# Large Language Model Prediction Capabilities: Evidence from a Real-World Forecasting Tournament

**Philipp Schoenegger**
London School of Economics
and Political Science

**Peter S. Park**
MIT

## Abstract

Accurately predicting the future would be an important milestone in the capabilities of artificial intelligence. However, research on the ability of large language models to provide probabilistic predictions about future events remains nascent. To empirically test this ability, we enrolled OpenAI's state-of-the-art large language model, GPT-4, in a three-month forecasting tournament hosted on the Metaculus platform. The tournament, running from July to October 2023, attracted 843 participants and covered diverse topics including Big Tech, U.S. politics, viral outbreaks, and the Ukraine conflict. Focusing on binary forecasts, we show that GPT-4's probabilistic forecasts are significantly less accurate than the median human-crowd forecasts. We find that GPT-4's forecasts did not significantly differ from the no-information forecasting strategy of assigning a 50% probability to every question. We explore a potential explanation, that GPT-4 might be predisposed to predict probabilities close to the midpoint of the scale, but our data do not support this hypothesis. Overall, we find that GPT-4 significantly underperforms in real-world predictive tasks compared to median human-crowd forecasts. A potential explanation for this underperformance is that in real-world forecasting tournaments, the true answers are genuinely unknown at the time of prediction; unlike in other benchmark tasks like professional exams or time series forecasting, where strong performance may at least partly be due to the answers being memorized from the training data. This makes real-world forecasting tournaments an ideal environment for testing the generalized reasoning and prediction capabilities of artificial intelligence going forward.

## 1 Introduction

In the field of artificial intelligence (AI), large language models (LLMs) have recently shown surprising capabilities in a multitude of economically relevant tasks[25] that were previously thought to require human cognition.[6] State-of-the-art LLMs are comprised of an extremely large amount of parameters—typically organized in terms of the Transformer architecture[37]—and trained on a large corpus of Internet-based text data. This corpus of training data is used to train LLMs to predict the next sequence of tokens given an input. Despite the simplicity of the general task for which LLMs are trained—next token prediction—the resulting transfer learning causes LLMs to also become ostensibly proficient at a wide variety of specific tasks, including reading comprehension,[30,38] summarization,[11] translation,[15] coding,[6] deception,[32] medical-license exams[6,27], and bar exams.[6,17]

The ostensibly impressive capabilities of LLMs come with several important caveats. One of the most important caveats to keep in mind is that when an LLM consistently outputs the true answers to questions in a task benchmark or exam, it is unclear whether the LLM's true answers reflect a genuine understanding of the task: and as such, are likely to generalize to out-of-distribution settings.[1] An alternative hypothesis is that the training dataset contains at least part of the task benchmark's questions and corresponding answers, and that the given LLM capability may thus not reflect a

genuine understanding that can generalize to analogous questions not in the training dataset.[2,4,8,19] For example, consider a scenario where an organic-chemistry exam question is inputted into the LLM, and it outputs the true answer. It remains unclear whether this is due to the LLM possessing a deep understanding of the relevant organic-chemistry concepts, as opposed to simply reproducing the answer to the specific question contained in its training dataset. It is not trivial to rigorously formalize the dichotomy between genuine understanding and training-data memorization, as genuine understanding also ultimately arises from the relevant content in the training dataset. But the very real phenomenon of generalizability—or lack thereof—seems to be at the heart of the dichotomy.[12]

Given the problem of distinguishing whether an LLM's proficiency at a task benchmark is due to genuine understanding or to memorizing the training data, we propose evaluating LLM capabilities in contexts where the true answers are unknown beforehand and as such cannot be part of the training data. An especially important such context is forecasting, a domain where the answers are initially unknown to everyone (even the human evaluators) until they are validated or falsified by future events.[33] We distinguish this context from the related but different forecasting context of zero-shot extrapolations of time series.[13] In contrast, our definition of 'forecasting'—producing accurate probabilistic predictions of future events—does not encounter the problem of distinguishing the benchmark data from the training data. As such, it allows for a more robust and generalizable evaluation of LLM capabilities: beyond the rote memorization of training data.

Our study evaluates the forecasting capabilities of GPT-4, a state-of-the-art LLM created by OpenAI,[28] by entering it in a forecasting tournament. Forecasting tournaments are competitions wherein individual forecasters provide probabilistic forecasts on questions concerning future events.[36] These predictions are then scored based on the accuracy of their forecasts: the closer one's prediction is to the truth, the more likely one is to be rewarded. The collective accuracy of the predictions resulting from such forecasting tournaments hinges on the 'wisdom of the crowd,' i.e., the observation that aggregated forecasts of a group are often more accurate than the forecasts of individuals,[7,20] even if judgments are correlated and biased.[10] This setting allows us to test LLM capabilities in a context where the answer is genuinely unknown beforehand and thus cannot be contained within the training dataset. This setting may also be critical for whether LLMs are or will be able to rival—or even outperform—humans at jobs that require prescient decision-making.[31] Accurately predicting the future is foundational to decision-making in most public and private sectors.[33] Consequently, comparing the forecasting capabilities of GPT-4 with those of humans promises to be a good test case for many potential economic applications of advanced AI.

As such, we have pre-registered—along with our analysis plans—the following null hypothesis:

> The mean accuracy of GPT-4 forecasts is not different from the mean accuracy of median human-crowd forecasts. $H_0 : \mu_{\text{GPT-4}} = \mu_{\text{Human}}$

Human crowds in forecasting tournaments are among the best options for producing accurate probabilistic predictions.[9,35] If a state-of-the-art LLM like GPT-4 is able to outperform a human crowd of forecasters, this would be consistent with the model having learned a deep understanding of the relevant capabilities, like probabilistic reasoning, generalization, and accurate prediction.

## 2 Methods

We conducted our study on Metaculus, a platform that hosts forecasting tournaments where members of the public can submit predictions, compete for prizes, and establish a forecasting track record. Metaculus' forecasting platform has been employed in various academic and policy prediction contexts, such as the monkeypox outbreak in 2022[22] and the COVID pandemic.[21] We leveraged the opportunity provided by the launch of the Quarterly Cup[24] on July 3, 2023, as this three-month tournament offered an ideal context to evaluate LLM prediction capabilities. This setting provided us with a context where an LLM could compete with human forecasters to make probabilistic predictions about various topical questions.

Our questions were vetted by the Metaculus moderation team, who aimed to resolve the questions in a consistent and accurate manner in order to ensure high-quality data. The forecasting questions studied here were on a wide array of topics, such as U.S. industrial action disputes, military interventions in Niger, outbreaks of Marburg virus disease, and the Black Sea grain deal. For examples of the forecasting questions used in our study, see Table 1. For a full set of questions, see Appendix A.

| |
|---|
| Will the United Auto Workers call a strike against any of the Big Three Detroit automakers before September 19, 2023? |
| Will Mohamed Bazoum, Nigerien President, return to power before August 31, 2023? |
| Will India's Chandrayaan-3 mission successfully land a rover on the moon? |
| Will the Black Sea grain deal be revived before October 1, 2023? |
| Will a non-proprietary LLM be in the top 5 of the chat.lmsys.org leaderboard on September 30, 2023? |

**Table 1:** Examples of forecasting questions used in our study.

Our study's questions were answered by both GPT-4 and a large set of human forecasters. The tournament concluded on October 4, with a total of 51 questions posed, and with 843 unique forecasters entering at least one prediction. For our analysis, we focused solely on the subset of binary questions, of which there were 23. This is because binary questions were the only types of questions that our LLM forecaster could straightforwardly answer without requiring additional human input. This could have biased the predictions, such as by drawing distributions based on quartile point estimates.

We used the web interface of GPT-4 at the default temperature value. Temperature is a hyperparameter that controls the randomness of a model's output, with higher temperature settings resulting in more random outputs, and lower temperature settings resulting in more deterministic outputs. Our prompt was crafted by drawing on established research in forecasting, with the aim of steering GPT-4 towards comfortably providing numerical predictions and of enhancing the accuracy of these predictions. First, we prompted the model to emulate a superforecaster,[36] aligning with the best-practice recommendations of how to prompt models to act as if they were domain experts.[39] This part of the prompt enabled us to consistently avoid getting responses of a different format than explicit probabilistic forecasts.

Second, we grounded our prompt on research from the forecasting literature that highly complex qualitative rationales and the use of base rates are associated with forecasting accuracy.[16] This part of the prompt was an attempt to not only elicit probabilistic forecasts, but to ensure that these are as highly accurate as possible.

The overall prompt is given by the following.

> **Prompt**: In this chat, you are a superforecaster that has a strong track record of accurate forecasts of the future. As an experienced forecaster, you evaluate past data and trends carefully and aim to predict future events as accurately as you can, even though you cannot know the answer. This means you put probabilities on outcomes that you are uncertain about (ranging from 0 to 100%). When the outcome is continuous, you give me 25th interquartile ranges. You also quickly outline your rationale. In your rationales, you carefully consider the reasons for and against your probability estimate, you will make use of comparison classes of similar events and probabilities and take into account base rates and past events as well as other forecasts and predictions. You will also consider different perspectives.

Each question of the tournament was predicted via a single chat log, with the prompt positioned at the outset. Following this, GPT-4 was supplied with the background information, the specific resolution criteria of the question, and the final question text exactly as presented on Metaculus. This information, provided by Metaculus users or staff, helped mitigate the temporal distance between the training data set cutoff data and the question launch, which is less fundamental of an information gap for the human comparison group. For a full set of this information for one question, see Appendix B. For some questions, the median human-crowd prediction was available at the time of forecast; while for others, it was not. In the case of the latter, we included the median human-crowd prediction in the description provided to GPT-4.

Subsequently, we collected the probabilistic prediction and inputted it into the question, in percent form. In cases where the model output was a probability range, the mean probability of the two

values was used instead. We abstained from making any revisions to the initial forecasts. We did so because determining when to update and what information to employ for this update would have necessitated a significant degree of human involvement in the process, which we sought to minimize given that updating is a crucial skill for human forecasters.[23] Since repeated updating is fundamental to the method of crowd-aggregated forecasting, it was important to account for GPT-4's inability to autonomously carry out this task in order to enable unbiased inferences from our results.

In order to control for these confounding factors, our design entailed recording both the initial model forecast and the crowd-aggregated median forecast at the close of the same day at which the model forecast was submitted. This allowed us to control for updating over extended question periods, which is done by humans, but is difficult to operationalize for GPT-4 without introducing bias in deciding when and how to prompt the update. Our comparison thus directly put GPT-4 on par with the human forecasters by only looking at initial forecasts. The full data encompassed for each question GPT-4's forecast, the median human-crowd forecast at the conclusion of the same day, the ground truth, whether a community prediction was visible at the time of prediction, the number of human forecasters at this point in time (for a median number of 37 forecasters), and the total duration, spanning the timeframe from when predictions were made to the eventual resolution of the question. This set of data was collected for all 23 binary forecasts.

# 3 Results

First, we examine the probabilistic forecasts of both GPT-4 and the human crowd for each question in a directional analysis. See Figure 1 for a paired comparison plot illustrating the respective probability forecasts, with the upper panel displaying questions that resolved positively and the lower panel displaying questions that resolved negatively. These data indicate that in 18 out of 23 questions, the median human-crowd forecasts were directionally closer to the truth than GPT-4's predictions, $\chi^2(1) = 12.52, p = .001$.

We also compared the frequency with which each of GPT-4's forecast and the aggregate human forecast was directionally correct, i.e., on the correct side of the probability midpoint. When doing so, we observed that the GPT-4 forecast was directionally correct 69.57% of the time, in contrast to 95.65% of the time for the aggregate human forecast, although the difference in distributions between the two was not statistically significant, $\chi^2(1) = 3.78, p = .052$. Similarly, the 69.57% proportion of directionally correct answers for GPT-4 also did not statistically deviate from the random 50% baseline, $Z = 1.88, p = .061$. These data suggest that while the aggregate human forecasts were comparatively more often closer to the truth than those of GPT-4, we do not find statistically significant results when testing this independently in relation to the probability scale midpoint.

The primary pre-registered outcome of interest for our analysis is forecasting accuracy. To compute this accuracy as our dependent variable, we employ Brier scores,[5] with the score for each individual provided below. Here, $f$ is the forecasted probability and $o$ is the observed outcome (which is 1 if the event occurs and 0 if it does not), given by $B = (f - o)^2$. We then aggregate the question-level Brier scores per condition as

$$\frac{1}{N} \sum_{i=1}^{N} (f_{i,\text{condition}} - o_i)^2. \tag{1}$$

A perfect accuracy would yield a Brier score of 0, while perfect inaccuracy would result in a Brier score of 1.

In our data, we calculate the mean accuracy by aggregating the the question-level Brier scores for GPT-4's predictions and the median human-crowd forecasts. We observe an average Brier score for GPT-4's predictions of $B = .20$ ($SD = .18$), while the human forecaster average Brier score was $B = .07$ ($SD = .08$). Initially, we test both GPT-4's and the human crowd's accuracy against a simple no-information baseline. This baseline is the Brier score of 0.25, equivalent to predicting 50% on each question. This serves as a first test of prediction accuracy. Our analysis reveals that GPT-4 does not exhibit predictive performance that is statistically distinct from the no-information baseline, $t(22) = -1.23, p = 0.23$. On the other hand, we find that the human crowd's accuracy is significantly superior to that of the no-information baseline, $t(22) = -10.69, p < .001$. These
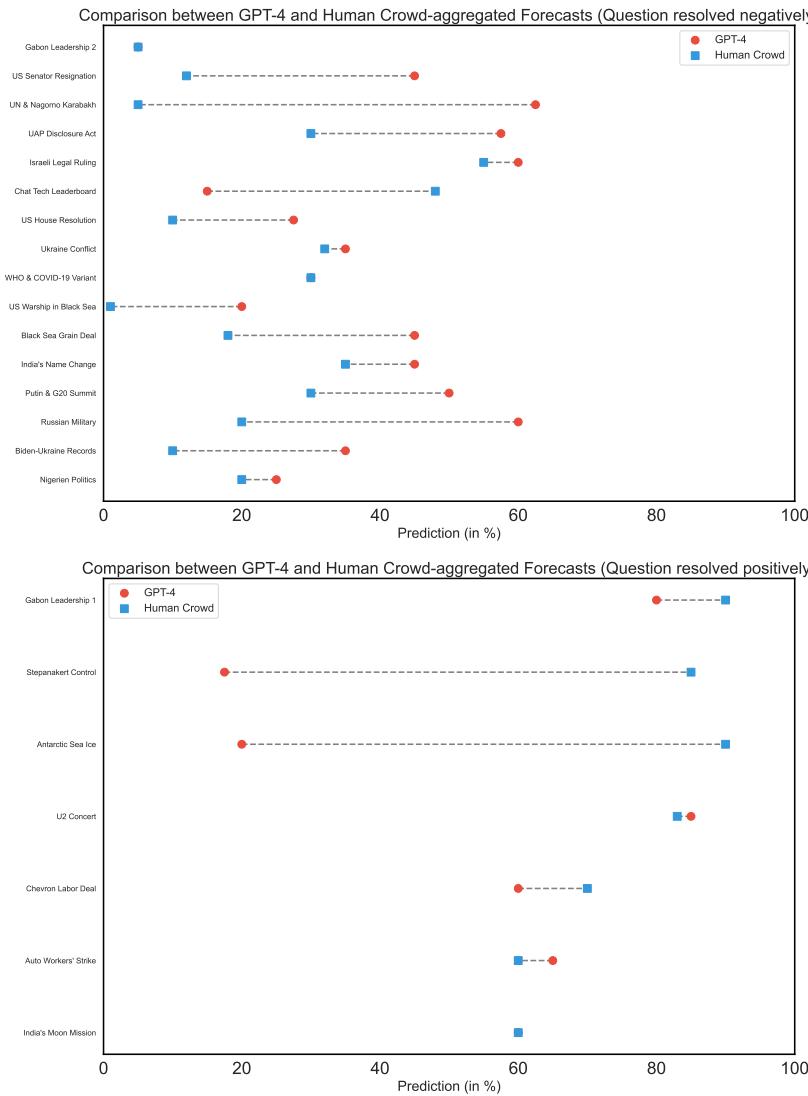
**Figure 1:** Paired comparison plot showing probability forecasts (in %) per question for the median human-crowd forecast and the GPT-4 forecast. The top panel lists questions that resolved negatively. The bottom panel lists questions that resolved positively.

findings suggest that while the aggregate human-crowd forecasts have high accuracy, the mean LLM forecasts do not significantly improve on the no-information baseline.

In order to test our pre-registered null hypothesis, we compare the mean accuracy of GPT-4 and that of the human crowd. This yields a total mean difference of $\Delta = .13$, which is statistically significant at the two-tailed test, $t(44) = 3.11, p = .003$, with an effect size of Cohen's $d = .94$. Based on this, we reject our null hypothesis, and conclude that GPT-4 significantly underperforms in this real-world forecasting tournament compared to the median human-crowd forecasts that are currently employed in forecasting tournaments. We also find that GPT-4's accuracy was not significantly different on questions that included the community prediction compared to those that did not, $t(21) = -1.81, p = 0.085$, suggesting that this result is explained by an internal factor of GPT-4's forecasting performance. For a graphical depiction of the distribution of Brier scores per condition, see Figure 2.

One potential reason that might explain this outcome could be that GPT-4's predictions, due to its reinforcement learning from human feedback,[40] may gravitate towards the probability scale
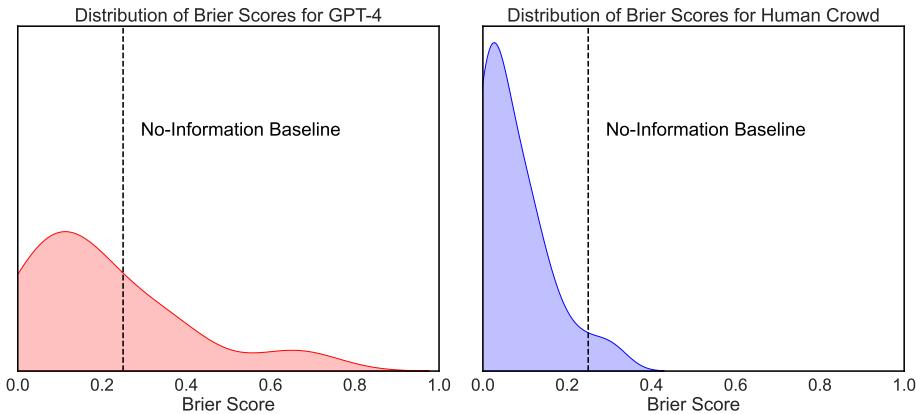
5

**Figure 2:** Kernel Density Estimation (KDE) plots of Brier scores for GPT-4's forecasts and median human-crowd forecasts. Black dotted line represents a Brier score of 0.250 for the 'No-Information Baseline'.

midpoint (50%) in a form of model conservatism. This could explain the low prediction accuracy in our forecasting tournament with a small-to-medium sample size of questions. To analyse this, we conducted the following exploratory analyses. First, we computed the coefficient of variation (CV) as a normalized measure of dispersion for the predictions made by both GPT-4 and the human crowd. Our analysis found a CV of 48.22% for GPT-4 and 73.67% for the human crowd when centered on the mean, suggesting that GPT-4 has less dispersion. Similarly, when anchoring the CV around the midpoint of the probability scale, the values were 42.14% for GPT-4 and 57.59% for the human crowd respectively, again showing the pattern of GPT-4 having less dispersion.

These analyses might suggest that GPT-4's predictions could be more densely clustered around both the mean and the midpoint compared to the human crowd, aligning with a more cautious prediction pattern: one that is less inclined towards extreme forecasts. However, our inferential analysis employing Levene's test for equality of variances at the 50% midpoint did not yield statistically significant evidence to substantiate differences in variance around the midpoint, $F(1, 44) = .54, p = .47$. Thus, our results do not provide grounds for asserting a difference in variance around the probability scale midpoint. See Figure 3 for a graphical representation of these findings.
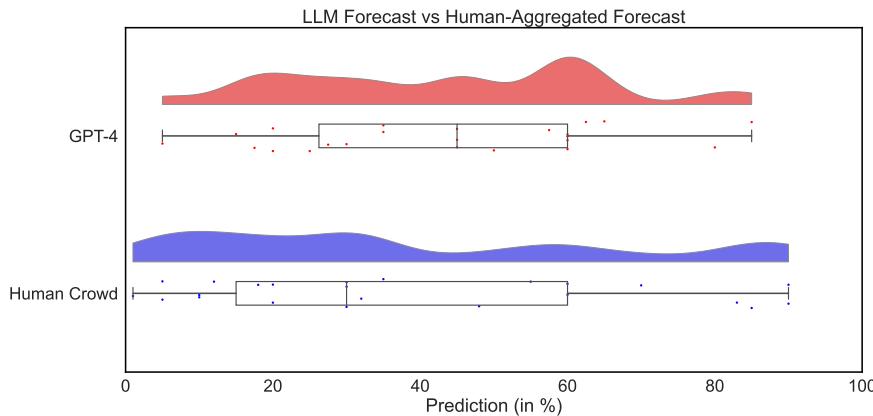


**Figure 3:** Raincloud plot of the distribution of probability forecasts (in %) made by GPT-4 and by the median human-crowd forecast for all questions. Box plots represent interquartile ranges.

Furthermore, in an additional exploratory analysis, we probe whether the relationship between the duration of a forecasting question remaining unresolved and the forecasting accuracy at the inception of this period significantly diverges between GPT-4's forecasts and human-crowd forecasts. This

is interesting because it may be that GPT-4's forecasts are especially good, or especially bad, at predicting questions that resolve very quickly compared to those that do not. See Figure 4 for a scatterplot illustrating accuracy and question duration with linear fits. Our results show that we do not detect statistically significant effects of duration, $ps > .36$.
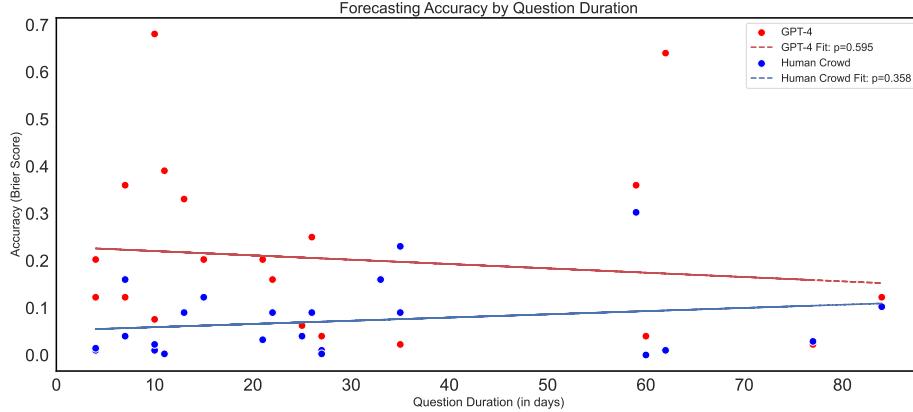


**Figure 4:** Scatterplot and linear fit for the relationship between forecasting accuracy (Brier score) and duration (in days) by condition.

To test this question directly, we conducted an Analysis of Covariance (ANCOVA) to assess whether the relationship between accuracy and question duration differs between GPT-4's forecasts and those made by human forecasters. The main effect of the method (human crowd vs. GPT-4) was statistically significant, $F(1, 42) = 9.38, p = .004$, indicating a notable difference in accuracy between the two forecasting methods. However, the main effect of duration was not significant, $F(1, 42) = 0.017, p = .898$, suggesting that question duration did not significantly affect the accuracy of forecasts. Moreover, the interaction between method and duration was also not significant, $F(1, 42) = 0.75, p = .392$, suggesting that the relationship between question duration and accuracy did not significantly vary between GPT-4 and the human crowd.

Additionally, there are also a variety of aggregation techniques that may prove useful in combining the aggregate human forecast and GPT-4's forecast. We report the exploratory results of a Bayesian Model Averaging (BMA) approach, which takes into account the uncertainty across models prior to combining the predictions from multiple models. We first calculate the likelihoods of each model (.82 for GPT-4 and .93 for the human crowd) before normalizing them and computing the posterior model probabilities. With these, we have the following weighted average:

$$\text{Brier}_{\text{BMA}} = (0.467 \cdot \text{Brier}_{\text{GPT-4}}) + (0.533 \cdot \text{Brier}_{\text{Human Crowd}}). \tag{2}$$

We compute $\text{Brier}_{\text{BMA}} = .13$ ($SD = .09$), which is significantly more accurate than the random baseline, $t(22) = -5.91, p < .001$.

We have presented this simple aggregation as an instance of how LLM predictions may feed into future aggregation approaches. We point out that given our finding of poor prediction accuracy, current models like GPT-4 may be unlikely to make worthwhile additions to aggregation algorithms. However, this may change with advances in LLM capabilities and LLM-forecasting technology.

## 4  Discussion

Our findings from entering GPT-4 into a real-world forecasting tournament on the Metaculus platform suggest that even this state-of-the-art LLM has unimpressive forecasting capabilities. Despite being prompted with established superforecasting techniques and best-practice prompting approaches, GPT-4 was heavily outperformed by the forecasts of the human crowd, and did not even outperform a no-information baseline of predicting 50% on every question. The robustness of this finding is suggested by the fact that the question set was diverse, drawing on a wide variety of topics like Big Tech, U.S. politics, viral outbreaks, and the Ukraine conflict. We argue that our data also provide

additional evidence for the predictive power of human crowd forecasting competitions, as their accuracy was impressive throughout the questions studied here. Our results thus suggest that current LLMs may not yet perform well in a variety of real-world prediction tasks that would necessitate probabilistic foresight.

One potential contributing factor to GPT-4's inadequate forecasting capability is the fact that its training data is subject to a knowledge cutoff after a certain point in time.[28] In contrast to human forecasters who keep up-to-date with recent events, GPT-4 is not updated with events that occur post-training. Given the dynamic nature of many world events, GPT-4's lack of real-time knowledge updating can be a significant barrier to accurate forecasting. Our design tried to address this concern by inputting the question's background information presented on Metaculus to the LLM, so that it would not be wholly unaware of a potentially novel context. While this does not fully bring the model up to speed with human forecasters that seek out prediction competition on Metaculus, it does help mitigate the information gap posed by the knowledge cutoff. Additional ways to mitigate the limitations to forecasting posed by LLMs' knowledge cutoff would be fruitful in future studies.[18]

OpenAI's mission is to create "highly autonomous systems that outperform humans at most economically valuable work."[29] Whether this AI-led future is on track to occur will likely be in large part determined by how capable LLMs—and AI systems in general—turn out to be at economically relevant tasks,[31] especially without human hand-holding. Forecasting is a task of especially high economic relevance, especially for a large proportion of white-collar fields—such as business, policy, and law—that rely heavily on accurate predictions across a myriad of contexts. As of now, our results suggest that even state-of-the-art AI systems are not yet posed to replace human expertise in these areas, due to the inadequate forecasting capabilities of these systems. However, it will be especially important to closely monitor AI capability advances in the domain of accurate forecasting, as foresight will remain an essential skill for autonomous systems.

Another relevant implication of AI forecasting capabilities—or lack thereof—is the risk arising from AI systems that are proficient at long-term planning. An AI system with robust long-term planning capabilities would be able to presciently pursue their goal, which can be a sizable and potentially catastrophic danger if the goal happens to be incompatible with the well-being of humans[26] (e.g., the goal of engineering a pandemic that kills as many people as possible over the long run). Proficiency in long-term planning requires the ability to accurately forecast future scenarios. For many real-world tasks, such accurate forecasting is a highly complex endeavor, at which even many (but not all) humans arguably perform poorly.[34] Our finding that GPT-4 has particularly poor forecasting capabilities bolsters the case that the threat of an AI system planning in the long term against human interests remains presently low. However, this should not be taken as a reason to be complacent. The pace at which AI capabilities advance suggests that it remains crucial to continually monitor the progression of AI systems in terms of their forecasting abilities to ensure that their development remains robustly safe.

Our results raise several avenues for further research. First, it may be argued that our findings might, at least in part, be explained by the human comparison group's ability to harness the wisdom of the crowd,[14] whereas the LLM prediction might be best understood as a single forecast only. One counterargument to this claim is that LLM predictions themselves may draw on a wisdom-of-the-crowd effect from their large and multifaceted training dataset. Future research on how to employ a larger ensemble of LLM forecasters that draw on diverse inputs, training datasets, prompts, temperature values, and other variations would be fruitful.

Second, there may be merit in implementing LLM forecasters that can access the internet for information and autonomously update their forecasts. While the design and implementation of such a setup pose challenges, particularly as most such design choices require human input which may introduce bias, it remains a noteworthy endeavor. Internet access may also help mitigate the design weakness outlined earlier, where the temporal distance between the training data cutoff and the forecasting question's start may impair contextual understanding of the question details.

Third, we outlined a simple Bayesian model-averaging approach for aggregating human and machine forecasts. Future work may delve deeper into alternative ways to combine such forecasts, such as exploring ways of using LLM forecasts in extremizing algorithms or other ways of selectively choosing to favor or disfavor LLM forecasts depending on their eventual performance and comparative strengths and weaknesses. However, it is worth pointing out that given the current poor prediction

performance, such approaches probably only become worthwhile once LLM foresight abilities are significantly improved.

Finally, it would be productive to investigate the integration of human and LLM forecasts, such as by examining whether hybrid forecasting models outperform the standalone baselines. Such hybrid forecasters (i.e., human forecasters that rely upon LLM outputs to make predictions) may be able to combine the respective strengths of both human cognition and LLM cognition, which open up the potential for new forecasting techniques based on the augmentation—rather than replacement—of humans by AI.[3]

Our data suggest a clear weakness in an otherwise impressive catalogue of LLM capabilities: predicting the future. This highlights a number of technical challenges and future research directions for how LLMs may be harnessed for forecasting, as well as for various real-world tasks that require or are helped by robust forecasting capabilities. Our results suggest that the wisdom of human crowds remains a powerful tool in providing probabilistic forecasts about the future. This context of a forecasting tournament may also prove to be an especially fruitful environment for probing the degree to which LLMs are capable of generalized reasoning and prediction.

## Acknowledgements

## Data availability

Our pre-registration and analysis plan can be found in our Open Science Foundation database: `https://osf.io/tfbvd/?view_only=bfc6d61152654e3b84c36446f1858fb2`

## References

[1] Sanjeev Arora and Anirudh Goyal. "A Theory for Emergence of Complex Skills in Language Models". In: *arXiv preprint arXiv:2307.15936* (2023).

[2] Emily M. Bender et al. "On the Dangers of Stochastic Parrots: Can Language Models be too Big?" In: *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*. FAccT '21. Virtual Event, Canada: Association for Computing Machinery, 2021, pp. 610–623. ISBN: 9781450383097. DOI: 10.1145/3442188.3445922.

[3] Daniel M. Benjamin et al. "Hybrid Forecasting of Geopolitical Events". In: *AI Magazine* (2023).

[4] Stella Biderman et al. *Emergent and Predictable Memorization in Large Language Models*. 2023. arXiv: 2304.11158 [cs.CL].

[5] Glenn W Brier. "Verification of Forecasts Expressed in Terms of Probability". In: *Monthly Weather Review* 78.1 (1950), pp. 1–3.

[6] Sébastien Bubeck et al. *Sparks of Artificial General Intelligence: Early Experiments with GPT-4*. 2023. arXiv: 2303.12712 [cs.CL].

[7] David V Budescu and Eva Chen. "Identifying Expertise to Extract the Wisdom of Crowds". In: *Management Science* 61.2 (2015), pp. 267–280.

[8] Nicholas Carlini et al. "Quantifying Memorization Across Neural Language Models". In: *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023. URL: https://openreview.net/pdf?id=TatRHT%5C_1cK.

[9] Emily Dardaman and Abhishek Gupta. "Asking Better Questions: The Art and Science of Forecasting". In: *CHI 2023 Designing Technology and Policy Simultaneously: Towards A Research Agenda and New Practice Workshop*. Hamburg, Germany: ACM, 2023. URL: https://doi.org/....

[10] Clintin P. Davis-Stober et al. "When is a Crowd Wise?" In: *Decision* 1.2 (2014), p. 79.

[11] Tanya Goyal, Junyi Jessy Li, and Greg Durrett. *News Summarization and Evaluation in the Era of GPT-3*. 2023. arXiv: 2209.12356 [cs.CL].

[12] Nathaniel P Grove and Stacey Lowery Bretz. "A Continuum of Learning: From Rote Memorization to Meaningful Learning in Organic Chemistry". In: *Chemistry Education Research and Practice* 13.3 (2012), pp. 201–208.

[13] Nate Gruver et al. *Large Language Models Are Zero-Shot Time Series Forecasters*. 2023. arXiv: 2310.07820 [cs.LG].

[14] Mark Himmelstein, David V. Budescu, and Ying Han. "The Wisdom of Timely Crowds". In: *Judgment in Predictive Analytics*. Springer, 2023, pp. 215–242.

[15] Wenxiang Jiao et al. *Is ChatGPT a Good Translator? Yes with GPT-4 as the Engine*. 2023. arXiv: 2301.08745 [cs.CL].

[16] Christopher W Karvetski et al. "What do Forecasting Rationales Reveal about Thinking Patterns of Top Geopolitical Forecasters?" In: *International Journal of Forecasting* 38.2 (2022), pp. 688–704.

[17] Daniel Martin Katz et al. "GPT-4 Passes the Bar Exam". In: *SSRN* (2023).

[18] Daliang Li et al. "Large Language Models with Controllable Working Memory". In: *Findings of the Association for Computational Linguistics: ACL 2023*. Toronto, Canada: Association for Computational Linguistics, 2023, pp. 1774–1793. DOI: 10.18653/v1/2023.findings-acl.112. URL: https://aclanthology.org/2023.findings-acl.112.

[19] Inbal Magar and Roy Schwartz. "Data Contamination: From Memorization to Exploitation". In: *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*. Dublin, Ireland: Association for Computational Linguistics, May 2022, pp. 157–165. DOI: 10.18653/v1/2022.acl-short.18. URL: https://aclanthology.org/2022.acl-short.18.

[20] Albert E Mannes, Jack B Soll, and Richard P Larrick. "The Wisdom of Select Crowds". In: *Journal of Personality and Social Psychology* 107.2 (2014), p. 276.

[21] Thomas McAndrew et al. "Chimeric Forecasting: Combining Probabilistic Predictions from Computational Models and Human Judgment". In: *BMC Infectious Diseases* 22.1 (2022), p. 833.

[22] Thomas McAndrew et al. "Early Human Judgment Forecasts of Human Monkeypox, May 2022". In: *The Lancet Digital Health* 4.8 (2022), e569–e571.

[23] Barbara Mellers et al. "Identifying and Cultivating Superforecasters as a Method of Improving Probabilistic Predictions". In: *Perspectives on Psychological Science* 10.3 (2015), pp. 267–281.

[24] Metaculus. *Quarterly Cup*. 2023. URL: https://www.metaculus.com/tournament/quarterly-cup-2023q3/.

[25] Humza Naveed et al. *A Comprehensive Overview of Large Language Models*. https://github.com/humza909/LLM_Survey.git. 2023.

[26] Richard Ngo, Lawrence Chan, and Sören Mindermann. *The Alignment Problem from a Deep Learning Perspective*. 2023. arXiv: 2209.00626 [cs.AI].

[27] Harsha Nori et al. *Capabilities of GPT-4 on Medical Challenge Problems*. 2023. arXiv: 2303.13375 [cs.CL].

[28] OpenAI. *GPT-4 Technical Report*. 2023. arXiv: 2303.08774 [cs.CL].

[29] OpenAI. *OpenAI Charter*. OpenAI, 2018. URL: https://openai.com/charter.

[30] Peter S. Park, Philipp Schoenegger, and Chongyang Zhu. *Diminished Diversity-of-Thought in a Standard Large Language Model*. 2023. arXiv: 2302.07267 [cs.HC].

[31] Peter S. Park and Max Tegmark. *Divide-and-Conquer Dynamics in AI-Driven Disempowerment*. 2023. arXiv: 2310.06009 [cs.CY].

[32] Peter S. Park et al. *AI Deception: A Survey of Examples, Risks, and Potential Solutions*. 2023. arXiv: 2308.14752 [cs.CY].

[33] Fotios Petropoulos et al. "Forecasting: Theory and Practice". In: *International Journal of Forecasting* 38.3 (2022), pp. 705–871.

[34] Philip E. Tetlock and Dan Gardner. *Superforecasting: The Art and Science of Prediction*. Random House, 2016.

[35] Philip E. Tetlock, Barbara A Mellers, and J Peter Scoblic. "Bringing Probability Judgments into Policy Debates via Forecasting Tournaments". In: *Science* 355.6324 (2017), pp. 481–483.

[36] Philip E. Tetlock et al. "Forecasting Tournaments: Tools for Increasing Transparency and Improving the Quality of Debate". In: *Current Directions in Psychological Science* 23.4 (2014), pp. 290–295.

[37] Ashish Vaswani et al. "Attention is All You Need". In: *Advances in Neural Information Processing Systems* 30 (2017).

[38] Joost C. F. de Winter. "Can ChatGPT Pass High School Exams on English Language Comprehension?" In: *International Journal of Artificial Intelligence in Education* (2023). ISSN: 1560-4292.

[39] Benfeng Xu et al. *ExpertPrompting: Instructing Large Language Models to be Distinguished Experts*. 2023. arXiv: 2305.14688 [cs.CL].

[40] Daniel M Ziegler et al. "Fine-tuning Language Models from Human Preferences". In: *arXiv preprint arXiv:1909.08593* (2019).

## Appendix A    Forecasting questions

| |
|---|
| Will India's Chandrayaan-3 mission successfully land a rover on the moon? |
| Will Mohamed Bazoum, Nigerien President, return to power before August 31, 2023? |
| Will the House Oversight Committee receive access to requested Joe Biden records related to Hunter Biden's Ukraine dealings before September 1st? |
| Will General Sergei Surovikin be stripped of his command by July 11th? |
| Will Putin attend the G20 summit in India? |
| Will the United Auto Workers call a strike against any of the Big Three Detroit automakers before September 19, 2023? |
| Will Chevron reach an agreement with the Offshore Alliance to end or prevent industrial actions before September 25, 2023? |
| Will a bill be introduced in the Indian Parliament to change the official name of the country to Bharat before September 23, 2023? |
| Will the Black Sea grain deal be revived before October 1, 2023? |
| Will a US warship enter the Black Sea before September 25, 2023? |
| Will the U2 concert at The Sphere on September 29, 2023 take place? |
| Will the WHO name BA.2.86 as a SARS-CoV-2 Variant of Interest before October 1, 2023? |
| Will the extent of Antarctic sea ice for every day in September 2023 be the lowest in recorded history? |
| Will Ukraine regain control of central Bakhmut by the end of September 2023? |
| Will a vote on a Republican-introduced resolution to vacate the Speaker of the House be held before October 1, 2023? |
| Will a non-proprietary LLM be in the top 5 of the chat.lmsys.org leaderboard on September 30, 2023? |
| Will the Israeli High Court issue a ruling on the 'reasonableness' law before October 1, 2023? |
| Will the US pass the Unidentified Anomalous Phenomena (UAP) Disclosure Act of 2023 before October 1, 2023? |
| Will the UN Security Council adopt a resolution related to the Nagorno Karabakh conflict between Azerbaijan and Armenia before October 1, 2023? |
| Will Stepanakert / Khankendi be under de facto Azerbaijani control on September 30, 2023? |
| Before October 1, 2023, will US Senator Bob Menendez announce that he is resigning? |
| Who will be the de facto leader of Gabon on September 30, 2023 (General Brice Oligui Nguema)? |
| Who will be the de facto leader of Gabon on September 30, 2023? (Albert Ondo Ossa) |

**Table 2:** Full list of forecasting questions used in our study.

## Appendix B    Example forecasting question and its corresponding information

**Question**

Will a vote on a Republican-introduced resolution to vacate the Speaker of the House be held before October 1, 2023?

**Resolution Criteria**

This question resolves as Yes if, before October 1, 2023, a member of the Republican Party introduces a resolution to remove Kevin McCarthy as Speaker of the House and a vote on the resolution is held. Both the introduction of the resolution and the vote must occur before October 1, 2023. Otherwise this question resolves as No. The outcome of the vote and any such resolutions introduced by representatives who are not Republicans are irrelevant for the purposes of this question.

**Background Information**

Kevin McCarthy was elected Speaker of the US House of Representatives on January 7, 2023, after 15 ballots, the first time since 1923 an election for Speaker required more than one ballot. The contentious election and concessions to the House Freedom Caucus — including rules that allow any member of the House to call for a vote that would oust the Speaker by simple majority — have weakened his position as Speaker.

The rules for the House of Representatives of the 118th Congress adopt those of the 117th Congress with some amendments. The relevant portion of the rules is Rule IX, the text of which is quoted below. The amended rules remove subparagraph (3) of clause 2(a) (shown in bold).

> 1. Questions of privilege shall be, first, those affecting the rights of the House collectively, its safety, dignity, and the integrity of its proceedings; and second, those affecting the rights, reputation, and conduct of Members, Delegates, or the Resident Commissioner, individually, in their representative capacity only.
>
> 2. (a)(1) A resolution reported as a question of the privileges of the House, or offered from the floor by the Majority Leader or the Minority Leader as a question of the privileges of the House, or offered as privileged under clause 1, section 7, article I of the Constitution, shall have precedence of all other questions except motions to adjourn. A resolution offered from the floor by a Member, Delegate, or Resident Commissioner other than the Majority Leader or the Minority Leader as a question of the privileges of the House shall have precedence of all other questions except motions to adjourn only at a time or place, designated by the Speaker, in the legislative schedule within two legislative days after the day on which the proponent announces to the House an intention to offer the resolution and the form of the resolution. Oral announcement of the form of the resolution may be dispensed with by unanimous consent.
>
> (2) The time allotted for debate on a resolution offered from the floor as a question of the privileges of the House shall be equally divided between (A) the proponent of the resolution, and (B) the Majority Leader, the Minority Leader, or a designee, as determined by the Speaker.
>
> **(3) A resolution causing a vacancy in the Office of Speaker shall not be privileged except if offered by direction of a party caucus or conference.**
>
> (b) A question of personal privilege shall have precedence of all other questions except motions to adjourn.

On September 19, 2023, journalist Matt Laslo claimed to have discovered a draft motion to vacate the office of Speaker of the House in a bathroom in the US Capitol, with Matt Gaetz as the member to submit the resolution.