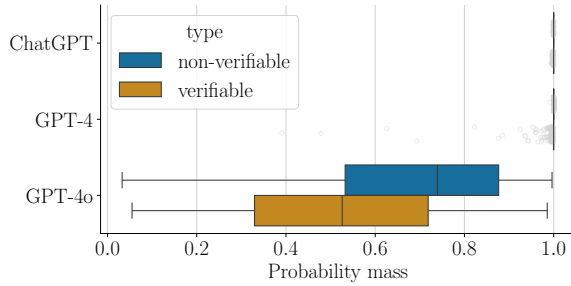
***Constructing the Greedy Histogram.*** Unlike traditional greedy decoding, we condition the selection of the arg-max prediction over numbers the top-k (k=20 for OpenAI). That is, we select the most likely number that is present in the top-20 predicted tokens. If no number is present in the top-20 tokens, we assign a default value of '-1'.

***Constructing the Probabilistic Histogram.*** Like the "full" probability approach, we construct the probabilistic histogram by accumulating the probabilistic information that we gather with each inference call: For a given triplet $(u, s, e)$, we make an API call and obtain probability about the next 20 tokens. If any of these strings represents a number between 0 and 100, we exponentiate it to obtain a probability value, and accumulate the probability in the corresponding bin. Any remaining probability that is not assigned to a number in the top-20 is accumulated in the '-1' bin. After accumulating the probability over all $(u, s, e)$ triplets, we normalize the histogram by dividing by the number of triplets.
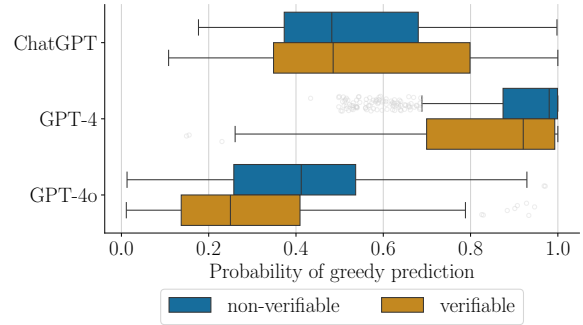
## D.3 Method 3: Sampling-based approach

Some models are provided through black-box APIs with no access to next-token probability information. Ideally, one could estimate the next-token probability information by continuously sample, but doing so would be too costly. In the main paper, we only run sampling based approaches using a single sample.

***Constructing the Greedy Histogram.*** This approach differs from the previous two in that it is unconstrained sampling-based approach, which means that the next immediate token may not be a number. Moreover, it is also agnostic to the tok-

(a) Total probability mass assigned to a number across models and settings.



(b) Probability mass of greedy prediction.

Figure 12: **Differences in the probability mass as determined by OpenAI models on the top-20 tokens**. We report these values across all statements (n=840) in the verifiable and non-verifiable settings. We observe that numbers account for the majority of probability mass in `ChatGPT` and `GPT-4`. Upon analysis we found lower probability mass assigned by `GPT-4` to be correlated with the appearance of the words "Given," "The," and "And".

enization. It many cases, the model may generate some text before it actually produces the answer "The answer is 100" or "**Sure, here is the answer: 55". To accommodate for such scenarios, we consider the arg-max prediction to be the first number between 0 and 100 to be mentioned in the model's response. Otherwise we assume the arg-max prediction to be '-1'.

***Constructing the Probabilistic Histogram.*** Due to budget constraints, we could not reliably estimate a probabilistic histogram via the sampling-based approach, as this would require thousands of requests. Future work could use strategies like self-consistency to estimate a probabilistic histogram.

### D.4 Which methodology applied to which model?

The following describes the list of methodologies used for each LLM:

- Greedy Sampling: `Gemini`, `LLama3 (70B)`, `Mixtral 8x7B`, `Mixtral 8x22B`. We opt for using greedy sampling (`max_tokens=200`, `temperature=0`) as opposed to standard sampling due to budget constraints. Given the nature of our experiments, faithfully estimating the empirical distributions over the 101 numbers would require hundreds or even thousands of calls. These calls are time-consuming and costly. We believe that using greedy decoding is still representative of how a model would behave in the majority of the cases.

- Full next-token probability distribution: `LLama3 (8B)`, `Gemma (2B)`, `OLMo (7B)`. We found these models to be particularly brittle to prompting.

- Next-token probability distribution: `ChatGPT`, `GPT-4`, `GPT-4o`. As of June 2024, OpenAI models only provide access to the next-token probabilities of the top-20 tokens. In the experiments in Section 5.4, we collect the information about the top 20 numbers.

## E Additional Results

In this section, we report the following additional results:

- In Section E.1 we report the mean and 95% confidence estimates for the proportional agreement (PA) and mean average error (MAE) metrics. Since the PA prioritizes mode-matching behavior, we also include a distribution-matching metric based on Wasserstein distance.

- In Section E.2 we analyze differences in model responses between true and false statements on a pairwise basis, pairing statements by the original question they were based on.

- In Section E.3 we comment on the variability of the distributions obtained when using probabilistic decoding vs simple greedy decoding.

- In Section E.4 we provide visualizations of the histogram distributions for the non-verifiable and verifiable settings.

### E.1 Mode- and Distribution-Matching Metrics

Tables 9 and 10 show the mean and 95% confidence intervals for the PA and MAE metrics. For completeness, we also include the PA score and

MAE metrics per expression in Sections E.1.1 and E.1.2.

### E.1.1 Proportional Agreement

As described in Section 4.3, the proportional agreement (PA) metric gauges the overall agreement between an agent's and a reference (population) distribution. We use the results of the human studies in the non-verifiable setting as our reference distribution (dubbed Human Mode) throughout the whole paper.

Tables 12 and 13 report the proportional agreement (PA) metric discriminated by uncertainty expression in the non-verifiable and verifiable settings, respectively. The values are reported with respect to the population-level human distribution described in Section 3, denoted Human+NV.

### E.1.2 Mean Absolute Error

Proposed in Section 4.3, Mean Absolute Error (MAE) measures the average agreement across uncertainty expressions between an agent's distribution and a reference (population) distribution. In this case, we also use the human results from the non-verifiable setting as our reference distribution throughout the paper.

### E.1.3 Wasserstein Distance

We use one-dimensional Wasserstein distance as a measure of the similarity between two conditional distributions. We use `scipy.stats.wasserstein_distance` and provide the 22 bins for each conditional distribution (*i.e.*, all 21 bins from $0, 5, ...100$ but also the -1) and the corresponding normalized counts that were used to estimate the conditional distributions. Two identical distributions will be assigned a Wasserstein-1 distance of 0, whereas two maximally distant distributions will be assigned a Wasserstein-1 distance of 101.

Table 11 summarizes the distributional differences between the LLM conditional distributions for both LLM+NV and LLM+V settings with respect to the reference distribution (estimated in Section 3 and denoted Human+NV). We observe that overall the obtained results correlate well with the results reported by the MAE metric.

Other uses of distributional differences in the verifiable setting: we compare how different the empirical distributions obtained in the main experiment differ with respect to the empirical distributions observed in the generalization experiment. These results are reported in Section F.

### E.1.4 Differences in Mean Numerical Responses

Figures 13 and 16 show the mean numerical response for the two verifiable datasets, discriminated by model and uncertainty expression. Overall, both plots show evidence of the large perceptual differences exhibited by different models according to the truthfulness of the evaluated statements.

### E.2 Pairwise Correctness Analysis

In our paper, we find that LLMs exhibit a systematic "knowledge bias," where the average numerical response is higher for true statements than for false statements. To provide statistical support to this claim, we run a one-sided paired Wilcoxon signed-rank test, comparing mean response between true and false statements. We perform a paired test, as a true and false statement were generated from each original question (e.g., "water's chemical formula is H2O" and "carbon monoxide's chemical formula is H2O" are paired). Overall, we find that the knowledge bias is statistically significant at $\alpha = 0.05$ for all models except OLMo (7B) and Gemma (2B)—see Tables 14, 15, and 16.

### E.3 Variability

In this section, we conduct a quantitative analysis of the variability of the LLMs' empirical distributions. This analysis is motivated by the observation that, visually, the empirical distributions obtained via greedy decoding appear to be less diverse than the population-level human distributions. In addition to drawing comparisons between humans and LLMs, we also inspect how the observed variability changes with different decoding algorithms, namely, we consider a probabilistic decoding algorithm `temperature=1`.

*Metric.* As a measure of the variability of an empirical distribution, we consider the *interquartile range (IQR)* (Mohr et al., 2022), *i.e.*, the difference between the 0.75 and 0.25 quantiles.[15] Intuitively, a smaller IQR value implies that the models' predictions tend to be concentrated in the same set of values, whereas a larger IQR value suggests a more uniform distribution. Because we have an empirical distribution for each uncertainty expression, we use the average IQR across uncertainty expressions to analyse the overall behavior.

---

[15]Other metrics could also be used to study the variability of the distributions, such as the entropy. However, we leave further analysis for future work.

Table 9: **Human-LLM agreement for non-verifiable statements**. Average Proportional Agreement (PA), PA as a fraction of the *Human Mode* results (% PA), and absolute error between mean responses (MAE). *Human Mode* represents the mode of the human NV distribution, whereas *Human Individual* represents the PA score of individual human responses relative to the population. 95% confidence intervals were computed using the adjusted bootstrap percentile method.

| | PA (95% CI) | % PA | MAE (95% CI) |
|---|---|---|---|
| Human Mode | 27.6 | — | — |
| Human Individual | $17.6_{(17.2,18.0)}$ | 63.8 | $8.91_{(8.08,9.17)}$ |
| ChatGPT | $19.7_{(19.1,20.1)}$ | 71.4 | $6.80_{(6.18,7.33)}$ |
| GPT-4 | $24.4_{(24.2,24.6)}$ | 88.4 | $4.64_{(4.53,4.74)}$ |
| GPT-4o | $18.9_{(18.4,19.5)}$ | 68.5 | $5.58_{(5.34,5.79)}$ |
| Gemini | $25.4_{(25.1,25.6)}$ | 92.0 | $4.09_{(3.92,4.23)}$ |
| LLama3 (8B) | $17.8_{(17.4,18.2)}$ | 64.5 | $11.99_{(10.92,13.13)}$ |
| LLama3 (70B) | $23.6_{(23.4,23.8)}$ | 85.5 | $5.56_{(5.34,5.80)}$ |
| Mixtral 8x7B | $21.8_{(21.4,22.0)}$ | 79.0 | $5.88_{(5.44,6.34)}$ |
| Mixtral 8x22B | $21.8_{(21.6,22.0)}$ | 79.0 | $7.20_{(6.56,7.81)}$ |
| OLMo (7B) | $11.1_{(10.7,11.5)}$ | 40.2 | $21.41_{(20.17,22.41)}$ |
| Gemma (2B) | $8.1_{(7.7,8.5)}$ | 29.3 | $20.17_{(19.11,21.17)}$ |

***Results.*** In the non-verifiable setting, we find that most models exhibit lower average IQR values than population-level human perceptions (see Table 17). In particular, the values reported for GPT-4, GPT-4o, and LLama3 (70B) are up to 13x smaller than the one reported for humans. These findings suggest that, when using greedy decoding, these LLMs are unable to replicate the diversity of human responses.

When considering the change in avg IQR values for probabilistic decoding in Table 17, we observe an increase in the spread for most models with respect to the greedy decoding values (+6.60 and +10.60 IQR increase on average for 2-shot and 0-shot, respectively). In fact, from the evaluated models, GPT-4 and LLama3 (70B) remain mostly insensitive to the change of decoding algorithm (between 0 and 1.20) whereas the IQR values for ChatGPT and GPT-4o increase +4.38 and +35.3, respectively. These results seem to suggest that for some models, the variability of the observed empirical distributions may be a function of the employed decoding algorithm.

### E.4 Histograms

Figure 14 depicts the empirical distributions for the non-verifiable experiments.

## F Generalization results

Diversity of grammatical and semantic structures is an important component of current evaluation practices in LLMs (Selvam et al., 2023; Seshadri et al., 2022), since it helps ensure that obtained results are not an artifact of the evaluation methodology and/or benchmarks used. The experiments described in the main paper were carefully crafted to cover various topics and situations where uncertainty expressions could be used. To further strengthen our analysis and validate our findings, we simultaneously run collect models perceptions of uncertainty expressions using a larger dataset. This dataset by the authors based on the AI2-ARC test set (Clark et al., 2018) — a popular question-answering dataset consisting of genuine grade-school level, multiple-choice science questions. Not only has this dataset been recently used to measure commonsense reasoning of current state-of-the-art LLMs (Jiang et al., 2023; Achiam et al., 2024; Beeching et al., 2023), but it is also composed of easier questions, a key aspect to our verifiable experiment setup.

The creation of this dataset mirrors the procedure described in Section 4. We manually repurposed 200 question-answer pairs from AI2-ARC (100 from the easy set and another 100 from the challenge set). For every statement, the authors produce a true statement and a false statement using the available information about the correct and incorrect multiple choices. The final dataset consists

Table 10: **Human-LLM agreement for verifiable statements**. Average Proportional Agreement (PA), absolute error between mean responses (MAE), and the difference between these scores and those from the non-verifiable statements (Table 1) (△ PA and △ MAE, respectively). Again *Human Mode* represents the mode of the human NV distribution, whereas *Human Individual* represents the average behavior across individual humans on the verifiable setting.

| | PA *(95% CI)* | △ PA | MAE *(95% CI)* | △ MAE |
|---|---|---|---|---|
| Human Mode | 27.6 | — | — | — |
| Human Individual | $16.7_{(16.3,17.1)}$ | -0.9 | $9.35_{(8.23,9.50)}$ | 0.44 |
| ChatGPT | $15.3_{(14.6,15.9)}$ | -4.4 | $8.57_{(6.81,10.07)}$ | 1.77 |
| GPT-4 | $22.1_{(21.7,22.4)}$ | -2.3 | $3.84_{(3.03,4.45)}$ | -0.80 |
| GPT-4o | $15.2_{(14.5,15.9)}$ | -3.7 | $7.05_{(6.45,7.62)}$ | 1.47 |
| Gemini | $21.3_{(20.7,21.8)}$ | -4.1 | $7.23_{(6.04,8.49)}$ | 3.14 |
| LLama3 (8B) | $10.1_{(9.5,10.7)}$ | -7.7 | $16.59_{(15.10,18.18)}$ | 4.60 |
| LLama3 (70B) | $18.9_{(18.1,19.6)}$ | -4.7 | $13.73_{(11.92,15.64)}$ | 8.17 |
| Mixtral 8x7B | $15.2_{(14.5,15.9)}$ | -6.6 | $12.23_{(10.37,14.18)}$ | 6.35 |
| Mixtral 8x22B | $18.6_{(18.3,19.0)}$ | -3.2 | $9.78_{(8.51,11.08)}$ | 2.58 |
| OLMo (7B) | $7.6_{(7.2,8.0)}$ | -3.5 | $33.66_{(31.82,35.11)}$ | 12.25 |
| Gemma (2B) | $5.3_{(5.0,5.6)}$ | -2.8 | $25.07_{(23.45,26.66)}$ | 4.9 |

of 200 true statements and 200 false statements.

To determine distributional differences between the conditional distributions obtained in the main paper and the ones obtained in the generalization set, we compare the Wasserstein-1 distance of the two empirical distributions. These values are reported in Table 18. In general, we find models that performed worse in the main paper, including Gemma (2B) and OLMo (7B), to exhibit the largest distributional differences with Wasserstein distances of 48.5 and 13.1 when averaged over uncertainty expressions. ChatGPT, LLama3 (70B), and Mixtral models all exhibit higher differences in expressions of higher certainty, e.g., "highly likely", "probable", "possible". On the other hand, the two GPT-4 models, as well as Gemini (Pro) suffer the least changes distributionally (1.9, 1.8, and 4.2 Wasserstein-1 distances on average, respectively), suggesting that these models were robust to changes in the statements.

In the main paper, we find it surprising that LLMs perception abilities differ significantly based on whether the uncertainty expressions are referring to someone's belief in a true or false statement. To test the generalization of this finding in a larger (and different) dataset, we repeat the same analysis and compare the observed mean response differences with that of humans obtained in the original setting (see Figures 15 and 16). We observe that in absolute sense the differences are smaller than

those observed in the original setting, but that models are affected by this knowledge gap to a greater extent than humans.

## G  Probabilistic Decoding

The findings reported in the main paper mostly concern the conditional distributions that are estimated using greedy decoding algorithm. Despite being a common decision in analyses papers (Yona et al., 2024; Steyvers et al., 2024), the use of greedy decoding may not provide the full picture of model behavior (Ivgi et al., 2024; Holtzman et al., 2020). Thus, to ensure that our results are not a degenerate behavior only observed when using greedy decoding algorithms, we conduct an analysis considering a probabilistic decoding (temperature=1): instead of estimating the conditional distributions for each uncertainty expression using the arg-max predicted numerical response, we use the available probability information to estimate the conditional distributions.

Most of our analyses on the study of OpenAI models— ChatGPT, GPT-4, and GPT-4o—using the top-k approach. Using OpenAI models has the benefit of being more cost-efficient than the other approaches, *i.e.*, it is less time-consuming and lower cost than the full probability methodology or the sampling-based approach. Moreover, even though theoretically the use of full probability information is better for estimating the conditional distributions,

Table 11: **Summary metrics averaged across uncertainty expressions for both NV and V settings**. All metrics are computed with respect to the human distribution in the non-verifiable setting (Human+NV). "PA" reports the general agreement between LLMs and the mode of the human distribution, reported in percentages. "MAE" reports the absolute error between the mean responses of LLMs and those of humans. Wasserstein-1 computes the distance between LLMs and human distributions.

| | | Avg PA (↑) | | Avg MAE (↓) | | Avg Wasserstein-1 (↓) | |
|---|---|---|---|---|---|---|---|
| | | NV | V | NV | V | NV | V |
| Human | Mode | 27.6 | 27.6 | — | — | — | — |
| | Individual | 17.6 | 16.7 | 8.91 | 9.35 | 12.35 | 12.99 |
| Baseline | Random | 5.1 | 5.1 | 27.72 | 27.72 | 28.16 | 28.16 |
| LLM | OLMo | 12.1 | 7.6 | 18.44 | 33.67 | 20.45 | 40.16 |
| | Gemma (2B) | 8.1 | 6.6 | 20.17 | 24.33 | 22.13 | 25.89 |
| | Llama3 8B | 17.8 | 10.1 | 11.99 | 16.59 | 14.11 | 18.35 |
| | Llama3 (70B) | 23.6 | 18.8 | 5.56 | 13.73 | 9.94 | 16.39 |
| | Mixtral 8x7B | 21.8 | 15.2 | 5.88 | 12.32 | 8.88 | 15.93 |
| | Mixtral 8x22B | 21.8 | 18.6 | 7.20 | 9.78 | 10.78 | 12.05 |
| | Gemini | 25.4 | 21.3 | 4.09 | 7.23 | 9.24 | 9.78 |
| | ChatGPT | 19.7 | 15.3 | 6.80 | 8.57 | 9.26 | 12.74 |
| | GPT-4 | 24.4 | 22.1 | 4.64 | 3.84 | 9.96 | 6.88 |
| | GPT-4o | 18.9 | 15.2 | 5.58 | 7.05 | 10.34 | 9.96 |

Table 12: **Proportional Agreement (PA) score per uncertainty expression in the non-verifiable setting**. The scores are with respect to the population-level human reference distribution (Human+NV).

| Methodology | OLMo (7B) full | Gemma (2B) full | LLama3 (70B) full | LLama3 (8B) full | ChatGPT top-k | GPT-4 top-k | GPT-4o top-k | Mixtral 8x22B sampling | Mixtral 8x7B sampling | |
|---|---|---|---|---|---|---|---|---|---|---|
| Average | 12.2 | 8.6 | 25.1 | 18.8 | 20.5 | 25.0 | 19.1 | 23.9 | 22.3 | 22.0 |
| Standard Deviation | 13.6 | 8.8 | 12.8 | 14.2 | 12.5 | 13.1 | 8.8 | 13.5 | 13.0 | 12.5 |
| almost certain | 55.0 | 0.8 | 60.0 | 60.6 | 55.9 | 60.6 | 35.5 | 58.7 | 60.6 | 59.7 |
| doubtful | 2.6 | 5.5 | 18.9 | 9.8 | 13.2 | 12.9 | 14.2 | 17.5 | 16.4 | 13.7 |
| highly likely | 20.5 | 1.0 | 18.0 | 14.9 | 16.1 | 22.1 | 17.4 | 8.3 | 23.9 | 21.5 |
| highly unlikely | 14.9 | 10.5 | 35.0 | 19.8 | 25.7 | 35.1 | 24.8 | 35.1 | 28.7 | 25.2 |
| not likely | 4.2 | 4.6 | 17.8 | 15.5 | 16.7 | 17.1 | 11.3 | 17.7 | 16.1 | 10.3 |
| possible | 6.7 | 14.3 | 15.3 | 7.7 | 14.2 | 15.4 | 6.5 | 15.0 | 8.6 | 14.8 |
| probable | 7.9 | 3.3 | 15.2 | 7.0 | 7.6 | 14.3 | 14.9 | 14.8 | 15.7 | 13.5 |
| somewhat likely | 9.5 | 2.7 | 16.5 | 4.7 | 6.8 | 17.9 | 12.8 | 16.5 | 18.2 | 15.4 |
| somewhat unlikely | 1.0 | 15.5 | 22.3 | 21.5 | 21.8 | 19.3 | 18.0 | 22.3 | 18.4 | 19.3 |
| uncertain | 6.2 | 34.0 | 35.1 | 35.1 | 33.1 | 35.1 | 35.1 | 35.1 | 35.1 | 35.1 |
| unlikely | 1.6 | 9.2 | 16.6 | 16.8 | 15.9 | 18.3 | 11.1 | 16.8 | 10.1 | 16.4 |
| very likely | 14.4 | 0.8 | 17.4 | 15.2 | 13.6 | 19.5 | 19.1 | 14.4 | 15.1 | 19.1 |
| very unlikely | 13.9 | 9.2 | 38.3 | 15.9 | 26.5 | 37.9 | 27.8 | 38.8 | 22.3 | 21.8 |

we found that the smaller evaluated open-source models (*e.g.*, OLMo (7B) and Gemma (2B)) were extremely sensitive to the prompt formatting and, frequently resulted in negligible probability mass being assigned to integers in the range [0, 100].

## H Ablation: 0-shot vs 2-shot

Table 20 reports the proportional agreement (PA) metric discriminated by uncertainty expression in the non-verifiable and verifiable settings, respectively. The values are reported with respect to the population-level human distribution described in Section 3, denoted Human+NV.

Table 13: **Proportional Agreement (PA) score per uncertainty expression in the verifiable setting**. The PA scores are with respect to the population-level human reference distribution (Human+NV).

| Methodology | OLMo (7B) full | Gemma (2B) full | LLama3 (8B) full | ChatGPT top-k | GPT-4 top-k | GPT-4o top-k | LLama3 (70B) sampling | Mixtral 8x22B sampling | Mixtral 8x7B sampling | Gemini sampling |
|---|---|---|---|---|---|---|---|---|---|---|
| Average | 7.8 | 7.0 | 10.3 | 15.9 | 22.8 | 15.3 | 19.5 | 19.4 | 15.6 | 22.1 |
| Standard Deviation | 12.2 | 8.9 | 5.1 | 11.5 | 13.8 | 10.8 | 14.6 |  | 9.3 | 13.2 |
| almost certain | 46.0 | 0.8 | 20.6 | 48.0 | 60.6 | 28.7 | 35.2 | 60.6 | 36.2 | 54.6 |
| doubtful | 0.9 | 4.2 | 11.2 | 6.2 | 11.5 | 11.6 | 15.2 | 13.8 | 9.3 | 16.4 |
| highly likely | 18.5 | 0.9 | 11.1 | 13.0 | 22.0 | 14.7 | 17.4 | 24.1 | 12.8 | 22.8 |
| highly unlikely | 8.3 | 1.2 | 4.0 | 22.1 | 32.8 | 17.2 | 33.3 | 28.7 | 23.9 | 33.3 |
| not likely | 2.6 | 1.5 | 14.0 | 12.8 | 15.0 | 10.9 | 15.5 | 13.3 | 7.5 | 15.7 |
| possible | 1.2 | 14.8 | 3.5 | 5.8 | 10.2 | 6.1 | 4.0 | 4.6 | 6.0 | 6.7 |
| probable | 1.8 | 3.3 | 6.5 | 3.8 | 11.5 | 11.3 | 6.8 | 5.8 | 5.7 | 9.9 |
| somewhat likely | 1.7 | 3.9 | 7.2 | 6.5 | 12.1 | 7.0 | 8.6 | 9.0 | 9.7 | 10.0 |
| somewhat unlikely | 0.8 | 15.0 | 10.6 | 15.4 | 18.8 | 15.0 | 17.3 | 11.8 | 19.6 | 15.4 |
| uncertain | 0.0 | 32.4 | 14.2 | 25.1 | 34.2 | 30.3 | 34.0 | 33.4 | 31.0 | 34.0 |
| unlikely | 0.4 | 11.3 | 17.7 | 13.6 | 16.8 | 10.6 | 15.0 | 9.4 | 11.6 | 14.3 |
| very likely | 11.6 | 0.7 | 9.3 | 10.1 | 17.2 | 17.2 | 14.9 | 14.9 | 11.1 | 17.6 |
| very unlikely | 7.5 | 1.7 | 4.1 | 24.2 | 34.1 | 18.8 | 35.9 | 22.3 | 18.6 | 36.2 |

Table 14: **Analysis of the pairwise differences between true and false statements in the verifiable dataset**. Using a one-sided Wilcoxon signed-rank test, we find that, for most models, numerical responses for true statements are statistically significantly larger than those for false statements (* = significant at $\alpha = 0.05$).

| Model | Methodology | Statistic | p-value |
|---|---|---|---|
| ChatGPT | top-k | 31180.00 | **<0.0001*** |
| GPT-4 | top-k | 34980.00 | **<0.0001*** |
| GPT-4o | top-k | 36185.50 | **<0.0001*** |
| Gemini | sampling | 25564.00 | **<0.0001*** |
| LLama3 (70B) | sampling | 34979.00 | **<0.0001*** |
| Mixtral 8x22B | sampling | 16051.00 | **<0.0001*** |
| OLMo (7B) | sampling | 5309.50 | **<0.0001*** |
| LLama3 (8B) | full | 17115.00 | **<0.0001*** |
| OLMo (7B) | full | 5309.50 | **<0.0001*** |
| Gemma (2B) | full | 786.50 | 0.9979 |

Table 15: **Analysis of the pairwise differences between true and false statements in the AI2-ARC (Easy) dataset**. Using a one-sided Wilcoxon signed-rank test, we find that, for most models, numerical responses for true statements are statistically significantly larger than those for false statements (* = significant at $\alpha = 0.05$).

| Model | Methodology | Statistic | p-value |
|---|---|---|---|
| ChatGPT | top-k | 183290.50 | **<0.0001*** |
| GPT-4 | top-k | 303749.50 | **<0.0001*** |
| GPT-4o | top-k | 346533.00 | **<0.0001*** |
| Gemini | sampling | 180349.50 | **<0.0001*** |
| LLama3 (70B) | sampling | 250250.00 | **<0.0001*** |
| Mixtral 8x22B | sampling | 129417.50 | **<0.0001*** |
| OLMo (7B) | sampling | 111281.00 | **<0.0001*** |
| Gemma (2B) | sampling | 46702.50 | 0.1148 |

Table 16: **Analysis of the pairwise differences between true and false statements in the AI2-ARC (Challenge) dataset**. Using a one-sided Wilcoxon signed-rank test, we find that, for most models, numerical responses for true statements are statistically significantly larger than those for false statements (* = significant at $\alpha = 0.05$).

| ChatGPT | top-k | 123271.00 | **<0.0001*** |
|---|---|---|---|
| GPT-4 | top-k | 194138.00 | **<0.0001*** |
| GPT-4o | top-k | 231543.50 | **<0.0001*** |
| Gemini | sampling | 106837.00 | **<0.0001*** |
| LLama3 (70B) | sampling | 135453.00 | **<0.0001*** |
| Mixtral 8x22B | sampling | 99689.50 | **<0.0001*** |
| OLMo (7B) | sampling | 48763.50 | 0.0333 |
| Gemma (2B) | sampling | 34161.50 | 0.0191 |

Table 17: **InterQuartile Range (IQR) of the empirical distributions averaged across all uncertainty expressions in the non-verifiable setting**. Reported values include both histograms created using greedy decoding (`temperature=0`) and random decoding (`temperature=1`). While most models exhibit become more diverse when using random decoding, models like `GPT-4` and `LLama3 (70B)` seem to be minimally affected by the change in decoding algorithm, suggesting that these models lead to less diverse results (when compared to humans).

| Models | Methodology | 2-shot prompt | | 0-shot prompt | |
|---|---|---|---|---|---|
| | | Greedy | Probabilistic | Greedy | Probabilistic |
| Humans | — | 15.00 | — | — | — |
| ChatGPT | top-k | 4.62 | 9.00 | 8.08 | 11.54 |
| GPT-4 | top-k | 1.15 | 1.15 | 1.15 | 2.69 |
| GPT-4o | top-k | 3.85 | 39.15 | 3.85 | 39.00 |
| LLama3 (70B) | full | 1.92 | 3.08 | 0.77 | 1.54 |
| LLama3 (8B) | full | 15.38 | 21.54 | 10.00 | 21.92 |
| OLMo (7B) | full | 36.54 | 39.23 | 36.15 | — |
| Gemma (2B) | full | 23.08 | 19.62 | 0.77 | 11.54 |
| Mixtral 8x22B | sampling | 5.00 | — | 5.77 | — |
| Mixtral 8x7B | sampling | 4.23 | — | 4.62 | — |

Table 18: **Analysis of the distributional differences between the conditional distributions estimated using greedy (`temperature=0`) *versus* probabilistic decoding (`temperature=1`)**. We use Wasserstein-1 distance to report the distributional differences between each conditional distribution: maximally distant distributions exhibit a score of 101.

| Model | 0-shot prompting | | 2-shot prompting | | | |
|---|---|---|---|---|---|---|
| | NV | V | NV | V | AI2-ARC (Easy) | AI2-ARC (Challenge) |
| ChatGPT | 3.28 | 3.40 | 3.32 | 2.75 | 2.87 | 3.08 |
| GPT-4 | 0.31 | 0.52 | 0.22 | 0.84 | 0.49 | 0.78 |
| GPT-4o | 19.88 | 15.51 | 15.38 | 19.86 | 20.70 | 20.97 |

Table 19: **Summary metrics averaged across uncertainty expressions for both NV and V settings when using probabilistic decoding (`temperature=1`)**. All metrics are computed with respect to the human distribution in the non-verifiable setting (`Human+NV`). "PA" reports the general agreement between LLMs and the mode of the human distribution, reported in percentages. "MAE" reports the absolute error between the mean responses of LLMs and those of humans. Wasserstein-1 computes the distance between LLMs and human distributions.

| | | Avg PA (↑) | | Avg MAE (↓) | | Avg Wasserstein-1 (↓) | |
|---|---|---|---|---|---|---|---|
| | | NV | V | NV | V | NV | V |
| Human | Mode | 27.6 | 27.6 | — | — | — | — |
| | Individual | 17.6 | 16.7 | 8.91 | 9.35 | 12.35 | 12.99 |
| LLM | ChatGPT | 16.4 | 12.8 | 6.40 | 8.32 | 7.65 | 12.14 |
| | GPT4 | 24.4 | 21.4 | 4.62 | 4.00 | 9.78 | 6.72 |
| | GPT4o | 12.9 | 8.7 | 19.01 | 26.07 | 19.69 | 26.14 |

(a) "almost certain"  (b) "highly likely"  (c) "highly unlikely"

(d) "possible"  (e) "very likely"  (f) "very unlikely"

(g) "probable"  (h) "likely"  (i) "not likely"

(j) "unlikely"  (k) "somewhat likely"  (l) "somewhat unlikely"
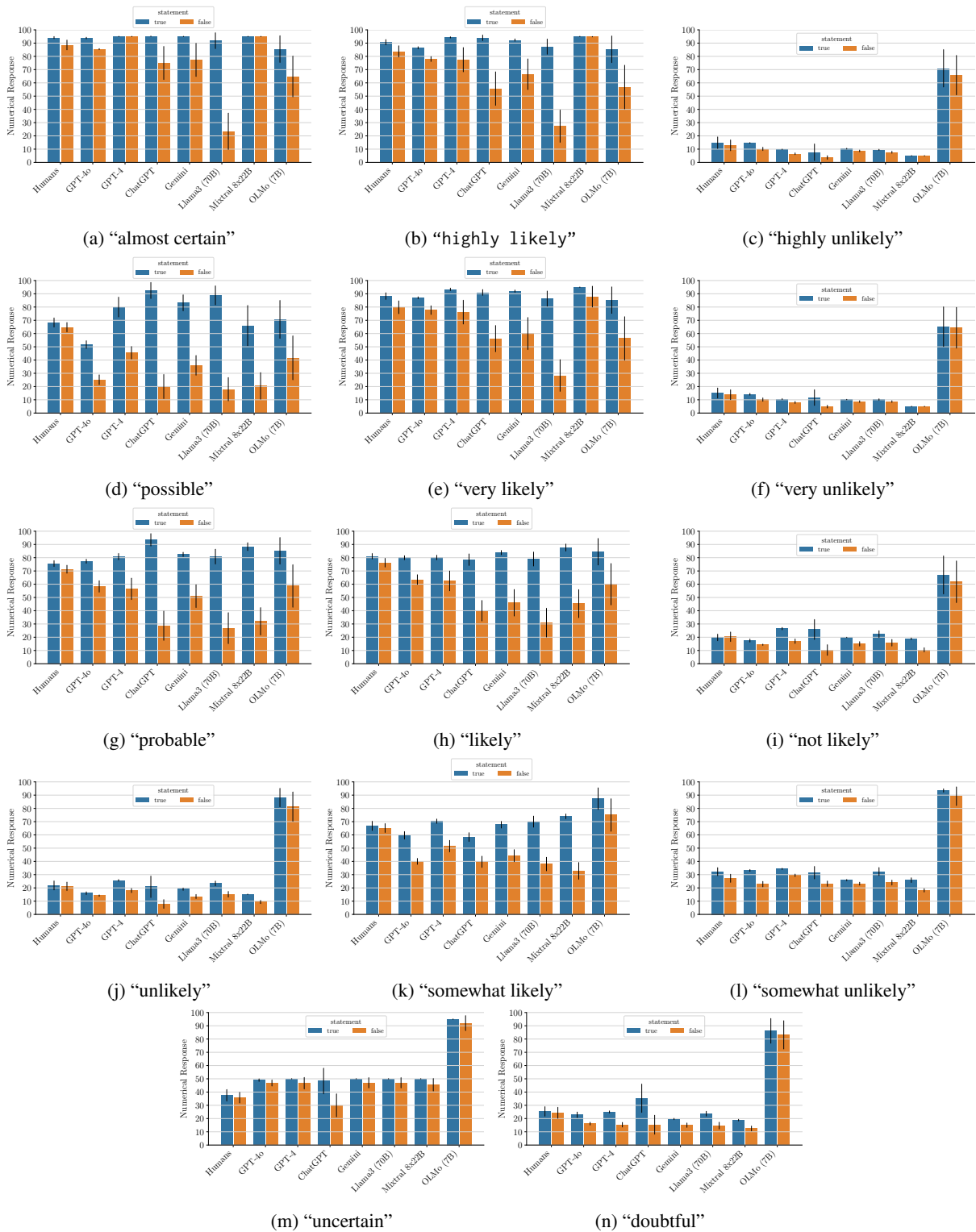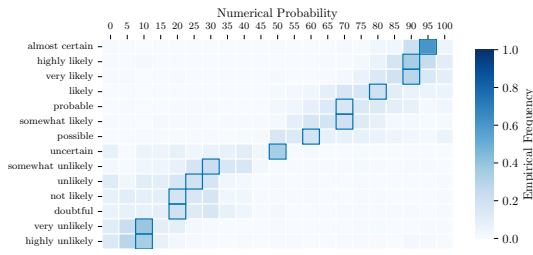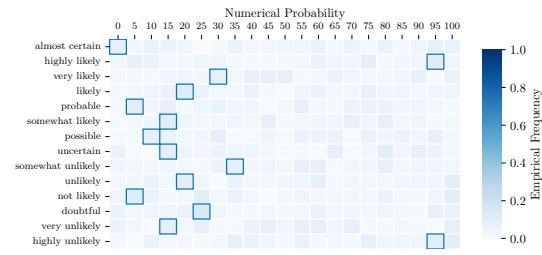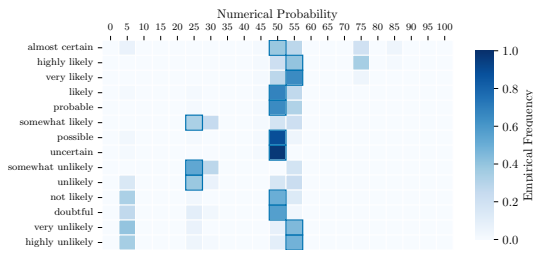
(m) "uncertain"  (n) "doubtful"

Figure 13: **Mean numerical response for the verifiable statements discriminated by truthfulness of the statements**. We observe that the differences in mean numerical responses differ by uncertainty expressions. Focusing on the *opposing* linguistic expressions (*e.g.*, "unlikely" *vs* "likely"), the mean numerical response gap observed for models seems to be systematically larger for positive expressions (*e.g.*, "likely") than for negative expressions (*e.g.*, "unlikely").

Figure 14: **Models (and Human) empirical distributions of numerical responses per uncertainty expression in the non-verifiable setting (NV)**. The empirical distributions are estimated using the greedy decoding algorithm (temperature=1). Highlighted boxes represent the mode value for each expression. With the exception of OLMo (7B) and Gemma (2B), LLMs are able to replicate the overall population-level human behavior: assigning lower numerical responses to uncertainty expressions (*e.g.*, "highly unlikely" and "unlikely") and progressively higher numerical responses until reaching the certainty expressions (*e.g.*, "almost certain").
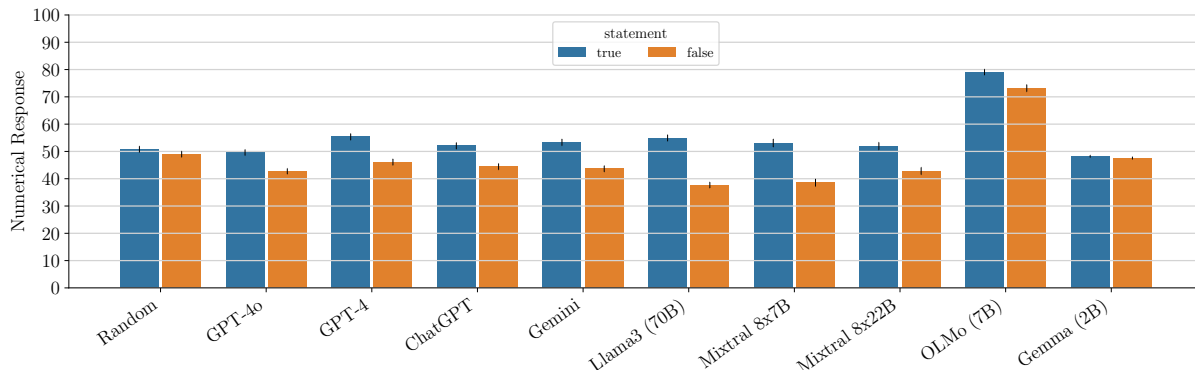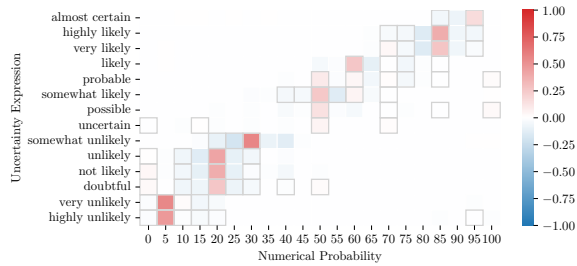
Figure 15: **Mean numerical response discriminated by truthfulness of statements for 400 verifiable statements derived from the AI2-ARC dataset**. The mean numerical responses produced by LLMs when evaluated in the context of true statements is significantly larger than when evaluated with the false statements. Although the numerical response gap is lower in magnitude than those observed in the verifiable experiment (see Figure 6), the observed gap is still statistically significant.

Table 20: **Proportional Agreement (PA) score per uncertainty expression in the non-verifiable setting**. The scores are with respect to the population-level human reference distribution (Human+NV).
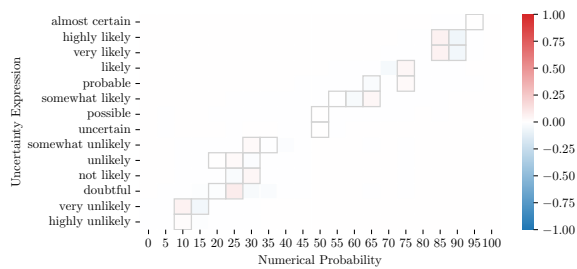
| Methodology | OLMo (7B) full | Gemma (2B) full | LLama3 (70B) full | LLama3 (8B) full | ChatGPT top-k | GPT-4 top-k | GPT-4o top-k | Mixtral 8x22B sampling | Mixtral 8x7B sampling |
|---|---|---|---|---|---|---|---|---|---|
| Average | 9.3 | 2.4 | 18.3 | 14.4 | 17.1 | 25.9 | 22.3 | 20.2 | 18.8 |
| Standard Deviation | 9.2 | 4.0 | 6.4 | 5.2 | 8.3 | 12.2 | 13.1 | 14.5 | 13.1 |
| almost certain | 35.3 | 0.2 | 27.8 | 15.5 | 33.0 | 58.0 | 57.4 | 60.6 | 59.0 |
| doubtful | 8.4 | 0.0 | 16.5 | 16.5 | 8.7 | 18.0 | 12.7 | 8.9 | 10.7 |
| highly likely | 21.2 | 0.4 | 8.7 | 5.1 | 12.0 | 27.2 | 19.4 | 25.6 | 16.3 |
| highly unlikely | 4.5 | 0.5 | 16.0 | 14.1 | 22.6 | 35.1 | 33.0 | 28.7 | 16.7 |
| not likely | 4.3 | 0.5 | 17.6 | 17.6 | 13.8 | 17.9 | 11.5 | 9.9 | 10.8 |
| possible | 3.1 | 12.7 | 15.2 | 12.8 | 14.5 | 15.4 | 10.7 | 9.5 | 14.9 |
| probable | 6.9 | 4.9 | 15.9 | 8.7 | 10.2 | 14.7 | 15.4 | 9.0 | 13.4 |
| somewhat likely | 2.2 | 9.7 | 16.0 | 10.9 | 8.1 | 16.6 | 11.0 | 16.5 | 6.1 |
| somewhat unlikely | 3.5 | 0.6 | 22.2 | 21.2 | 18.4 | 22.3 | 19.5 | 7.5 | 10.8 |
| uncertain | 3.0 | 0.6 | 35.1 | 25.5 | 34.5 | 35.1 | 34.1 | 34.5 | 32.6 |
| unlikely | 6.0 | 0.0 | 17.0 | 16.8 | 13.6 | 16.7 | 11.0 | 9.2 | 15.2 |
| very likely | 16.6 | 0.9 | 14.5 | 8.9 | 11.3 | 21.1 | 21.4 | 20.7 | 18.5 |
| very unlikely | 6.1 | 0.0 | 15.6 | 14.2 | 21.2 | 38.8 | 32.1 | 22.3 | 19.2 |

(a) almost certain

(b) "highly likely"

(c) "highly unlikely"

(d) "possible"

(e) "very likely"

(f) "very unlikely"

(g) "probable"

(h) "likely"

(i) "not likely"

(j) "unlikely"

(k) "somewhat likely"

(l) "somewhat unlikely"

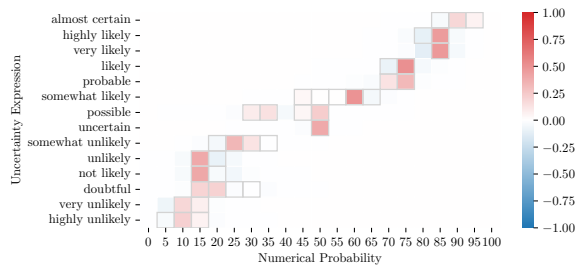(m) "uncertain"

(n) "doubtful"

Figure 16: **Mean numerical response for verifiable statements discriminated by truthfulness of statements for the AI2-ARC dataset**. While the observed gaps are smaller in magnitude than in the main experiment, there is still a significant difference between mean numerical responses depending on the truthfulness of the statements. We hypothesize that the observed magnitude differences between the two datasets may be explained by semantic differences between the two QA datasets used to curate the verifiable statements. Future work should investigate this behavior more closely.

(a) ChatGPT



(b) GPT-4



(c) GPT-4o

Figure 17: **Absolute Difference between the models' conditional probability distribution when estimated using greedy (`temperature=0`) *versus* probabilistic decoding (`temperature=1`)**. The results refer to the non-verifiable dataset. Gray rectangles indicate bins whose empirical distribution with the greedy information has non-zero value. Bins are colored blue if the bin's estimated probability is larger when using the probabilistic approach and are red if the probability was larger using the greedy approach. While we observe that the decoding choice has a larger impact for `ChatGPT` and `GPT-4o`, we find that it has minimal impact in `GPT-4`'s distributions, suggesting that `GPT-4`'s predictions tend to be very confident in its predictions to begin with.
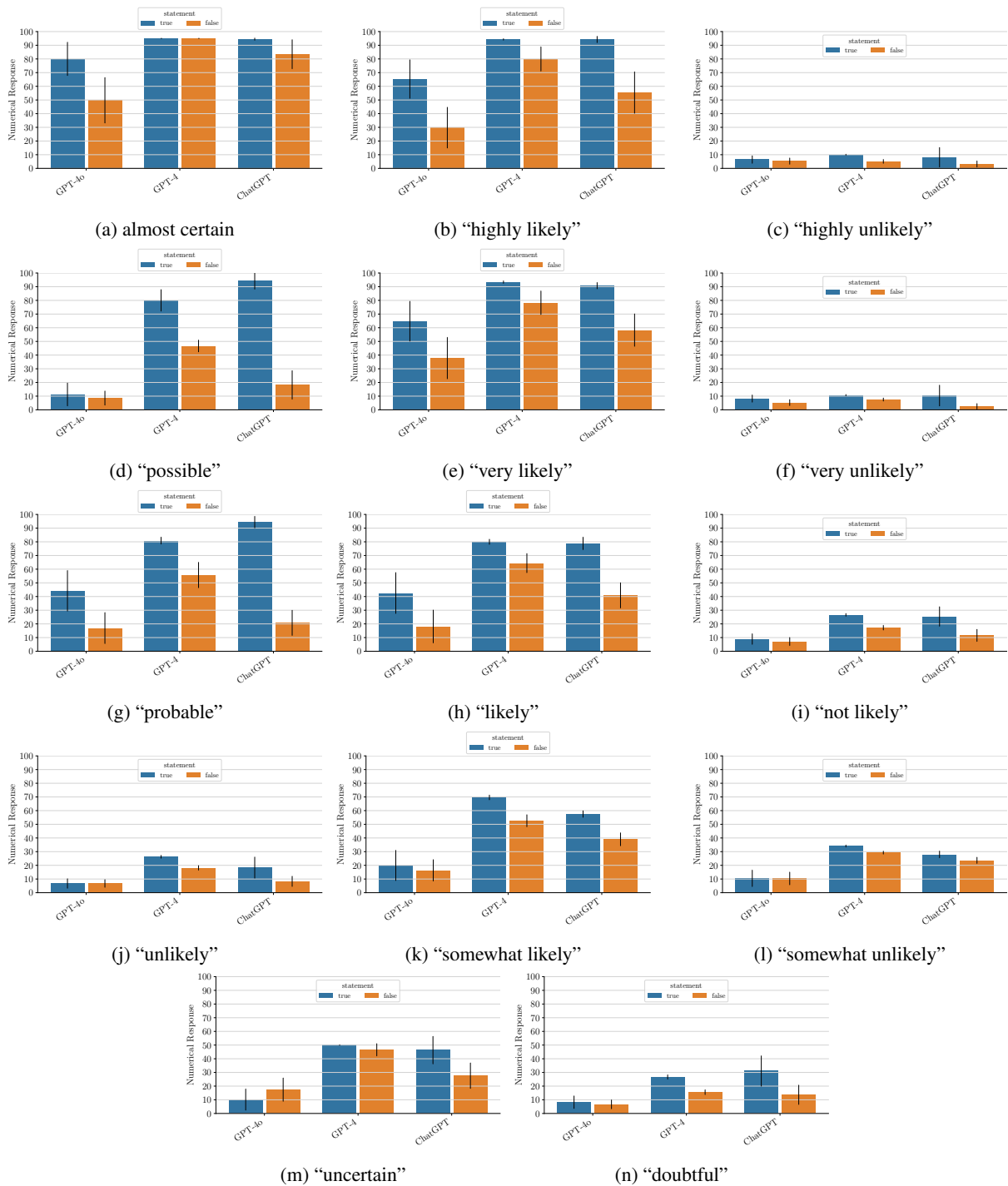
Figure 18: **Mean numerical response discriminated by truthfulness of verifiable statements, estimated using probabilistic decoding (`temperature=1`).** Overall, we still observe statistically meaningful gaps depending on the truthfulness of the statement.