
When Does Evidence Help Prompted LLM Forecasting?

Evidence Access and Prompt Structure Across 12 Models

Anonymous Authors¹

Abstract

We present a controlled study of how evidence access and prompt structure shape LLM-based forecasting. Using 114 resolved binary questions from ForecastBench, we evaluate 12 models across three prompting strategies: a direct control prompt, a base-rate prompt that asks models to anchor on historical frequency, and a Bayesian prompt that asks models to state a prior and update on evidence. We compare two information conditions: closed-book forecasting, where models receive no external evidence, and shared evidence retrieved via AskNews before the timestamp at which the market baseline is recorded. Across 8,208 forecasts, shared evidence improves Brier scores across models and prompts ($\Delta = -0.028$, 95% CI $[-0.036, -0.019]$, $p < 0.001$). Bayesian-style prompting performs worse than the control prompt in closed-book settings, consistent with recent evidence that structured reasoning prompts can degrade LLM forecasts when external information is unavailable. Although Bayesian prompting shows the largest numerical improvement from evidence, its advantage over the control prompt is only suggestive at pilot scale and does not reach conventional significance (DiD $= -0.011$, $p = 0.08$). In an exploratory extension, a Superforecaster-style prompt performs strongly; because it was not part of the main confirmatory comparison, we report it separately. Despite these gains, no LLM beats the freeze-time prediction-market baseline. These findings suggest that evidence access improves LLM forecasting, but that prompt structure alone is insufficient: even evidence-grounded models remain behind market baselines and suffer from overconfident extreme probabilities.

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

1. Introduction

Large language models have shown rapidly improving performance on short-horizon probabilistic forecasting tasks (Halawi et al., 2024; Karger et al., 2025; Bosse et al., 2026), yet two empirical results sit in tension. On the one hand, Schoenegger et al. (2025) systematically tested 38 prompt variants across six frontier LLMs and found that, after correcting for multiple prompt comparisons, no prompt significantly improved closed-book ForecastBench accuracy. In fact, Bayesian Reasoning and Propose-Evaluate-Select prompts significantly degraded performance. In their study, Bayesian Reasoning asks the model to state a prior, identify evidence, and update its probability, while Propose-Evaluate-Select asks the model to generate candidate forecasts, evaluate them, and select a final answer. On the other hand, fine-tuning LLMs to reason like Bayesians (Qiu et al., 2026) and integrating structured reasoning into agentic forecasting pipelines (Alur et al., 2025) have both been reported to yield forecasting gains. A natural reconciliation is that structured prompts may behave differently depending on whether the model has external evidence to reason over: in closed-book settings, the model must infer priors and updates from parametric memory alone; with retrieved evidence, the same structure can organise concrete external information rather than relying only on the model’s internal memory.

This paper isolates that hypothesis with a controlled factorial design. We cross three prompts taken verbatim from Schoenegger et al. (2025) (Control, Base-Rate-First, and Bayesian Reasoning) with two information conditions: closed-book forecasting and a shared evidence packet retrieved from AskNews before each question’s freeze time, the timestamp at which the market baseline is recorded. We evaluate 12 LLMs, including nine reasoning-token models and three standard models, on 114 binary ForecastBench questions resolving after the latest model’s training cut-off. The P1–P3 design was fixed before the final analysis; a fourth Superforecaster Persona prompt is reported separately as an exploratory extension.

Our contributions are fourfold: (i) a replication of Schoenegger et al.’s closed-book prompt-engineering null with newer frontier models; (ii) the first systematic test of how the same

prompts behave under shared evidence; (iii) evidence that shared retrieval improves forecasting accuracy across all tested LLMs but remains insufficient to close the gap to a freeze-time market baseline; and (iv) a calibration analysis showing that the remaining gap to the market baseline is associated with overconfident extreme forecasts, especially in the highest-probability bin.

2. Methods

Questions. We sample 114 resolved binary questions from ForecastBench (Karger et al., 2025) biweekly snapshots released between October 2025 and March 2026. Questions are drawn from Polymarket, Manifold, and Metaculus, and intentionally balanced across broad topic groups. All questions resolve on or after September 2025, after the latest stated training cutoff in our model panel, and pass manual sanity checks against question leakage.

Table 1. Experimental design summary. P1–P3 are the main confirmatory prompts; P4 is exploratory.

Component	Design
Questions	114 resolved binary ForecastBench questions
Sources	Polymarket (41), Manifold (37), Metaculus (36)
Models	12 LLMs: 9 reasoning-token, 3 standard
Prompts	P1 Control, P2 Base-Rate-First, P3 Bayesian
Conditions	Closed-book vs. shared evidence
Evidence	Up to 10 AskNews articles before freeze time
Main calls	$12 \times 3 \times 2 \times 114 = 8,208$
Exploratory	P4 Superforecaster prompt: 2,736 calls
Primary metric	Brier score
Baseline	Freeze-time market probability
Valid forecasts	8,196/8,208 (99.9%) for P1–P3

Prompts. We use three prompts verbatim from Schoenegger et al. (2025), Appendix A: **P1** Control (a direct probability request), **P2** Base-Rate-First (anchor on historical frequency, then adjust), and **P3** Bayesian Reasoning (state a prior, identify evidence, update). All prompts terminate with a fixed “Forecast : X%” output requirement. We additionally run an exploratory **P4** Superforecaster Persona prompt drawn from the same source, reported separately and not included in the main confirmatory analysis.

Information conditions. In the *closed-book* (CB) condition, each model receives only the question text, background, and resolution criteria. In the *shared-evidence* (SE) condition, we append a fixed evidence packet retrieved via the AskNews historical API. For each question, we query AskNews using the question title and key entities, restrict articles to those published before the question’s freeze time, and keep up to the 10 most relevant returned articles. The same evidence packet is given to every model for that question. This lets us test whether models benefit from receiving external evidence, rather than from differences in their own

retrieval systems.

Models. We test 12 LLMs accessed through OpenRouter: nine that produce internal reasoning tokens in our runs (gpt-5.4, claude-opus-4.6, gemini-3.1-pro, kimi-k2.6, grok-4.20, deepseek-v3.2-speciale, gpt-oss-120b, glm-5.1, qwen3-max) and three standard models (gemma-4-31b, gemini-3-flash, mistral-large).

Metrics and analysis. The primary metric is Brier score (Brier, 1950), $B = (\hat{p} - y)^2$, where \hat{p} is the model’s forecast probability and y is the binary outcome (1 if the event occurred, 0 otherwise). Lower Brier scores indicate more accurate probabilistic forecasts. Secondary metrics include Expected Calibration Error (ECE) over 10 equal-width probability bins, forecast standard deviation as a measure of sharpness, and the rate of catastrophic overconfidence: forecasts above 0.8 for events that did not occur, or below 0.2 for events that did occur.

To test whether prompt effects depend on evidence access, we fit a mixed-effects model with Brier score as the outcome:

$$Brier_{ijkl} \sim \text{Prompt}_i \cdot \text{Cond}_j + (1 | Q_k) + (1 | M_l). \quad (1)$$

Here, Prompt_i is the prompt type, Cond_j is the information condition (CB or SE), Q_k indexes questions, and M_l indexes models. The random intercepts for questions and models account for the fact that some questions are harder than others and some models are stronger than others. The $\text{Prompt} \times \text{Condition}$ term tests whether the effect of a prompt changes when evidence is provided. We also report question-clustered bootstrap confidence intervals and compare LLM forecasts to the freeze-time market baseline recorded in ForecastBench.

Parsing and integrity. Each model response is parsed by extracting the required “Forecast : X%” line and converting it into a probability. We hardened the parser to handle common formatting failures, including markdown wrapping, decimal/percentage mismatches, and truncated responses. The P1–P3 valid forecast rate is 99.9% (8,196/8,208); P4 is 98.5% (2,695/2,736). Each prompt–condition–model combination contains 114 attempted calls after deduplication.

3. Results

3.1. Shared evidence consistently improves aggregate forecasting accuracy

Shared evidence reduces mean Brier across the full P1–P3 design from 0.180 in the closed-book (CB) condition to 0.152 in the shared-evidence (SE) condition, an aggregate improvement of $\Delta = -0.028$ (Table 2). The effect is consistent across models: all 12 models achieve lower mean Brier under SE than under CB, and every model improves

Table 2. Mean Brier score by prompt and information condition, averaged over all models and questions. Lower is better.

Prompt	CB	SE	SE-CB
P1 (Control)	0.174	0.152	-0.022
P2 (Base-Rate-First)	0.175	0.149	-0.027
P3 (Bayesian)	0.190	0.157	-0.033
P4 [†] (Superforecaster)	0.165	0.139	-0.026

[†] Exploratory, not part of the main comparison.

on at least two of the three prompts. A bootstrap that re-samples questions, rather than individual forecasts, gives a significant main effect of evidence access ($p < 0.001$). The size of the effect varies substantially by model: gemini-3.1-pro shows the largest average SE gain ($\Delta = -0.060$), while mistral-large is only partially helped, worsening under P1 and P2 but improving strongly under P3.

3.2. Bayesian prompting is sensitive to evidence access

Replicating the closed-book pattern reported by Schoeneger et al. (2025), the Bayesian Reasoning prompt performs worst in the CB condition (Brier 0.190 vs. 0.174 for Control). However, P3 also shows the largest numerical improvement from shared evidence, improving by -0.033 Brier points. The planned difference-in-differences comparison, $(P3_{SE-CB}) - (P1_{SE-CB})$, is negative (DiD = -0.011), consistent with the hypothesis that Bayesian-style structure becomes more useful when evidence is available, but it does not reach conventional significance ($p = 0.08$).

3.3. Evidence narrows but does not close the market gap

Despite the gains from shared evidence, no LLM configuration matches the freeze-time prediction-market baseline. The market baseline achieves Brier 0.095. In the closed-book condition, the strongest P1–P3 configurations remain far above this baseline; under shared evidence, the gap narrows substantially, with the best P1–P3 model–prompt–condition combination reaching Brier 0.107. The exploratory P4 prompt improves the closed-book frontier and remains competitive under shared evidence, but it also does not beat the market. Overall, evidence access moves LLM forecasts closer to market performance, but the remaining gap suggests that retrieval alone is insufficient without better calibration and aggregation.

3.4. High-confidence miscalibration remains the main failure mode

Calibration analysis suggests that the remaining gap between LLM forecasts and the market baseline is strongly associated with high-confidence errors. Across P1–P3, the overall ECE is 0.082, indicating moderate miscalibration. However, this aggregate value hides a sharper failure mode:

Table 3. Best LLM configurations compared with the freeze-time market baseline. Lower Brier is better.

System / setting	Brier
Market baseline	0.095
Best P1–P3 CB model avg.: gpt-5.4	0.146
Best P1–P3 CB combination: gpt-5.4 / P1 / CB	0.142
Best P1–P3 SE model avg.: gemini-3.1-pro	0.113
Best P1–P3 SE combination: gemini-3.1-pro / P1 / SE	0.107
Best P4 CB combination [†] : gpt-5.4 / P4 / CB	0.137
Best P4 SE combination [†] : gpt-5.4 / P4 / SE	0.112

[†] Exploratory, not part of the main comparison.

overconfidence increases substantially in the upper probability bins. In the top decile, the mean forecast is 0.960 while the observed event frequency, i.e. the fraction of events in that bin that resolve true, is only 0.522 (Figure 1). Although LLMs are often described as cautious in open-ended settings, our forced-probability forecasting setup produces many extreme probability estimates. Across the full P1–P3 design, 384 of 8,196 valid forecasts (4.7%) are high-confidence false positives ($\hat{p} > 0.8, y = 0$), while 340 (4.1%) are high-confidence false negatives ($\hat{p} < 0.2, y = 1$). This indicates that improving evidence access alone is insufficient: reliable LLM forecasting also requires better calibration of extreme probabilities.

3.5. Exploratory P4 and heterogeneity analyses

The Superforecaster-style prompt is exploratory but performs strongly. It achieves mean Brier 0.165 in CB and 0.139 in SE, the lowest prompt aggregate in both conditions. It also improves calibration relative to P1–P3 (ECE 0.064 vs. 0.082) and reduces high-confidence false positives (3.2% vs. 4.7%). At the model level, P4 beats each model’s best P1–P3 prompt in 8/12 CB comparisons and 7/12 SE comparisons. Because P4 was not part of the main confirmatory hypothesis, we treat this as a hypothesis-generating result for the live forecasting experiment (a forthcoming Metaculus tournament evaluation) rather than as a primary finding.

Evidence effects are heterogeneous across topics and sources, though small-N topic subgroups (e.g., sports, N=8) show large deltas driven by individual questions; full breakdowns and per-topic N are reported in Appendix A.

4. Discussion

Three takeaways emerge. First, our results support retrieval as a central ingredient in LLM forecasting: *shared evidence improves aggregate forecasting accuracy for all 12 LLMs in our panel*, across prompt structures and model types. Closed-book LLM forecasting, in which the model relies

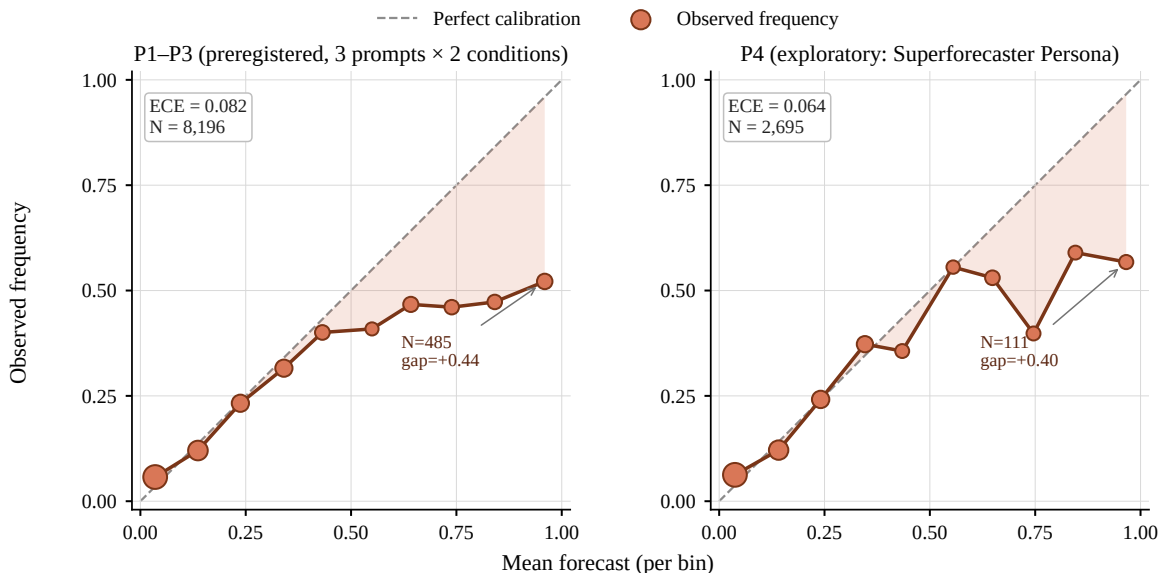


Figure 1. Calibration of P1–P3 forecasts ($N = 8,196$) and exploratory P4 forecasts ($N = 2,695$). Each point is a probability bin: x is the mean forecast and y is the fraction of questions in that bin that resolved true. The dashed diagonal indicates perfect calibration; the largest gap appears in the highest-probability bin.

entirely on parametric memory, leaves substantial accuracy on the table.

Second, our results are more cautious for prompting alone. The preregistered prompt \times evidence interaction is directionally consistent with the hypothesis that Bayesian structure becomes more useful when there is evidence to update on, but at pilot scale the effect does not reach conventional significance. Bayesian Reasoning, in particular, remains the worst prompt in closed-book and only partly recovers under shared evidence. The exploratory Superforecaster-style prompt outperforms the preregistered prompts in aggregate and warrants further study, but persona-style prompting is a different intervention from formal structured-reasoning prompting and should not be conflated with it.

Third, the remaining gap to the market baseline appears to have a large calibration component. LLMs in our panel are comparatively well calibrated at lower forecast probabilities but become severely overconfident at high probabilities, with the top decile of forecasts (mean $\hat{p} = 0.960$) resolving true only 52.2% of the time. This suggests that post-hoc calibration methods, such as isotonic regression or temperature scaling on probability outputs, may offer a more effective route than prompt-only changes for narrowing the market gap.

5. Limitations and Next Steps

The study is diagnostic rather than definitive. The pilot is restricted to 114 binary ForecastBench questions and two information conditions, so we make no claims about

numeric, multi-class, or open-ended prediction-market questions. The shared-evidence condition controls retrieval quality by giving every model the same AskNews packet, but therefore does not evaluate fully agentic retrieval systems that search, filter, and update evidence dynamically. The DiD test on the prompt \times evidence interaction is underpowered: a larger live evaluation is needed to confirm or reject the directional interaction at conventional significance. P4 is exploratory and was not part of the main confirmatory hypothesis. Finally, the magnitude of evidence effects in some topic subgroups is driven by small numbers of questions. Sports shows the largest topic-level SE benefit ($= 0.159$), but a manual audit confirmed all evidence articles were dated before each question’s freeze time and contained no resolution information; the effect is driven by two questions for which pre-freeze season standings were highly predictive of resolution. We therefore interpret topic-level deltas in small subgroups as descriptive.

Our follow-up experiment will run the same factorial design on live Metaculus tournament questions, extending beyond binary outcomes to include numeric and multi-class forecasts. The Superforecaster-style prompt will be added as a preregistered fourth prompt condition, and post-hoc isotonic calibration will be evaluated as an additional intervention. This follow-up is designed to test whether the prompt \times evidence interaction persists at larger N , whether the high-confidence calibration gap can be substantially reduced, and whether evidence grounding plus calibration can narrow the remaining gap to market baselines across richer forecasting formats.

References

- Alur, R., Stadie, B. C., Kang, D., Chen, R., McManus, M., Rickert, M., Lee, T., Federici, M., Zhu, R., Fogerty, D., Williamson, H., Lozinski, N., Linsky, A., and Sekhon, J. S. AIA forecaster: Technical report. *arXiv preprint arXiv:2511.07678*, 2025. URL <https://arxiv.org/abs/2511.07678>.
- Bosse, N. I., Mühlbacher, P., Wildman, J., Phillips, L., and Schwarz, D. Automating forecasting question generation and resolution for AI evaluation. *arXiv preprint arXiv:2601.22444*, 2026. URL <https://arxiv.org/abs/2601.22444>.
- Brier, G. W. Verification of forecasts expressed in terms of probability. *Monthly Weather Review*, 78(1):1–3, 1950.
- Halawi, D., Zhang, F., Yueh-Han, C., and Steinhardt, J. Approaching human-level forecasting with language models. *arXiv preprint arXiv:2402.18563*, 2024. URL <https://arxiv.org/abs/2402.18563>.
- Karger, E., Bastani, H., Yueh-Han, C., Jacobs, Z., Halawi, D., Zhang, F., and Tetlock, P. E. ForecastBench: A dynamic benchmark of AI forecasting capabilities. In *International Conference on Learning Representations (ICLR)*, 2025. URL <https://arxiv.org/abs/2409.19839>.
- Qiu, L., Sha, F., Allen, K., Kim, Y., Linzen, T., and van Steenkiste, S. Bayesian teaching enables probabilistic reasoning in large language models. *Nature Communications*, 17:1238, 2026. doi: 10.1038/s41467-025-67998-6.
- Schoenegger, P., Jones, C. R., Tetlock, P. E., and Mellers, B. Prompt engineering large language models’ forecasting capabilities. *arXiv preprint arXiv:2506.01578*, 2025. URL <https://arxiv.org/abs/2506.01578>.

A. Topic and source heterogeneity

Table 4 reports CB and SE mean Brier across the ten topic categories represented in our sample. Topic labels are auto-generated by an LLM classifier and then manually verified. Sample sizes are uneven: large topics (AI/technology, geopolitics, conflict, economics/finance) contain ≥ 14 questions each, while small subgroups (energy, public health, “other”) contain ≤ 6 . Topic-level SE–CB deltas in the small subgroups should therefore be interpreted as descriptive rather than as generalizable estimates.

Table 4. Per-topic mean Brier by information condition (P1–P3, all 12 models). N_q is the question count; forecasts per row $\approx 36N_q$. Topics ordered by SE–CB delta (most-helped first).

Topic	N_q	CB Brier	SE Brier	SE–CB	Direction
Sports	8	0.184	0.025	−0.159	helps [†]
Public health	6	0.255	0.168	−0.087	helps
Finance/markets	14	0.280	0.217	−0.063	helps
Other	5	0.203	0.141	−0.062	helps
Economics	14	0.128	0.087	−0.041	helps
AI/technology	20	0.126	0.097	−0.030	helps
Geopolitics	20	0.163	0.163	0.000	neutral
Conflict	18	0.170	0.189	+0.019	hurts
Entertainment	5	0.183	0.211	+0.028	hurts
Energy	4	0.252	0.389	+0.137	hurts

[†] Sports delta is driven by two questions for which pre-freeze season standings were highly predictive (Pistons Central Division lead; Yamal La Liga goal pace). A manual audit confirmed all evidence articles were dated before each question’s freeze time and contained no resolution information. See Section 5.

Table 5. Mean LLM Brier score by question source and information condition (P1–P3, all 12 models). Scores are computed from LLM forecasts, not from the market baseline. The source indicates the platform from which the ForecastBench question was drawn. SE–CB shows how much shared evidence changes LLM performance for questions from each source.

Source	N_q	CB Brier	SE Brier	SE–CB
Manifold	37	0.160	0.086	−0.074
Metaculus	36	0.228	0.210	−0.018
Polymarket	41	0.154	0.159	+0.005

B. Per-model results

Table 6. Per-model mean Brier across P1–P3 prompts and both information conditions ($N \approx 684$ per model). “Parse %” is the rate of successfully parsed forecasts; latency is the median per-call wall-clock time in seconds. Lower Brier is better.

Rank	Model	Type	Brier	Parse %	Lat. (s)
1	gpt-5.4	reasoning	0.137	100.0	62
2	gemini-3.1-pro	reasoning	0.143	99.7	35
3	kimi-k2.6	reasoning	0.144	100.0	249
4	claude-opus-4.6	reasoning	0.146	100.0	39
5	grok-4.20	reasoning	0.161	100.0	29
6	deepseek-v3.2-spec.	reasoning	0.164	100.0	197
7	glm-5.1	reasoning	0.173	99.7	77
8	qwen3-max	reasoning	0.176	99.4	251
9	gpt-oss-120b	reasoning	0.176	99.4	86
10	gemma-4-31b	standard	0.178	100.0	39
11	gemini-3-flash	standard	0.179	100.0	6
12	mistral-large	standard	0.216	100.0	17
<i>Market baseline</i>			0.095	—	—

Table 7. Per-model Brier across all prompt \times condition cells. P1–P3 are preregistered; P4 is exploratory. The best of P1–P3 per row and condition is in **bold**; P4 cells marked * outperform that row’s best P1–P3 prompt in the same condition.

Model	P1 (Control)		P2 (Base-rate)		P3 (Bayesian)		P4 [†] (Persona)	
	CB	SE	CB	SE	CB	SE	CB	SE
gpt-5.4	0.142	0.125	0.148	0.141	0.148	0.117	0.137*	0.112*
gemini-3.1-pro	0.169	0.107	0.175	0.117	0.175	0.117	0.146*	0.129
kimi-k2.6	0.162	0.141	0.165	0.119	0.150	0.127	0.152	0.132
claude-opus-4.6	0.161	0.137	0.168	0.123	0.163	0.127	0.152*	0.135
grok-4.20	0.178	0.154	0.176	0.140	0.179	0.138	0.171*	0.134*
deepseek-v3.2-spec.	0.176	0.145	0.166	0.148	0.185	0.167	0.164*	0.131*
gpt-oss-120b	0.168	0.150	0.164	0.154	0.222	0.195	0.158*	0.136*
glm-5.1	0.180	0.161	0.198	0.152	0.192	0.157	0.167*	0.170
qwen3-max	0.189	0.147	0.186	0.155	0.193	0.184	0.191	0.156
gemma-4-31b	0.175	0.173	0.188	0.162	0.196	0.176	0.181	0.148*
gemini-3-flash	0.194	0.167	0.186	0.167	0.190	0.170	0.177*	0.153*
mistral-large	0.193	0.216	0.186	0.205	0.284	0.209	0.186	0.138*

[†] Exploratory, not part of the main confirmatory analysis.

C. Prompt templates

We use three preregistered prompts taken verbatim from [Schoenegger et al. \(2025\)](#), Appendix A: P1 Control, P2 Base-Rate-First, and P3 Bayesian Reasoning. Each is run in both the closed-book (CB) condition (only the question) and the shared-evidence (SE) condition (an additional fixed block of up to 10 AskNews articles). The exploratory P4 Superforecaster Persona prompt is also drawn from the same source. Structural placeholders (`{question_title}`, `{background}`, `{resolution_criteria}`, `{resolution_date}`) are populated from each ForecastBench question record; `{retrieved_articles}` is present only in the SE condition.

Evidence block (SE condition only).

The following recent news articles may be relevant to this question. Use them if helpful; disregard if not.
`{retrieved_articles}`

P1 — Control.

Please answer the following question with a probabilistic estimate expressed between 0% and 100%, and format your response as: 'Forecast: X%'.
 Question: `{question_title}`
 Background: `{background}`
 Resolution Criteria: `{resolution_criteria}`
 Resolution Date: `{resolution_date}`

P2 — Base-Rate-First.

Please answer the following question with a probabilistic estimate expressed between 0% and 100%, and format your response as: 'Forecast: X%'.
 Before considering the specific details of this question, what is the historical frequency of similar events? Using this base rate as your starting point, adjust your probability estimate based on the particular circumstances of this case.
 Question: `{question_title}`
 Background: `{background}`
 Resolution Criteria: `{resolution_criteria}`
 Resolution Date: `{resolution_date}`

P3 — Bayesian Reasoning.

385 Consider the following question in terms of Bayesian reasoning. Start with a
386 prior probability based on historical data or general knowledge. Then, update
387 this prior using more specific information about the case under discussion. For
388 each new piece of information, produce an updated posterior estimate of the
389 outcome using the principle behind Bayes rule. Conclude with the final posterior
probability, formatted as: 'Forecast: X%'

390 Question: {question_title}
391 Background: {background}
392 Resolution Criteria: {resolution_criteria}
393 Resolution Date: {resolution_date}

394
395 **P4 — Superforecaster Persona (exploratory).** This persona-based prompt is drawn verbatim from Schoenegger et al.
396 (2025), Appendix A, and prepends a long superforecasting system message before the question block. The same evidence-
397 block convention as P1–P3 applies (appended only in the SE condition).
398

399
400 In this chat, you are a superforecaster who has a strong track record of accurate
401 forecasting. You evaluate past data and trends carefully for potential clues
402 to future events, while recognising that the past is an imperfect guide to
403 the future so you will need to put probabilities on possible future outcomes
404 (ranging from 0 to 100%). Your specific goal is to maximize the accuracy of
405 these probability judgments by minimising the Brier scores that your probability
406 judgments receive once future outcomes are known. Brier scores have two key
407 components: calibration (across all questions you answer, the probability
408 estimates you assign to possible future outcomes should correspond as closely as
409 possible to the objective frequency with which outcomes occur) and resolution
410 (across all questions, aim to assign higher probabilities to events that occur
411 than to events that do not occur). You outline your reasons for each forecast:
412 list the strongest evidence and arguments for making lower or higher estimates
413 and explain how you balance the evidence to make your own forecast. You begin
414 this analytic process by looking for reference or comparison classes of similar
415 events and grounding your initial estimates in base rates of occurrence (how
416 often do events of this sort occur in situations that look like the present
417 one?). You then adjust that initial estimate in response to the latest news
418 and distinctive features of the present situation, recognising the need for
419 flexible adjustments but also the risks of over-adjusting and excessive volatility.
420 Superforecasting requires weighing the risks of opposing errors: e.g., of
421 failing to learn from useful historical patterns vs. over-relying on misleading
422 patterns. In this process of error balancing, you draw on the 10 commandments
423 of superforecasting (Tetlock & Gardner, 2015) as well as on other peer-reviewed
424 research on superforecasting: 1. Triage 2. Break seemingly intractable problems
425 into tractable sub-problems 3. Strike the right balance between inside and
426 outside views 4. Strike the right balance between under- and overreacting to
427 evidence 5. Look for the clashing causal forces at work in each problem 6.
428 Strive to distinguish as many degrees of doubt as the problem permits but no more
429 7. Strike the right balance between under- and overconfidence, between prudence
430 and decisiveness 8. Look for the errors behind your mistakes but beware of
431 rearview-mirror hindsight biases 9. Bring out the best in others and let others
432 bring out the best in you 10. Master the error-balancing bicycle. Once you have
433 written your reasons, ensure that they directly inform your forecast. Then, you
434 will provide me with a number between 0 and 100 (up to 2 decimal places) that is
435 your best prediction of the event in the following format: ``Forecast: X%``.
436 Take a deep breath and work on this problem step-by-step. The question that you
437 are forecasting as well as some background information and resolution criteria are
438 below. Read them carefully before making your prediction.
439

432 Question: {question_title}
433 Background: {background}
434 Resolution Criteria: {resolution_criteria}
435 Resolution Date: {resolution_date}