The Future Outcome Reasoning and Confidence Assessment Benchmark

Anonymous ACL submission

Abstract

Forecasting is an important task in many domains, such as technology and economics. 002 However existing forecasting benchmarks 003 004 largely lack comprehensive confidence assessment, focus on limited question types, 005 and often consist of artificial questions that 006 do not align with real-world human forecasting needs. To address these gaps, we 009 introduce FORECAST (Future Outcome 010 Reasoning and Confidence Assessment), a benchmark that evaluates models' ability to 011 make predictions and their confidence in 012 them. FORECAST spans diverse forecast-013 ing scenarios involving Boolean questions, 014 015 timeframe prediction, and quantity estimation, enabling a comprehensive evaluation 016 of both prediction accuracy and confidence 017 calibration for real-world applications. 018

1 Introduction

019

Recent advances in large language models (LLMs) have significantly improved their perfor-021 mance across a wide range of natural language processing (NLP) tasks. Alongside these devel-023 opments, various benchmarks and datasets have 024 025 been introduced to effectively assess the capabilities of LLMs, particularly in terms of knowl-026 edge and reasoning (Zellers et al., 2019; Guo 027 et al., 2023). Fact-based benchmarks, like TruthfulQA (Lin et al., 2022) evaluate LLMs based

on factual correctness, focusing on tasks like retrieving and verifying facts that are known. 030

031

053

054

055

056

057

Forecasting is a crucial yet challenging task across various domains, including technology, 033 economics, and public policy. Unlike tasks that 034 rely on retrieving and verifying existing knowl-035 edge, forecasting requires predicting plausible 036 outcomes for future events, often under uncer-037 tainty and incomplete information. This makes 038 forecasting particularly difficult, as models must 039 infer trends, assess probabilities, and adapt to new information. Several datasets have been 041 introduced to evaluate LLMs' forecasting ca-042 pabilities. ForecastQA (Jin et al., 2020) uses 043 a multiple-choice format where models predict future outcomes, but it lacks confidence assess-045 ment. AutoCast (Zou et al., 2022) incorporates confidence intervals, however its confidence es-047 timates are not designed for forecasting. Other 048 datasets such as ExpTime (Yuan et al., 2024a) are artificially generated from structured data, focusing on explainable event forecasting based 051 on temporal knowledge graphs.

All of these aforementioned benchmarks ignore a crucial aspect of forecasting: confidence evaluation. Confidence plays a central role in forecasting, as predictions about unresolved events inherently lack definitive correctness at the time of evaluation. Predictions made with absolute certainty are undesirable, even if they

Туре	Question	Resolution	Confidence
Boolean Question	Will a Frontier AI lab be established in China by 2026?	Yes	0.73
Timeframe Prediction	When will OpenAI announce GPT-5?	2024-08-01	0.85
Quantity Estimation	How many spacecrafts will land on the moon in 2025?	3	0.65

Table 1: Examples of forecasting questions with their resolutions and confidence scores.

ultimately prove to be correct, because they fail 060 to account for the uncertain nature of future 061 events. Moreover, miscalibrated confidence can 062 lead to poor decision-making: overconfident yet 063 incorrect forecasts may result in costly errors, 064 while underconfident but accurate predictions 065 can erode trust in the model. Therefore, well-066 calibrated confidence scores are as crucial as the 067 accuracy of the predictions themselves. 068

To address these gaps, we present FORE-069 CAST: Future Outcome Reasoning and Confidence Assessment. FORECAST focuses 071 on three distinct types of forecasting questions, 072 shown in Table 1: (1) Boolean questions, such as "Will there be a Frontier AI lab in China be-074 fore 2026?"; (2) Timeframe Prediction, such as "When will OpenAI announce GPT-5?"; and (3) Quantity Estimation, such as "How many space-077 crafts will land on the moon in each of the fol-078 lowing years?" We conduct experiments using 079 a range of models differing in size, training objectives, and cutoff times, and explore multiple 081 082 methods for estimating model confidence. Our results reveal that forecasting remains highly 083 challenging for current LLMs, particularly in confidence evaluation, with no direct correlation 085 between prediction performance and confidence 086 calibration, and while larger models sometimes 087 improve performance, the effect is inconsistent.

089

090

2 FORECAST: Problem Formulation

System responses in FORECAST consist of (1) a prediction answering a question given the

available information and (2) a confidence score in the prediction. This ensures a comprehensive assessment of forecasting, accounting for both correctness and confidence calibration. Questions belong to three types. (1) Boolean Ques-096 tions, which ask yes/no questions about the occurrence of future events (sometimes within a certain timeframe). Boolean questions are simple to evaluate, and they can still be surprisingly 100 challenging (Clark et al., 2019). (2) Timeframe 101 Prediction, which requires predicting a specific timeframe for an event, and are essential for 103 applications where knowing whether an event 104 will happen or not without a timeframe is insuf-105 ficient. (3) Quantity Estimation, which involves providing numerical estimates related to future 107 events, e.g., economic indicators or trends.

Formally, let Q represent a question about a future event, and let M denote a system with access to information up to time t_{train} (e.g., the system's knowledge cutoff point). The objective is for M to produce an answer A in A and an associated confidence score C, where:

$$M(Q) \to (A, C), \quad C \in [0, 1].$$
 (1)

109

110

111

112

113

114

115

116

117

$$A = \arg\max_{a \in \mathcal{A}} P(X = a | Q, \mathcal{K}(t_{train})).$$
 (2)

Here, $\mathcal{K}(t_{train})$ represents the knowledge accessible to the model up to time t_{train} . The answer118sible to the model up to time t_{train} . The answer119space \mathcal{A} depends on the type of forecasting question: for Boolean Questions, $\mathcal{A} = \{\text{Yes}, \text{No}\};$ 121for Timeframe Prediction, \mathcal{A} consists of a single122date in the YYYY-MM-DD format; and for Quantity123Estimation, $\mathcal{A} = \mathbb{R}$, representing real numbers.124

3 Evaluating Predictions and Confidence

125

126

139

154

Evaluating Boolean Questions For Boolean 127 questions, where the answer space is A =128 {Yes, No}, prediction performance is evaluated 129 using standard classification metrics, including accuracy and F1-score. Confidence calibration 131 is assessed using a modified version of the Brier 132 score (Brier, 1950), which measures the mean 133 squared error between predicted confidence and 134 the gold confidence provided in the dataset, 135 which we assume is provided, and represents 136 the likelihood of an event occurring. The modi-137 fied Brier score is defined as: 138

Brier =
$$\frac{1}{N} \sum_{i=1}^{N} (C_i^{\text{pred}} - C_i^{\text{gold}})^2$$
, (3)

140 where C_i^{pred} is the model's predicted confi-141 dence and C_i^{gold} is the gold confidence. This 142 modification ensures that models are evaluated 143 based on their ability to match the likelihood of 144 an event. Lower Brier scores indicate better cal-145 ibration, reflecting how well the predicted confi-146 dence aligns with the likelihood of the event.

147Evaluating Timeframe PredictionFor time-148frame prediction, where the answer space con-149sists of specific dates in the YYYY-MM-DD format,150predictive accuracy is measured using absolute151day error (ADE). Given a predicted date D_i^{pred} 152and the gold date D_i^{gold} , we compute the normal-153ized error as:

$$E_i^{\text{ADE}} = \frac{2}{1 + e^{-\alpha |D_i^{\text{pred}} - D_i^{\text{gold}}|}} - 1, \quad (4)$$

155where α is a scaling factor that controls how156sharply large errors are penalized. This transfor-157mation ensures that extreme deviations do not158disproportionately dominate the evaluation.

For confidence calibration, rely on the Contin-159 uous Ranked Probability Score (CRPS) (Matheson and Winkler, 1976), which is a generalisa-161 tion of the Mean Absolute Error to probabilistic 162 forecasts, and extend it to compare the predicted probability distribution with a gold distribution. 164 Specifically, we assume that both the predicted 165 and the gold confidence predictions follow Gaus-166 sian distributions, namely $\mathcal{N}(D_i^{\text{pred}}, \sigma_i^{\text{pred}})$ and 167 $\mathcal{N}(D_i^{\mathrm{gold}},\sigma_i^{\mathrm{gold}})$ respectively, where the stan-168 dard deviations are computed as:

$$\sigma_i^{\text{pred}} = \sigma_{\max} \cdot (1 - C_i^{\text{pred}}) + \sigma_{\min} \cdot C_i^{\text{pred}}, \ (5)$$

170

$$\sigma_i^{\text{gold}} = \sigma_{\max} \cdot (1 - C_i^{\text{gold}}) + \sigma_{\min} \cdot C_i^{\text{gold}}.$$
 (6)

Here, C_i^{pred} is the model's predicted confidence for the *i*th question, and C_i^{gold} is the correspond-173 174 ing gold confidence provided in our dataset. The 175 parameters $\sigma_{\rm max}$ and $\sigma_{\rm min}$ define the upper and 176 lower bounds for the standard deviation. Intu-177 itively, when confidence is low ($C \approx 0$), un-178 certainty is high, leading to $\sigma \approx \sigma_{\text{max}}$, while 179 when confidence is high ($C \approx 1$), uncertainty is 180 low, resulting in $\sigma \approx \sigma_{\min}$. We then compute 181 the CRPS as the integrated squared difference between the cumulative distribution functions (CDFs) of the predicted and gold distributions:

$$CRPS = \frac{1}{N} \sum_{i=1}^{N} \int \left(F_i^{\text{pred}}(d) - F_i^{\text{gold}}(d) \right)^2 \, \mathrm{d}d,$$
(7)

where F_i^{pred} and F_i^{goid} denote the CDFs of the186predicted and gold Gaussian distributions, re-187spectively. A lower CRPS indicates better cal-188ibration, as it reflects a closer match between189the predicted uncertainty and the uncertainty as190specified by the gold confidence.191

Evaluating Quantity EstimationFor quan-tity estimation, where the answer space consists193of non-negative real numbers ($\mathcal{A} = \mathbb{R}_{\geq 0}$), we194

195 196

197

- 198 199
- 200
- 201
- 202

215

216

217

218

219

220

221

222

223

224

225

226

227

evaluate prediction performance using two error metrics: absolute percentage error (APE) and mean absolute error (MAE). Given a predicted quantity Q_i^{pred} and the gold quantity Q_i^{gold} , the normalized errors are computed as:

 $E_i^{\text{APE}} = \frac{2}{1+e^{-\alpha \frac{|Q_i^{\text{pred}} - Q_i^{\text{gold}}|}{Q_i^{\text{gold}} + \epsilon}}} - 1, \quad (8)$

$$E_i^{\text{MAE}} = \frac{2}{1 + e^{-\alpha |Q_i^{\text{pred}} - Q_i^{\text{gold}}|}} - 1.$$
(9)

203 Here, ϵ is a small constant to prevent division by zero, and α is a scaling factor that controls how 204 sharply large errors are penalized, similar to the 205 timeframe prediction evaluation. Confidence 206 calibration is assessed using CRPS, following 207 the same Gaussian assumption as in timeframe 208 prediction. The predicted quantity is modeled as 209 a Gaussian distribution $\mathcal{N}(Q_i^{\text{pred}}, \sigma_i^{\text{pred}})$, and the gold quantity as $\mathcal{N}(Q_i^{\text{gold}}, \sigma_i^{\text{gold}})$. The standard deviations σ_i^{pred} and σ_i^{gold} are computed using 210 211 212 the same formulation as in timeframe prediction. 213

214 4 FORECAST Construction

4.1 Data Source and Question Selection

FORECAST is constructed from Metaculus,¹ an online forecasting platform where forecasters submit probabilistic predictions to questions across various domains. Metaculus aggregates individual forecasts into a continuously updated community prediction, which is finalized just before resolution. Each question has predefined resolution criteria, ensuring verifiable outcomes. To ensure dataset reliability, we include only questions with a definitive resolution and at least 100 forecasts to maintain statistical reliability. Ambiguous or subjectively resolved questions

¹www.metaculus.com. Examples in Appendix A.

are excluded, and we remove those whose out-
comes depend on arbitrary or uncontrollable fac-
tors. These steps ensure that FORECAST con-
sists of high-quality, well-formed forecasting
questions with verifiable outcomes.228229231231231232232

233

234

4.2 Extracting Confidence from Crowdsourced Forecasts

Forecasting distinguishes between physical 235 probabilities-objective likelihoods derived 236 from statistical or scientific models-and hu-237 man beliefs (Sanders, 1963) about future events. 238 While physical probabilities can be useful, they 239 are often unavailable, particularly for questions 240 involving human behavior, economics, or so-241 ciopolitical outcomes. Instead, collective hu-242 man forecasts offer a more practical confidence 243 estimate, integrating expert reasoning, contex-244 tual knowledge, and evolving evidence. For in-245 stance, predicting a technological breakthrough 246 depends more on expert assessment and current 247 trends than on rigid probabilistic models. There-248 fore, confidence in FORECAST is derived from 249 Metaculus community forecasts, which aggregate predictions from a diverse pool of forecast-251 ers. While human predictions are sometimes incorrect, they still serve as a valuable proxy 253 for uncertainty, as they reflect the best available 254 reasoning given the information at the time. 255

Formal Definition of Gold Confidence. Gold 256 confidence in FORECAST is derived from the fi-257 nal Metaculus community prediction before res-258 olution. Instead of directly using the predicted 259 probability for the correct outcome, we compute a log score relative to a uniform baseline, 261 ensuring that confidence reflects how much the 262 forecast deviates from random guessing. This 263 transformation prevents extreme probabilities in 264 inherently uncertain scenarios and makes confi-265 dence scores more comparable across different 266 forecasting tasks. The final score is mapped to (0, 1) using a sigmoid function.

267

268

269

270

271

273

274

276

277

278

279

280

281

282

283

284

285

286

287

289

290

291

292

293

294

295

296

297

298

299

300

301

302

For Boolean questions, where the humanforecasted probability for the correct outcome is P^{gold} , gold confidence is computed as:

272
$$C^{\text{gold}} = \sigma \left(\frac{\ln P^{\text{gold}} - \ln 0.5}{\ln 2} \right), \quad (10)$$

where $\sigma(x)$ is the sigmoid function.

Similarly, for timeframe prediction and quantity estimation, where the human-forecasted probability density function (PDF) assigns probability to a continuous outcome x^{gold} , gold confidence is computed as:

$$C^{\text{gold}} = \sigma\left(\frac{\ln f(x^{\text{gold}}) - \ln f_{\text{uniform}}}{2}\right), \quad (11)$$

where $f(x^{\text{gold}})$ is the forecasted probability density at the resolved outcome, and f_{uniform} is the uniform baseline density over the valid range of values. The denominator 2 ensures numerical stability and scales confidence appropriately.

4.3 Dataset Statistics and Comparison

FORECAST consists of 2256 forecasting questions, spanning domains such as politics, economics, science, and technology. Each question includes a resolved outcome, a gold confidence score, and a final Metaculus community forecast before resolution. To facilitate model development and evaluation, we split the dataset into 65% training, 10% validation, and 25% test. The full dataset statistics is shown in Appendix B.

Table 2 provides a comparison between FORECAST and existing forecasting benchmarks. Compared to prior datasets, FORECAST uniquely emphasizes both forecasting accuracy and confidence calibration, includes a diverse set of forecasting tasks, and is constructed from a well-established crowdsourced platform with rigorous resolution criteria.

Benchmark	Question Types	Natural Questions	Confidence
ForecastQA	MCQ	1	X
AutoCast	Various	1	X
ExpTime	Boolean	×	X
FORECAST	Various	1	1

Table 2: Comparison of key features across our benchmark variants, highlighting our evaluation of confidence across different question types.

5 Experiments on FORECAST

5.1 Experimental Setup

Models.We evaluate a diverse set of large305language models (LLMs) with varying training306data cutoffs, model sizes, and instruction tuning.307To analyze the impact of knowledge recency, we308group models by family and assume the cutoff309date is the 1st of the stated month. The models310used are shown in Table 3.311

303

304

Model Family	Model Variants	Cutoff Date
GPT-2	GPT-2, GPT-2 XL	2017-12-01
Pythia	14M, 160M, 2.8B	2020-03-01
BLOOM	560M, 7B1	2021-12-01
LLaMA	LLaMA-7B	2022-08-01
OLMo	1B, 7B, 7B-Instruct	2023-03-01
OLMo-2	7B, 7B-Instruct	2023-12-01

Table 3: Models used in our experiments, grouped by family and ordered by training data cutoff date, including: GPT-2 (Radford et al., 2019), Pythia (Biderman et al., 2023), BLOOM (Scao et al., 2023), LLaMA (Touvron et al., 2023), OLMo (Groeneveld et al., 2024), and OLMo-2 (OLMo et al., 2024).

Inference. We use 1-shot in-context learning 312 to provide models with a structured example 313 of how to answer forecasting questions. For 314 instruction-tuned models, we add an extra line of 315 instruction to align with their training paradigm. 316 To ensure fair comparison, baseline prompts are 317 kept minimal while maintaining clarity. Confi-318 dence is estimated using logit-based normalized 319

heuristics.	Full prompts and hyperparameters	3
are provide	d in Appendix C and Appendix D.	

Model	Ν	Accuracy (\uparrow)	F1 (†)	Brier (\downarrow)
GPT2	401	0.5835	0.3650	0.4199
GPT2-XL	401	0.6708	0.3158	0.4540
Pythia-14m	343	0.5860	0.1374	0.6245
Pythia-160m	343	0.6093	0.2584	0.5273
Pythia-2.8b	343	0.5452	0.3906	0.4426
Bloom-560m	263	0.4905	0.4071	0.4521
Bloom-7b1	263	0.6426	0.3286	0.3200
Llama-7b	226	0.5708	0.3433	0.5201
OLMo-1B	188	0.2340	0.1591	0.7545
OLMo-7B	188	0.2074	0.1687	0.8199
OLMo-7B-Instruct	188	0.6543	0.1408	0.3811
OLMo-2-7B	145	0.5103	0.2444	0.4377
OLMo-2-7B-Instruct	145	0.5862	0.3023	0.4078

Table 4: Performance of forecasting models on Boolean questions in FORECAST. Reported metrics include the number of evaluated questions (N), accuracy, F1 score (both higher is better), and the Brier score (lower is better).

5.2 **Results and Findings**

Our experiments on the FORECAST dataset reveal that forecasting remains highly challeng-324 ing for current LLMs, particularly in confidence estimation. While models achieve reason-326 able accuracy in point predictions, their uncertainty estimates-captured by Brier score and CRPS—vary significantly. This suggests that confidence evaluation must be treated as a separate challenge from prediction assessment.

Boolean Ouestions: Table 4 shows that while some models achieve reasonable accuracy, their 333 calibration remains inconsistent. For example, GPT2-XL achieves an accuracy of 0.67 but 335 has a relatively high Brier score of 0.45, while 336 337 OLMo-7B-Instruct achieves similar accuracy (0.65) with a lower Brier score of 0.38, sug-338 gesting better confidence estimation. Notably, 339 models with later training cutoffs do not always 340 outperform older ones; for instance, LLaMA-7B 341

has lower accuracy (0.57) and a worse Brier score (0.52) than some earlier models. This indicates that forecasting accuracy depends on multiple factors beyond knowledge recency.

342

343

344

345

359

Model	Ν	ADE (\downarrow)	CRPS (\downarrow)
GPT2	26	0.9944	0.9884
GPT2-XL	26	1.0000	1.0000
Pythia-14m	25	1.0000	1.0000
Pythia-160m	25	1.0000	1.0000
Pythia-2.8b	25	0.9650	0.9634
Bloom-560m	12	1.0000	1.0000
Bloom-7b1	12	0.9984	0.9976
Llama-7b	10	0.9843	0.9769
OLMo-1B	6	1.0000	1.0000
OLMo-7B	6	1.0000	1.0000
OLMo-7B-Instruct	6	1.0000	1.0000
OLMo-2-7B	4	0.9981	0.9970
OLMo-2-7B-Instruct	4	0.8696	0.8197

Table 5: Performance of forecasting models on Timeframe Prediction tasks in FORECAST. Metrics include the number of evaluated questions (N), the normalized absolute days error (ADE), and the Continuous Ranked Probability Score (CRPS), where lower values indicate better performance.

Timeframe Prediction: Table 5 reveals that 346 most models struggle with predicting event 347 timing, as both ADE and CRPS remain near 348 the worst-case scenario of 1.0. Across all 349 models, only one achieves a CRPS below 0.9: OLMo-2-7B-Instruct (CRPS = 0.82), which 351 still indicates substantial uncertainty. Even models with more recent training data, such as OLMo-2-7B (CRPS = 0.99), fail to make well-354 calibrated temporal forecasts. These results suggest that even when models correctly anticipate whether an event will happen, quantifying when 357 it will occur remains a major challenge.

Quantity Estimation: Table 6 highlights that models differ widely in their ability to esti-

320

323

325

327

328

329

330

331

Model	Ν	APE (\downarrow)	MAE (\downarrow)	CRPS (\downarrow)
GPT2	81	0.2296	0.8697	0.8575
GPT2-XL	81	0.0295	0.7461	0.7169
Pythia-14m	77	0.1696	0.8753	0.8538
Pythia-160m	77	0.0274	0.7800	0.7620
Pythia-2.8b	77	0.0742	0.8106	0.7851
Bloom-560m	43	0.0462	0.7273	0.7082
Bloom-7b1	43	0.0644	0.7461	0.7176
Llama-7b	32	0.0564	0.6422	0.6126
OLMo-1B	23	0.2124	0.7883	0.7623
OLMo-7B	23	0.2157	0.8365	0.8201
OLMo-7B-Instruct	23	0.2358	0.7987	0.7703
OLMo-2-7B	20	0.0206	0.6457	0.6206
OLMo-2-7B-Instruct	20	0.0872	0.5968	0.5686

Table 6: Performance of forecasting models on Quantity Estimation tasks in FORECAST. Metrics include the number of evaluated questions (N), normalized absolute percentage error (APE), mean absolute error (MAE), and Continuous Ranked Probability Score (CRPS), where lower values are better.

mate numerical values. Some achieve relatively 361 low APE, such as GPT2-XL (APE = 0.03) and 362 Pythia-160M (APE = 0.02), yet their confi-363 dence calibration, CRPS, does not always align 364 with their point prediction performance. For 365 instance, GPT2-XL has a CRPS of 0.72, while 366 Pythia-160M has a slightly higher CRPS of 367 0.76, despite achieving lower APE. In contrast, OLMo-2-7B-Instruct, which has a relatively 369 higher APE of 0.08, achieves the lowest CRPS 370 (0.57) among all models. These results indicate 371 that prediction quality and confidence calibra-372 tion do not necessarily improve together, rein-373 forcing the complexity of numerical forecasting.

Impact of Model Size on Forecasting Per-375 formance Larger models do not consistently 376 improve forecasting performance. Within 377 378 the Pythia family, Pythia-2.8b shows occasional gains in accuracy and calibration over 379 Pythia-14m and Pythia-160m, but the im-380 provements are not uniform across all metrics. 381 Similarly, while OLMo-2-7B variants sometimes

achieve lower CRPS and MAE in quantity es-
timation, these benefits often come with trade-
offs in point prediction. These results suggest
that model size alone is not a reliable predic-
tor of forecasting performance. This highlights
the importance of developing task-specific tech-
niques rather than relying on ever-larger models
to solve the forecasting problem.383
383
384

391

392

393

394

395

396

397

398

400

401

402

403

404

405

406

407

408

409

410

411

Impact of Instruction Tuning on Forecasting Performance Table 7 compares base and instruct-tuned variants of OLMo-7B and OLMo-2-7B across Boolean, timeframe, and quantity forecasting tasks. For Boolean questions, OLMo-7B-Instruct achieves higher accuracy (0.65 vs. 0.21) and a lower Brier score (0.38 vs. 0.82), indicating better confidence calibra-In timeframe prediction, the instructtion. tuned OLMo-2-7B-Instruct improves uncertainty estimation, with an ADE of 0.87 and CRPS of 0.82, compared to 0.9981 and 0.9970 for the base OLMo-2-7B. For quantity estimation, instruct-tuned models have slightly higher APE but lower MAE and CRPS, suggesting better uncertainty calibration. These results indicate that instruction tuning enhances confidence estimation, even if it does not always improve point prediction accuracy. This suggests a trade-off where instruct-tuned models prioritize more reliable uncertainty quantification.

Impact of Aggregation Methods on Forecast-412 ing Performance Table 8 compares different 413 aggregation methods for deriving the final pre-414 diction and confidence estimate from the top 415 10 outputs of Llama-7B. Bayesian Aggregation 416 achieves the highest accuracy (0.5796) and F1 417 score (0.3485), suggesting it is the most effec-418 tive at identifying correct outcomes. However, 419 Weighted Average yields a significantly lower 420 Brier score (0.2914), indicating superior confi-421 dence calibration compared to other methods. 422

Model	Accuracy (†)	F1 (†)	Brier (\downarrow)	ADE (\downarrow)	CRPS (T) (\downarrow)	APE (\downarrow)	$\mathrm{MAE}\left(\downarrow\right)$	$CRPS\left(Q\right)\left(\downarrow\right)$
OLMo-7B	0.2074	0.1687	0.8199	1.0000	1.0000	0.2157	0.8365	0.8201
OLMo-7B-Instruct	0.6543	0.1408	0.3811	1.0000	1.0000	0.2358	0.7987	0.7703
OLMo-2-7B	0.5103	0.2444	0.4377	0.9981	0.9970	0.0206	0.6457	0.6206
OLMo-2-7B-Instruct	0.5862	0.3023	0.4078	0.8696	0.8197	0.0872	0.5968	0.5686

Table 7: Effect of instruction tuning on forecasting performance across boolean, timeframe, and quantity estimation questions. Metrics include accuracy, F1, and Brier score for binary questions; ADE and CRPS (T) for timeframe prediction; and APE, MAE, and CRPS (Q) for quantity estimation.

Aggregation Method	Ν	Accuracy (†)	F1 (†)	Brier (\downarrow)
Majority Vote	226	0.5575	0.2353	0.5287
Highest Confidence	226	0.5708	0.3433	0.5201
Weighted Average	226	0.5575	0.2353	0.2914
Logit Mean Probability	226	0.5575	0.2353	0.6195
Bayesian Aggregation	226	0.5796	0.3485	0.5493

Table 8: Ablation study on different aggregation methods for extracting predictions and confidence from Llama-7B outputs.

Majority Vote, Highest Confidence, and Logit 423 Mean Probability produce comparable accuracy 424 and F1 scores but have noticeably higher Brier 425 scores, suggesting weaker uncertainty estima-426 tion. These results highlight that even when 427 point prediction performance is similar, aggre-428 gation methods substantially impact confidence 429 reliability. The challenge remains in develop-430 ing techniques that optimize both accuracy and 431 calibration simultaneously, emphasizing the im-432 portance of uncertainty-aware forecasting. 433

6 Related Work

434

Recent forecasting benchmarks focus on event 435 prediction but largely overlook confidence cal-436 ibration. OpenForecast (Wang et al., 2025) in-437 troduces a large-scale dataset for open-ended, 438 multi-step event forecasting but does not as-439 440 sess model confidence. ForecastBench (Karger et al., 2024) evaluates binary (Yes/No) forecast-441 ing by prompting models for direct probability 442 estimates, but since it queries each option sepa-443 rately, the assigned probabilities do not necessar-444

ily sum to 1, leading to potential inconsistencies. Neither benchmark systematically evaluates confidence calibration, a crucial aspect for reliable forecasting in real-world applications. 445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

Beyond forecasting, several benchmarks assess language models' reasoning and inference capabilities. COPA (Roemmele et al., 2011) evaluates causal reasoning by presenting a premise and two alternatives, requiring models to select the more plausible cause or effect. HellaSwag (Zellers et al., 2019) challenges models with sentence completion tasks that demand commonsense reasoning, where models must choose the most sensible continuation of a given scenario. PRobELM (Yuan et al., 2024b) assesses models' capacity to rank scenarios by plausibility, bridging the gap between factual accuracy and world knowledge. While these benchmarks provide insights into models' reasoning abilities, they do not address the challenges of forecasting future events.

7 Conclusion

We introduce FORECAST, a benchmark for 467 evaluating both forecasting accuracy and con-468 fidence calibration in language models. Un-469 like existing datasets, FORECAST explicitly 470 assesses confidence alongside predictions. Our 471 results show that current models struggle with 472 both prediction and well-calibrated confidence, 473 underscoring the need for improved uncertainty 474 estimation and confidence calibration. 475

476	Limitations	Language Models Across Training and Scaling.	514
477	While FORECAST represents a significant step	ing, pages 2397–2430, PMLR.	515 516
478	toward evaluating forecasting accuracy and con-		
479	fidence calibration, there are inherent limita-	Glenn W Brier. 1950. Verification of forecasts ex-	517
480	tions that we view as opportunities for future	pressed in terms of probability. <i>Monthly weather</i>	518
481	research rather than fundamental shortcomings	review, $78(1):1-3$.	519
482	of our work. First, our dataset is constructed	Christopher Clark, Kenton Lee, Ming-Wei Chang,	520
483	solely from Metaculus, which may not fully rep-	Tom Kwiatkowski, Michael Collins, and Kristina	521
484	resent the global diversity of forecasting prac-	Toutanova. 2019. Boolq: Exploring the surpris-	522
485	tices or question domains. Second, our method	ing difficulty of natural yes/no questions. arXiv	523
486	for deriving gold confidence relies on commu-	preprint arXiv:1905.10044.	524
487	nity forecasts and heuristic transformations that	Dirk Groeneveld, Iz Beltagy, Pete Walsh, Ak-	525
488	might not capture all nuances of human uncer-	shita Bhagia, Rodney Kinney, Oyvind Tafjord,	526
489	tainty. Lastly, our focus on English-language	Ananya Harsh Jha, Hamish Ivison, Ian Magnus-	527
490	forecasts limits the benchmark's applicability	son, Yizhong Wang, et al. 2024. OLMo: Accel-	528
491	across different languages and cultural contexts.	preprint arXiv:2402.00838.	529 530
492	Addressing these issues is part of our future	$r \cdot r \cdot r$	
493	work agenda	Zishan Guo, Renren Jin, Chuang Liu, Yufei Huang,	531
-55	work ugendu.	Dan Shi, Linhao Yu, Yan Liu, Jiaxuan Li, Bojian	532
494	Ethical Statement	Xiong, Deyi Xiong, et al. 2023. Evaluating large	533
		preprint arXiv:2310.19736.	534 535
495	FORECAST is built from data sourced exclu-	1 1	
496	sively from Metaculus, an English-language	Woojeong Jin, Rahul Khanna, Suji Kim, Dong-Ho	536
497	forecasting platform. As a result, the dataset	Lee, Fred Morstatter, Aram Galstyan, and Xiang	537
498	may embody the linguistic, cultural, and social	Ren. 2020. Forecastqa: A question answering	538
499	biases inherent in its user community. These	data arXiv preprint arXiv:2005.00792	539
500	biases could affect both question selection and	uata. urxiv preprint urxiv.2005.00792.	540
501	confidence judgments. We acknowledge these	Ezra Karger, Houtan Bastani, Chen Yueh-Han,	541
502	concerns and stress that our benchmark is in-	Zachary Jacobs, Danny Halawi, Fred Zhang, and	542
503	tended as an initial step toward more inclusive	Philip E Tetlock. 2024. Forecastbench: A dy-	543
504	forecasting evaluations. Future efforts should	namic benchmark of ai forecasting capabilities.	544

aim to incorporate data from a broader range

of platforms and languages to mitigate these

Stella Biderman, Hailey Schoelkopf, Quentin Gre-

gory Anthony, Herbie Bradley, Kyle O'Brien,

Eric Hallahan, Mohammad Aflah Khan, Shivan-

shu Purohit, USVSN Sai Prashanth, Edward Raff,

et al. 2023. Pythia: A Suite for Analyzing Large

505

506

507

508

509

510

511

512

513

biases.

References

Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. TruthfulQA: Measuring How Models Mimic Human Falsehoods. In Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 3214–3252, Dublin, Ireland. Association for Computational Linguistics.

545

546

547

548

549 550

551

552

arXiv preprint arXiv:2409.19839.

James E. Matheson and Robert L. Winkler. 1976. 553 Scoring rules for continuous probability distribu-554 tions. *Management Science*, 22(10):1087–1096. 555 Team OLMo, Pete Walsh, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Shane Arora, Akshita Bhagia, Yuling Gu, Shengyi Huang, Matt Jordan, et al. 2024. 2 olmo 2 furious. *arXiv preprint arXiv:2501.00656*.

556

557

558

559

560

561

562

563

564

565

566

567

568 569

570

571

572 573

574

575

576

577

578 579

580

581

582

583 584

585

586

587

588

589

590

591

592

593 594

595

597 598

- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language Models are Unsupervised Multitask Learners. *OpenAI blog*, 1(8):9.
- Melissa Roemmele, Cosmin Adrian Bejan, and Andrew S Gordon. 2011. Choice of plausible alternatives: An evaluation of commonsense causal reasoning. In 2011 AAAI spring symposium series.
 - Frederick Sanders. 1963. On subjective probability forecasting. *Journal of Applied Meteorology and Climatology*, 2(2):191–201.
- Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, et al. 2023. BLOOM: A 176B-Parameter Open-Access Multilingual Language Model. Working paper or preprint.
 - Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and Efficient Foundation Language Models. *arXiv preprint arXiv:2302.13971*.
 - Zhen Wang, Xi Zhou, Yating Yang, Bo Ma, Lei Wang, Rui Dong, and Azmat Anwar. 2025. Openforecast: A large-scale open-ended event forecasting dataset. In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 5273–5294.
 - Chenhan Yuan, Qianqian Xie, Jimin Huang, and Sophia Ananiadou. 2024a. Back to the future: Towards explainable temporal reasoning with large language models. In *Proceedings of the ACM on Web Conference 2024*, pages 1963–1974.
- Moy Yuan, Eric Chamoun, Rami Aly, Chenxi Whitehouse, and Andreas Vlachos. 2024b. PRobELM: Plausibility ranking evaluation for language models. In *First Conference on Language Modeling*.

- Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali600Farhadi, and Yejin Choi. 2019. HellaSwag: Can601a Machine Really Finish Your Sentence? In Pro-602ceedings of the 57th Annual Meeting of the As-603sociation for Computational Linguistics, pages6044791–4800, Florence, Italy. Association for Com-605putational Linguistics.606
- Andy Zou, Tristan Xiao, Ryan Jia, Joe Kwon, Man-
tas Mazeika, Richard Li, Dawn Song, Jacob Stein-
hardt, Owain Evans, and Dan Hendrycks. 2022.608Forecasting future world events with neural net-
works. Advances in Neural Information Process-
ing Systems, 35:27293–27305.611

613 **A**

618

619

620

621

622

623

624

A Example Metaculus Questions

To illustrate how human forecasts evolve over time, we present two questions from different domains: Q1 is in the business and geopolitics domain and Q2 is in the technology domain.

Q1: Will TikTok become available in the US on both the App Store and Google Play before April 5, 2025?

Q2: When will a SpaceX Starship reach orbit?

For Q1, Figure 1 shows how community forecasts changed over time, while Figure 2 presents the histogram of the final forecast distribution.



Figure 1: Community prediction trend for a Metaculus question on TikTok's availability in the US.



Figure 2: Histogram of final community forecasts.

For Q2, Figure 3 tracks forecast updates, while Figure 4 shows the final probability density function (PDF) of predicted launch dates.



Figure 3: Community prediction trend for SpaceX Starship's first orbital launch.



Figure 4: Probability density function of final community forecasts for SpaceX Starship reaching orbit.

Split	Boolean	Timeframe	Quantity	Total
Training	1142	90	223	1465
Validation	175	13	35	223
Test	441	36	91	568
Total	1758	139	349	2256

Table 9: Dataset statistics for FORECAST, showing the distribution of questions across different forecasting types, with the overall total in the last column.

B Dataset Statistics

Table 9 presents detailed dataset statistics, in-cluding the total number of questions and theirdistribution across Boolean Questions, Time-frame Prediction, and Quantity Estimation tasks.

628

629

630

631

632

633

C Prompts

To ensure a fair and consistent evaluation 634 across models, we use simple one-shot prompts 635 with structured outputs in JSON format. For 636 instruction-tuned models, we provide an addi-637 tional instruction line specifying the task. The 638 prompts are designed for three forecasting ques-639 tion types: Boolean Questions (Yes/No), Time-640 frame Prediction (YYYY-MM-DD), and Quan-641 tity Estimation (numeric values). 642



643	C.1 Instruction-Tuned Models
644	For models containing "Instruct" in their name, we use the following prompts:
645	Quantity Estimation
646	You are an AI assistant providing precise numerical forecasts.
647	Answer the following question with a single numeric value in JSON format.
648	
649	Example:
650	Q: How much global photovoltaic energy generation was deployed
651	by the end of 2020?
652	A: { "value": 738 }
653	
654	Q: \$question
655	A: { "value": "
656	Timeframe Prediction
657	You are an AI assistant providing precise date forecasts.
658	Answer the following question with a single date in YYYY-MM-DD format in JSON.
659	
660	Example:
661	Q: When did an AI system achieve a significant victory against
662	a professional human in Starcraft 2?
663	A: { "value": "2019-01-24" }
664	
665	Q: \$question
000	
667	Boolean Questions
668	You are an AI assistant providing binary (Yes/No) answers.
669	Answer the following question with "Yes" or "No" in JSON format.
670	
671	Example:
672	Q: Will we confirm evidence for megastructures orbiting the
673	star KIC 8462852?
674	A: { Value : NO }
676	Ω • α
677	A. $\{$ "value". "
V 11	
678	C.2 Base Models
679	For non-instruction-tuned models, we use the same examples but without additional instructions:
680	Quantity Estimation

12

Q:	How much global photovoltaic energy generation was deployed	681
	by the end of 2020?	682
A:	{ "value": 738 }	683
		684
Q:	\$question	685
A:	{ "value": "	686
Tiı	meframe Prediction	687
Q:	When did an AI system achieve a significant victory against	688
	a professional human in Starcraft 2?	689
A:	{ "value": "2019-01-24" }	690
		691
Q:	\$question	692
A:	{ "value": "	693
Bo	polean Questions	694
Q:	Will we confirm evidence for megastructures orbiting the	695
	star KIC 8462852?	696
A:	{ "value": "No" }	697
		698
Q:	\$question	699
A:	{ "value": "	700

701

702

705

706

707

708

713

D Hyperparameter Settings

D.1 Generation Hyperparameters

703 We generate responses using temperature-based704 sampling with the following hyperparameters:

- $max_length = 200$
- do_sample = True

• $top_k = 50$

• $top_p = 0.9$

Among the generated outputs, we select the one
with the highest confidence as the final prediction. All experiments are conducted using full
precision on an NVIDIA RTX 8000 GPU.

D.2 Evaluation Hyperparameters

714The scaling factor α in Equation 4, Equation 8,715and Equation 9 is set to 0.05. For Equation 5716and Equation 6, we set σ_{max} to 30 and σ_{min} to7171 for Timeframe Prediction, while for Quantity718Estimation, σ_{max} is 20 and σ_{min} is 1. These values ensure that evaluation metrics appropriately720scale errors and confidence calibration.

721 E Additional Results

722 This section provides extended results categorized by the training data cutoff date of each 723 model. Forecasting performance depends on 724 model architecture, scale, and knowledge re-725 cency, so we evaluate models with different cut-726 off dates to examine how access to more recent 727 information influences prediction accuracy and 728 confidence calibration. 729

Models trained after certain event resolutions
may have indirectly encountered outcomerelated information, potentially affecting evaluation fairness. This should be considered when
interpreting results.

Detailed model-specific performance metrics for	735
Boolean Questions, Timeframe Prediction, and	736
Quantity Estimation are presented in Table 10	737
to Table 15.	738
These results highlight trends in forecasting ac-	739
curacy and confidence calibration across models	740
with different knowledge recency.	741

Model	Accuracy (†)	F1 (†)	Brier (\downarrow)	ADE (\downarrow)	CRPS (T) (\downarrow)	APE (\downarrow)	$\mathrm{MAE}\left(\downarrow\right)$	$CRPS\left(Q\right)\left(\downarrow\right)$
GPT2	0.5835	0.3650	0.4199	0.9944	0.9884	0.2296	0.8697	0.8575
GPT2-XL	0.6708	0.3158	0.4540	1.0000	1.0000	0.0295	0.7461	0.7169
Pythia-14m	0.5960	0.1548	0.6200	0.9998	0.9996	0.1754	0.8604	0.8399
Pythia-160m	0.6135	0.2537	0.5264	0.9998	0.9998	0.0267	0.7551	0.7387
Pythia-2.8b	0.5337	0.3746	0.4562	0.9661	0.9646	0.0758	0.7946	0.7708
Bloom-560m	0.4763	0.4199	0.4764	1.0000	1.0000	0.0585	0.7695	0.7533
Bloom-7b1	0.6284	0.3196	0.3514	0.9681	0.9549	0.0503	0.7645	0.7381
Llama-7b	0.5536	0.3678	0.5278	0.9390	0.9294	0.0837	0.7223	0.6999
OLMo-1B	0.2269	0.2041	0.7906	1.0000	1.0000	0.1521	0.8430	0.8245
OLMo-7B	0.2070	0.1739	0.7986	1.0000	1.0000	0.2693	0.8791	0.8687
OLMo-7B-Instruct	0.6658	0.3333	0.3871	0.9201	0.9138	0.2935	0.7997	0.7814
OLMo-2-7B	0.5362	0.3307	0.4995	0.9225	0.9081	0.0363	0.6983	0.6752
OLMo-2-7B-Instruct	0.5935	0.4240	0.3992	0.8407	0.8264	0.1281	0.6948	0.6725

Table 10: Combined forecasting performance for cutoff 2017-12-01. CRPS (T) denotes the Continuous Ranked Probability Score for Timeframe Prediction, while CRPS (Q) denotes the Continuous Ranked Probability Score for Quantity Estimation.

Model	Accuracy (†)	F1 (†)	Brier (\downarrow)	ADE (\downarrow)	CRPS (T) (\downarrow)	APE (\downarrow)	$\text{MAE}\left(\downarrow\right)$	$CRPS\left(Q\right)\left(\downarrow\right)$
GPT2	0.5918	0.3750	0.4136	0.9942	0.9879	0.2342	0.8852	0.8749
GPT2-XL	0.6735	0.3190	0.4465	1.0000	1.0000	0.0277	0.7667	0.7366
Pythia-14m	0.5860	0.1374	0.6245	1.0000	1.0000	0.1696	0.8753	0.8538
Pythia-160m	0.6093	0.2584	0.5273	1.0000	1.0000	0.0274	0.7800	0.7620
Pythia-2.8b	0.5452	0.3906	0.4426	0.9650	0.9634	0.0742	0.8106	0.7851
Bloom-560m	0.5015	0.4393	0.4562	1.0000	1.0000	0.0481	0.7832	0.7651
Bloom-7b1	0.6297	0.3280	0.3477	0.9668	0.9531	0.0508	0.7877	0.7601
Llama-7b	0.5510	0.3677	0.5311	0.9365	0.9266	0.0634	0.7239	0.7012
OLMo-1B	0.2332	0.1928	0.7843	1.0000	1.0000	0.1568	0.8696	0.8491
OLMo-7B	0.2041	0.1538	0.8023	1.0000	1.0000	0.2568	0.8855	0.8731
OLMo-7B-Instruct	0.6647	0.2968	0.3975	0.9169	0.9103	0.2746	0.8002	0.7818
OLMo-2-7B	0.5335	0.3256	0.5057	0.9194	0.9045	0.0267	0.7097	0.6848
OLMo-2-7B-Instruct	0.5948	0.4085	0.3987	0.8343	0.8195	0.1226	0.7044	0.6821

Table 11: Combined forecasting performance for cutoff 2020-03-01. CRPS (T) denotes the Continuous Ranked Probability Score for Timeframe Prediction, while CRPS (Q) denotes the Continuous Ranked Probability Score for Quantity Estimation.

Model	Accuracy (†)	F1 (†)	Brier (\downarrow)	ADE (\downarrow)	CRPS (T) (\downarrow)	APE (\downarrow)	$\text{MAE}\left(\downarrow\right)$	$CRPS\left(Q\right)\left(\downarrow\right)$
GPT2	0.6084	0.3522	0.3754	1.0000	1.0000	0.2649	0.8513	0.8384
GPT2-XL	0.6730	0.2586	0.4463	1.0000	1.0000	0.0313	0.7427	0.7035
Pythia-14m	0.6008	0.1319	0.6531	1.0000	1.0000	0.2187	0.8598	0.8376
Pythia-160m	0.6312	0.2400	0.5190	1.0000	1.0000	0.0220	0.7278	0.7016
Pythia-2.8b	0.5171	0.3280	0.4376	1.0000	1.0000	0.0686	0.7692	0.7412
Bloom-560m	0.5015	0.4393	0.4521	1.0000	1.0000	0.0481	0.7273	0.7082
Bloom-7b1	0.6297	0.3280	0.3200	0.9984	0.9976	0.0644	0.7461	0.7176
Llama-7b	0.5741	0.3625	0.5184	0.9869	0.9807	0.0473	0.6532	0.6232
OLMo-1B	0.2395	0.1849	0.7721	1.0000	1.0000	0.2159	0.8414	0.8203
OLMo-7B	0.2243	0.1754	0.7954	1.0000	1.0000	0.2253	0.8482	0.8354
OLMo-7B-Instruct	0.6654	0.1782	0.4045	0.9234	0.9205	0.2189	0.7629	0.7372
OLMo-2-7B	0.5323	0.2692	0.4872	0.9266	0.9221	0.0270	0.6676	0.6427
OLMo-2-7B-Instruct	0.6084	0.3758	0.3858	0.8590	0.8385	0.1039	0.6382	0.6114

Table 12: Combined forecasting performance for cutoff 2021-12-01. CRPS (T) denotes the Continuous Ranked Probability Score for Timeframe Prediction, while CRPS (Q) denotes the Continuous Ranked Probability Score for Quantity Estimation.

Model	Accuracy (↑)	F1 (†)	Brier (\downarrow)	ADE (\downarrow)	CRPS (T) (\downarrow)	APE (\downarrow)	$\mathrm{MAE}\left(\downarrow\right)$	$CRPS\left(Q\right)\left(\downarrow\right)$
GPT2	0.6150	0.3556	0.3589	1.0000	1.0000	0.2576	0.8125	0.7984
GPT2-XL	0.6814	0.2653	0.4392	1.0000	1.0000	0.0345	0.7430	0.6980
Pythia-14m	0.6062	0.1558	0.6537	1.0000	1.0000	0.1946	0.8696	0.8472
Pythia-160m	0.6504	0.2524	0.5285	1.0000	1.0000	0.0221	0.7275	0.6998
Pythia-2.8b	0.5177	0.3230	0.4470	1.0000	1.0000	0.0843	0.7769	0.7435
Bloom-560m	0.5044	0.3978	0.4308	1.0000	1.0000	0.0528	0.7232	0.7058
Bloom-7b1	0.6460	0.3333	0.3061	0.9981	0.9971	0.0796	0.7765	0.7449
Llama-7b	0.5708	0.3433	0.5201	0.9843	0.9769	0.0564	0.6422	0.6126
OLMo-1B	0.2257	0.1800	0.7762	1.0000	1.0000	0.2206	0.8267	0.8032
OLMo-7B	0.2257	0.1702	0.8156	1.0000	1.0000	0.2238	0.8724	0.8575
OLMo-7B-Instruct	0.6460	0.1176	0.4043	0.9081	0.9046	0.2058	0.7799	0.7512
OLMo-2-7B	0.5398	0.2576	0.4698	0.9119	0.9065	0.0300	0.6814	0.6535
OLMo-2-7B-Instruct	0.5973	0.3259	0.3992	0.9283	0.8983	0.1038	0.6517	0.6216

Table 13: Combined forecasting performance for cutoff 2022-08-01. CRPS (T) denotes the Continuous Ranked Probability Score for Timeframe Prediction, while CRPS (Q) denotes the Continuous Ranked Probability Score for Quantity Estimation.

Model	Accuracy (†)	F1 (†)	Brier (\downarrow)	ADE (\downarrow)	CRPS (T) (\downarrow)	APE (\downarrow)	$\mathrm{MAE}\left(\downarrow\right)$	$CRPS\left(Q\right)\left(\downarrow\right)$
GPT2	0.5957	0.3559	0.3800	1.0000	1.0000	0.2642	0.7407	0.7222
GPT2-XL	0.6543	0.2353	0.4183	1.0000	1.0000	0.0383	0.7087	0.6530
Pythia-14m	0.6011	0.1562	0.6417	1.0000	1.0000	0.2195	0.8489	0.8243
Pythia-160m	0.6330	0.2581	0.4971	1.0000	1.0000	0.0220	0.6945	0.6611
Pythia-2.8b	0.5213	0.3382	0.4390	1.0000	1.0000	0.1080	0.7581	0.7199
Bloom-560m	0.5266	0.4183	0.4124	1.0000	1.0000	0.0654	0.6986	0.6713
Bloom-7b1	0.6649	0.3883	0.2964	0.9968	0.9952	0.1022	0.7897	0.7574
Llama-7b	0.5532	0.3214	0.5113	1.0000	1.0000	0.0566	0.6014	0.5689
OLMo-1B	0.2340	0.1591	0.7545	1.0000	1.0000	0.2124	0.7883	0.7623
OLMo-7B	0.2074	0.1687	0.8199	1.0000	1.0000	0.2157	0.8365	0.8201
OLMo-7B-Instruct	0.6543	0.1408	0.3811	1.0000	1.0000	0.2358	0.7987	0.7703
OLMo-2-7B	0.5106	0.2414	0.4438	0.9987	0.9980	0.0214	0.6464	0.6160
OLMo-2-7B-Instruct	0.5851	0.3276	0.4107	0.9100	0.8739	0.1150	0.6067	0.5748

Table 14: Combined forecasting performance for cutoff 2023-03-01. CRPS (T) denotes the Continuous Ranked Probability Score for Timeframe Prediction, while CRPS (Q) denotes the Continuous Ranked Probability Score for Quantity Estimation.

Model	Accuracy (†)	F1 (†)	Brier (\downarrow)	ADE (\downarrow)	CRPS (T) (\downarrow)	APE (\downarrow)	$\text{MAE}\left(\downarrow\right)$	$CRPS\left(Q\right)\left(\downarrow\right)$
GPT2	0.5931	0.3059	0.3673	1.0000	1.0000	0.2490	0.7319	0.7167
GPT2-XL	0.7103	0.3000	0.4207	1.0000	1.0000	0.0405	0.6887	0.6392
Pythia-14m	0.6207	0.1395	0.6479	1.0000	1.0000	0.2487	0.8481	0.8293
Pythia-160m	0.6690	0.3143	0.5090	1.0000	1.0000	0.0218	0.6731	0.6463
Pythia-2.8b	0.5241	0.2887	0.4481	1.0000	1.0000	0.1189	0.7315	0.6914
Bloom-560m	0.5241	0.3784	0.4103	1.0000	1.0000	0.0262	0.6740	0.6479
Bloom-7b1	0.6828	0.3611	0.3007	1.0000	1.0000	0.1140	0.7831	0.7528
Llama-7b	0.5379	0.2716	0.5452	0.9621	0.9443	0.0632	0.6173	0.5901
OLMo-1B	0.2345	0.2059	0.7507	1.0000	1.0000	0.2195	0.7639	0.7370
OLMo-7B	0.2207	0.2333	0.8271	1.0000	1.0000	0.2195	0.8218	0.8066
OLMo-7B-Instruct	0.6828	0.1569	0.3642	1.0000	1.0000	0.2682	0.8082	0.7842
OLMo-2-7B	0.5103	0.2444	0.4377	0.9981	0.9970	0.0206	0.6457	0.6206
OLMo-2-7B-Instruct	0.5862	0.3023	0.4078	0.8696	0.8197	0.0872	0.5968	0.5686

Table 15: Combined forecasting performance for cutoff 2023-12-01. CRPS (T) denotes the Continuous Ranked Probability Score for Timeframe Prediction, while CRPS (Q) denotes the Continuous Ranked Probability Score for Quantity Estimation.