

Martingale Score: An Unsupervised Metric for Bayesian Rationality in LLM Reasoning

Zhonghao He*

University of Cambridge
zh378@cam.ac.uk

Tianyi Qiu*

Peking University
qiutianyi.qty@gmail.com

Hirokazu Shirado†

Carnegie Mellon University
shirado@cmu.edu

Maarten Sap†

Carnegie Mellon University
maartensap@cmu.edu

Abstract

Recent advances in reasoning techniques have substantially improved the performance of large language models (LLMs), raising expectations for their ability to provide accurate, truthful, and reliable information. However, emerging evidence suggests that iterative reasoning may foster belief entrenchment, rather than enhancing truth-seeking behavior. In this study, we propose a systematic evaluation framework for *belief entrenchment* in LLM reasoning by leveraging the Martingale property from Bayesian statistics. This property implies that, under rational belief updating, the expected value of future beliefs should remain equal to the current belief, i.e., belief updates cannot be predicted from solely the current belief. We propose the unsupervised, regression-based *Martingale Score* to measure violations of this property, signaling a deviation from the Bayesian ability of updating on new evidence. In open-ended problem domains, including event forecasting, value-laden questions, and academic paper review, we found such violations to be widespread across models, reasoning paradigms, problem domains, and system prompts, where the future beliefs are consistently predictable from the model’s current belief, a phenomenon which we term *belief entrenchment*. Through comprehensive experiments, we identify the models (e.g., GPT-4o), reasoning techniques (e.g., chain of thought), and domains (e.g., forecasting) more prone to belief entrenchment. Finally, we validate the Martingale Score by showing that it predicts ground-truth accuracy on problem domains where ground truth labels are available. This indicates that, while designed as an unsupervised metric that operates even in domains without access to ground truth, the Martingale Score is a useful proxy of the truth-seeking ability of the LLM reasoning process.

1 Introduction

Consider a stubborn person who refuses to be shown wrong, or a king surrounded by sycophants. Human beings, in their efforts to seek true beliefs, often end up entrenching their pre-existing beliefs, whether due to their own confirmation bias [Klayman, 1995] or due to external

*Equal contribution.

†These authors jointly supervised this work.

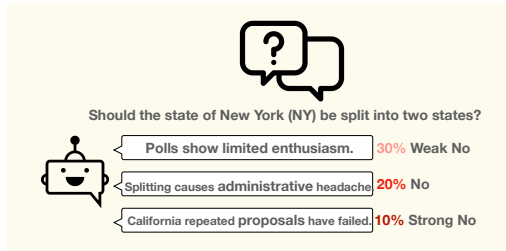


Figure 1: Example of Belief Entrenchment: LLM progressively updates beliefs in favor of its prior belief. Its belief update is highly predictable from the prior, violating the Martingale property.

confirmatory influence [Cinelli et al., 2021, Shirado et al., 2020]. Research suggests that such an entrenchment of belief lies at the heart of social epistemic problems such as polarization and misinformation [Del Vicario et al., 2017, Modgil et al., 2024, Lefebvre et al., 2024, Calhoun, 2004], and may be responsible for a wide range of downstream cognitive biases [Oeberst and Imhoff, 2023]. The entrenchment happens in the information processing of humans, a.k.a. their *reasoning* [Oeberst and Imhoff, 2023].

Having learned patterns of human reasoning from human-generated text data at an Internet scale, large language models (LLMs) are designed to help humans by answering their questions [Luo et al., 2022] or assisting with their tasks [Zheng et al., 2023]. Following their initial success in sophisticated cognitive tasks [Brown et al., 2020, Zhao et al., 2023], a range of *reasoning techniques* have emerged based on these models. From the earliest Chain-of-Thought (CoT) technique [Wei et al., 2022] to the more recent paradigm of inference-time scaling via reinforcement learning [Zelikman et al., 2024, Jaech et al., 2024, Guo et al., 2025] (*reinforced reasoning* henceforth), these methods aim to help language models in *truth-seeking*, i.e., gaining a correct understanding via argumentation, evidence-seeking, and trial-and-error at inference time.

However, such reasoning in language models can deviate from truth-seeking due to *belief entrenchment* — a systematic tendency to update beliefs *in favor of* prior opinions rather than *in response to* new evidence (Figure 1). When this occurs, it can degrade the accuracy of a model’s conclusions — mirroring well-documented effects in human cognition [Park et al., 2010] — and mislead users with an unjustified sense of confidence [Shi et al., 2024, Zhou et al., 2024a]. Although recent studies suggest the presence of belief entrenchment in LLMs [Schmidgall et al., 2024, Shi et al., 2024, Sumita et al., 2024], demonstrating it rigorously remains challenging. In any single example, it is difficult to distinguish a justified, prior-consistent update (e.g., a prediction later validated by evidence) from a biased update (e.g., representing evidence to maintain an unfalsifiable prior) [Atallah et al., 2021, Ji, 2023]. As a result, prior work relies on highly synthetic tasks or domain-specific setups where belief entrenchment is easy to detect [Schmidgall et al., 2024]. However, these methods fall short in assessing reasoning failures as they *unfold*, especially in open-ended tasks where the ground truth is ambiguous or unavailable.

To address this challenge, we propose a statistical measure of a reasoning model’s tendency toward belief entrenchment. Specifically, we define **belief entrenchment** as a violation of *the Martingale property* — a core principle in Bayesian reasoning which, informally, states that *the direction of belief updates should not be predictable in advance* [Chamley, 2004, Molavi, 2021]. This motivates using the Martingale violation as a principled and unsupervised indicator of irrational belief entrenchment. Concretely, if a model’s future belief can be reliably predicted from its prior belief across reasoning iterations—quantified via the goodness-of-fit of a regression model—this indicates a deviation from the Martingale property and thus the presence of belief entrenchment.

We empirically validate our measure’s usefulness through experiments on three domains: event forecasting, value-laden questions, and academic paper review (hereafter, Forecasting, r/ChangeMyView, and OpenReview, respectively). Across a wide range of models (e.g., GPT-4o, Llama 4), reasoning techniques (e.g., CoT, Debate), and system prompts, we find that belief entrenchment is a pervasive phenomenon. Our results show that future belief updates are consistently predictable from prior beliefs, a clear violation of the Martingale property. Critically, we validate the Martingale Score by demonstrating that it strongly correlates with a drop in ground-truth accuracy (measured by the Brier Score) in domains where such labels are available.

Contributions

- **The Belief Entrenchment Problem and the Martingale Score.** We define belief entrenchment, a statistical property that quantifies confirmation bias in LLM reasoning. We then introduce the Martingale Score, the predictability of belief updates based solely on prior, as an unsupervised and domain-agnostic measure of belief entrenchment.
- **Uncovering Widespread Belief Entrenchment in LLM Reasoning.** We use the Martingale Score to evaluate mainstream LLMs and find that belief entrenchment is a pervasive phenomenon across different domains (Forecasting, r/ChangeMyView, and OpenReview), prompts (prior-conforming, no-prompt, critical-thinking), and model families (GPT, DeepSeek, Gemini, Llama, etc).

- **Connecting Belief Entrenchment to Accuracy Loss.** We find that belief entrenchment consistently predicts an accuracy drop in problem domains where ground truth labels exist (e.g., forecasting). This suggests that the Martingale Score serves as a proxy for reasoning quality, even in settings where the ground truth is unavailable or still unfolding.

2 Related Work

Bayesian Rationality in Language Models In the LLM context, to evaluate Bayesian rationality is to evaluate the capacity to incorporate evidence via in-context learning in a way that is consistent with Bayes’s theorem. Gupta et al. [2025] presents that LLMs could follow Bayesian update when given a large amount of coin flips, despite a biased prior. However, more negative results appear in realistic settings with complex natural-language features: [Qiu et al., 2025a] demonstrated that LLMs do not update their beliefs of user preferences in multi-round interactions as expected by the Bayesian framework. [Falck et al., 2024] falsifies the hypothesis that LLM in-context learning is Bayesian. [Zhao et al., 2021, Wang et al., 2023] demonstrated that LLMs are sensitive to the arrangements of examples in prompts. Solutions are proposed to make the process of in-context learning or reasoning more Bayesian: mimicking the predictions of an ideal Bayesian reasoner [Qiu et al., 2025a], abstract reasoning [Zhou et al., 2024b], and combining abduction and deduction [Feng et al., 2024]. However, they tend to require auxiliary structures such as Bayesian networks, limiting their practical use.

Cognitive Biases in Language Models and in Humans LLMs suffer from a variety of cognitive biases as they are trained on a large corpus of human data [Echterhoff et al., 2024]. They also appear to employ biases distinct from those of humans when being deployed as judges, such as position bias (favoring certain positions), verbosity bias (favoring long answers), sentiment bias (preferring positive expressions) [Ye et al., 2024], and certainty bias [Zhou et al., 2024a]. Among all the cognitive biases, confirmation bias is of our particular attention as it is believed to be the root of polarization [Atallah et al., 2021]. It is also hypothesized that most cognitive biases are variants of confirmation bias [Oeberst and Imhoff, 2023]. Given these, we focus our attention on confirmation bias and use *belief entrenchment* to operationalize it in the LLM context.

Truthful AI and Truth-seeking AI Research in Truthful AI focuses on developing AI systems that output truth of the real world [Lin et al., 2021]. Previous work proposed interventions to improve “truthfulness” such as discovering learned activations aligned with truth [Burns et al., 2022] or shifting model activations during inference time [Li et al., 2023], and a trade-off between truthfulness and utility [Su et al., 2025]. Honest AI emphasizes that AI does not “intentionally” assert anything that does not align with its beliefs in pursuit of rewards [Evans et al., 2021]. This is closely related to the concept of deception in LLMs [Hubinger et al., 2024, Su et al., 2025]. In contrast, we are interested in a less investigated concept: truth-seeking AI [Koralus, 2025], the process of weighing evidence and discovering novel findings. Truth-seeking may become an important objective of LLMs as they are becoming more agentic [Chan et al., 2023], where the decision of what further information to seek is dependent on how the LLM interprets the current situation. Truth-seeking AI explicitly aims for gaining truth, but “truth-seeking” may not be a single metric to optimize against. Being Bayesian could be one target, and coherence might be another [Wen et al., 2025].

3 The Problems of Belief Entrenchment

3.1 Definitions

Belief entrenchment is the *systematic* tendency to update one’s beliefs in favor of one’s existing leanings rather than against, regardless of evidence. It is closely related to confirmation bias in cognitive science, where agents proactively seek, interpret, or recall evidence that is in favor of existing beliefs [Nickerson, 1998]. Confirmation bias is concerned with information processing at the *individual* instance level, and measuring it is difficult. In psychology, the measure of confirmation bias is task-specific and of low reliability. More importantly, to measure confirmation bias, the tendency of interpreting information that supports one’s beliefs and values, requires reliable access to and a measure of one’s prior beliefs [Berthet, 2021]. This is a significant challenge by itself.

In light of the difficulties, we use the *statistical violation* of the Martingale property from Bayesian statistics as the empirical measurement of belief entrenchment. Martingale property states that the expectation of the posterior, conditional on priors, should always equal the expectation of prior [Molavi, 2021]. In other words, under the Martingale property of Bayesian statistics, reasoners should never make predictable belief updates. Different from confirmation bias, because of how belief entrenchment is measured, it is defined as a statistical property, rather than an individual tendency.

Belief The term “beliefs” is used loosely in this study. We do not discuss questions such as whether LLMs could hold beliefs the way humans do. Rather, what we are concerned with is LLMs’ expressed “beliefs” through their output — when LLMs are entrenched by their own “beliefs,” they express high confidence in their own stated output, unjustified sense of confidence would mislead LLM users.

Belief update Relatedly, belief update refers to the *change in confidence* in LLM’s output. And we use LLM judge to assess confidence, as detailed in Section 5.1. We treat the process of having LLMs do chain-of-thought reasoning as a belief-updating process, similar to having a human being extensively reason about certain problems. In this process, the very beginning of model output is treated as “prior belief”, whereas the very end of model output, after extensive reasoning or engagement with external evidence, is treated as “posterior belief”.

3.2 Risks Brought on by Belief Entrenchment

We identify three levels at which *belief entrenchment* poses risks: the model’s reasoning capability, reasoning evaluation, and human-AI interactions.

Bayesian Reasoning and Truth-Seeking Previous research presents mixed results to the question whether LLM *in-context learning* is Bayesian [Gupta et al., 2025, Falck et al., 2024]. Meanwhile, specific reasoning failure instances were reported: sycophancy [Sharma et al., 2023], inverse scaling [Gema et al., 2025], following group majority rather than sticking to the truth [Weng et al., 2025]. Those specific failure cases of Bayesian reasoning compromise LLM task performance.

Failures in Reasoning Evaluation Empirically, LLM reasoning improves performance on a variety of tasks, including mathematics [Lu et al., 2023] and coding [Gu et al., 2024]. Reasoning evaluation is usually conducted by measuring ground truth accuracy [Sawada et al., 2023, Phan et al., 2025]. Such *outcome-based* evaluation falls short of our expectations because it does not tell us *how* LLM reasoning achieves superior performance, hence whether it would generalise out-of-distribution. If it is achieved by recalling parametric beliefs acquired from pre-training, LLM reasoning would have limited utility in real-world tasks because, in real-world tasks, problem-solvers need to take contextual information into account. We need process-based metrics to answer the question of “how” [Mondorf and Plank, 2024a].

Misleading Human-AI Interactions Research on sycophancy [Oeberst and Imhoff, 2023], and conformity [Weng et al., 2025] demonstrates that LLMs may favor pleasing human users or following the group majority over truth. This may mislead users into entrenchment of their own fallacies or biases. On the other hand, it is also known that human reasoning suffers from confirmation bias, and it causes downstream social epistemic problems, such as polarization [Lefebvre et al., 2024] and misinformation [Shirado et al., 2020]. The feedback loops between humans and LLMs (that humans acquire beliefs from LLMs that are trained on data containing human beliefs) may lead to lock-in of false beliefs collectively [Burton et al., 2024, Qiu et al., 2025b, Weidinger et al., 2023].

4 Measuring Belief Entrenchment with Martingale Score

To address these problems, we propose the Martingale Score as a principled, unsupervised measure of belief entrenchment in LLM reasoning. Effective reasoning requires the capacity to update beliefs in response to new evidence, a property aligned with Bayesian rationality. The Martingale Score quantifies the degree to which belief updates can be predicted from prior beliefs alone; a higher score indicates stronger entrenchment and weaker responsiveness to new information (Figure 2).

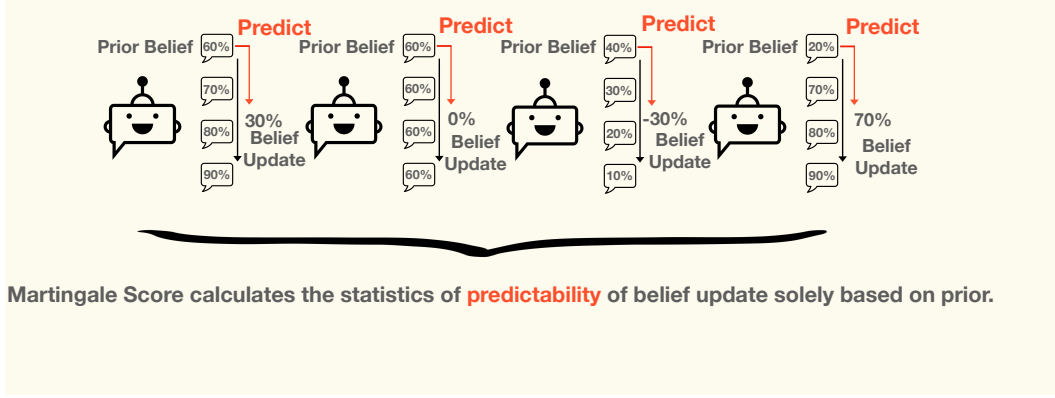


Figure 2: An illustration of Martingale Score calculation in our setting. “Prior Belief” refers to the expressed beliefs in most immediate LLM output; whereas “Posterior Belief” usually refers to the terminal beliefs after extended reasoning or engagements with external evidence. “Prior” and “Posterior” are relative concepts and their difference is taken as “Belief Update”. We estimate the linear coefficient when running a linear regression between belief updates and prior beliefs. The positive value of the linear coefficient is our practical choice of the Martingale Score, measuring the predictability of belief update solely based on prior belief.

4.1 Defining the Martingale Score

To compute the Martingale Score, we perform the regression $\Delta b = \beta_1 \cdot b_{\text{prior}} + \beta_0 + \epsilon$, where b_{prior} are the prior probabilities, $\Delta b = b_{\text{posterior}} - b_{\text{prior}}$, and ϵ is the error term.

We define the sample estimate $\hat{\beta}_1$ of the linear coefficient as the Martingale Score M , with the Ordinary Least Squares (OLS) method. Equivalently, when there are n samples,

$$M = \hat{\beta}_1 = \frac{\sum_{i=1}^n (\Delta b_i - \bar{\Delta b})(b_{\text{prior},i} - \bar{b}_{\text{prior}})}{\sum_{i=1}^n (b_{\text{prior},i} - \bar{b}_{\text{prior}})^2} \quad (1)$$

Martingale Score M measures the extent to which the prior belief b_{prior} positively (or negatively, if $M < 0$) predicts belief update Δb . Using OLS allows us to test the statistical significance of M , assessing whether the relationship between Δb and b_{prior} is distinguishable from zero (e.g., via a t-test with $p < 0.05$).

We choose the $M = \hat{\beta}_1$ definition for its simplicity, lack of confounders (as opposed to R^2 (coefficient of determination) of a logistic regression that introduces confounders such as the intrinsic variance of belief update not attributable to prior belief), and empirical reliability (as opposed to logistic regression on the binary *direction* of update, which neglects the magnitude of belief updates and produces random-seeming results).

4.2 Theoretical Justification for the Martingale Score

The **Martingale property** states that the expectation over one’s posterior, conditional on their prior, should always be equal to the prior [Molavi, 2021]. Formally,

$$\mathbb{E}[\Delta b \mid b_{\text{prior}} = p] = 0, \quad \forall p \in [0, 1]. \quad (2)$$

This implies that the direction of a Bayesian agent’s belief update (whether positive or negative) should not be predictable from the prior alone. Indeed, the Martingale property has been shown to be the defining characteristic of Bayesian rationality [Molavi, 2021].

The Martingale property (2) implies that the prior b_{prior} is statistically *exogenous* to the belief update Δb [Hayashi, 2011]. This exogeneity yields two desirable properties: the expected coefficient $\mathbb{E}[\hat{\beta}_1] = 0$ in the regression (1), and consistency of the estimator, i.e., $\hat{\beta}_1 \rightarrow 0$ in probability as the number of samples approaches $+\infty$ [Hayashi, 2011]. These properties justify **the Martingale Score**, defined in (1), as a principled measure of violations of the Martingale property (2). Formally,

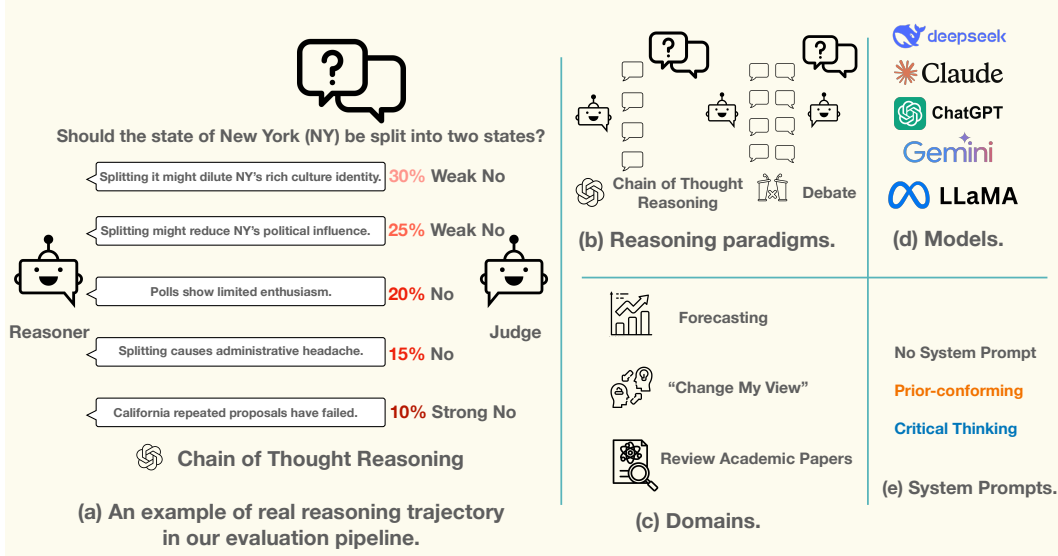


Figure 3: All experiment setups. **(a)** Example reasoning trajectory. **(b)** CoT reasoning and debate reasoning, where in the latter, one model is instructed to debate with its clone. **(c)** Problem domains: forecasting questions from Metaculus and Polymarket; value-laden questions from *r/ChangeMyView*, where owners of posts share statements they hold strong beliefs in, expecting counterarguments; acceptance decisions for ICLR submissions given the abstract and reviews. **(d)** Models evaluated. **(e)** System prompts used. “Prior-conforming prompt” instructs the model to fixate on their prior beliefs, whereas “critical thinking” prompt encourages models to challenge their prior beliefs. The two prompts represent the extreme behaviors we intend to demonstrate.

Proposition 1 *If the Martingale property holds, the population coefficient β_1 is 0 and the sample estimate $\hat{\beta}_1$ is an unbiased and consistent estimator of β_1 . In this specific case where $\beta_1 = 0$, this implies $E(M) = 0$ and $M \xrightarrow{p} 0$, as $n \rightarrow \infty$.³*

As the Martingale property is a prerequisite for Bayesian reasoning (i.e., ensuring that updates are driven by new evidence and not by prior views), its violation indicates belief updates are systematically predictable from priors. And the violation of Martingale property, in turn, defines the core of **belief entrenchment**. We thus operationalize belief entrenchment as the degree to which a model violates the Martingale property.

5 Experiment Setups

We set up experiments to evaluate belief entrenchment with the Martingale Score for a broad coverage of tasks (detailed in Section 5.2, and how it affects accuracy in domains where there is ground truth. The implementation details can be seen in Appendix C.

5.1 Using LLM Judge to Measure Beliefs

LLMs are often poorly calibrated: their confidence scores—whether expressed through token probability or self-reported scores—do not reliably reflect their true degree of belief [Pawitan and Holmes, 2024, Zhou et al., 2023]. To address this, we adopt an “expressed belief” approach rather than relying on “internal belief,” which remains an open research problem [Hase et al., 2021, Scherrer et al., 2023, Herrmann and Levinstein, 2024]. Expressed beliefs are inferred from the model’s outputs and are more consistent with how users experience and interpret LLM responses in practice [Zhou et al., 2024a]. To extract these beliefs, we employ a separate “judge” model (e.g., GPT-4o) that assesses each model’s reasoning steps and assigns a belief score $b \in [0, 1]$. To ensure robustness, we evaluate

³Refer to Appendix A for the proof.

multiple judge models and confirm that the results are consistent across them. Details can be seen in Section 6.2.

5.2 Problem Domains

To study belief entrenchment in LLM reasoning—specifically, how models incorporate new evidence during the reasoning process—we select domains that meet the following criteria:

- **Not solvable by memorization.** The domain should include questions that cannot be answered using information seen during pretraining. For example, we target events or facts that were resolved after the model’s knowledge cut-off. If a model was trained up to August 2024, it cannot know who won the 2024 U.S. presidential election without access to external tools.
- **Contain new evidence that could shift beliefs.** The domain must include incoming evidence that a Bayesian reasoner would use to revise its beliefs. This allows us to evaluate whether LLMs appropriately update their beliefs in response to new information, or whether they remain anchored to prior assumptions.
- **Ground truth becomes available after models’ knowledge cut-off.** The domain should provide verifiable ground truth labels after the model’s knowledge cut-off date. This enables us to assess whether belief entrenchment correlates with a drop in final accuracy when models fail to adapt to post-cut-off evidence.

Based on these criteria, we choose three domains for evaluating belief entrenchment: forecasting, value-laden questions, and academic paper review.

5.2.1 Forecasting

We source forecasting questions from Metaculus [Metaculus, 2015] and Polymarket [Polymarket, 2020] to test belief entrenchment in LLMs. We choose this domain for two key reasons: (1) forecasting questions come with ground truth labels once resolved; and (2) achieving high accuracy on a set of forecasting questions reflects Bayesian-like reasoning, as it requires seeking evidence and proportionally updating beliefs. While forecasting differs from factual question answering, accurate forecasting nonetheless requires well-calibrated, unbiased belief updates. This makes it a strong proxy for rational reasoning under uncertainty.

5.2.2 Value-Laden Questions

To assess belief entrenchment in subjective or controversial domains, we use questions from the *r/ChangeMyView* subreddit [Tan et al., 2016]. These discussions are explicitly designed to explore whether individuals (or in this case, LLMs) can revise their opinions when presented with counter-arguments. This allows us to examine whether LLMs update their value-oriented stances during multi-step reasoning or remain anchored to prior views.

5.2.3 Academic Paper Review

Scientific peer review is another setting that satisfies all three criteria for evaluating belief entrenchment. We use the open-access ICLR submission dataset from OpenReview [Höpner et al., 2025], which includes paper abstracts, bibliographies, reviewer comments, rebuttals, and final acceptance decisions. In our setup, the model is prompted to act as an area chair, making a final acceptance decision based on the abstract and the arguments presented in the reviews and rebuttals. This task allows us to evaluate how reasoning unfolds when prior impressions (e.g., from the abstract) may conflict with later-stage evidence (e.g., critiques and responses).

5.3 Reasoning Techniques

We evaluate belief entrenchment with two reasoning techniques: Chain-of-Thought (CoT) [Wei et al., 2022] and debate [Khan et al., 2024]. Debate lets two clones of the same expert model hold opposing positions on a topic and participate in debate, allowing a less informed observer to arrive at a better understanding of the subject.

Table 1: Martingale Scores under different setups. *CT* is short for *critical thinking*, while *PC* is short for *prior-conforming*. A positive Martingale Score M indicates that per unit increase in b_{prior} , there is an M -unit increase in Δb . Entries in this table measures belief entrenchment *per reasoning step*, and the bias may add up to much high levels during the full reasoning trajectory. Martingale Scores whose t-test produces $p < 0.05$ are marked with *.

		Forecasting		ChangeMyView		OpenReview	
		CoT	Debate	CoT	Debate	CoT	Debate
GPT-4o (May 7)	No Prompt	+0.0018	−0.0439	+0.0671*	+0.0941	+0.0734*	+0.1891*
	CT Prompt	+0.0156*	−0.0233	+0.0659*	+0.0822	+0.1030*	+0.1770*
	PC Prompt	+0.0896*	−0.0227	+0.1455*	+0.1572*	−0.0859*	+0.1718*
DeepSeek R1	No Prompt	+0.0207*	+0.0559	+0.0502*	+0.0845	+0.0676*	+0.0366
	CT Prompt	+0.0119*	+0.0121	+0.0511*	−0.0622	+0.0595*	+0.1860*
	PC Prompt	+0.0450*	+0.0487	+0.0526*	+0.0961	+0.0689*	+0.0299
DeepSeek V3	No Prompt	+0.0335*	−0.0929	+0.1155*	+0.0739	+0.1028*	+0.1337
	CT Prompt	+0.0348*	−0.0064	+0.0990*	+0.0179	+0.0865*	+0.0743
	PC Prompt	+0.0763*	−0.0216	+0.0879*	+0.0511	−0.1493*	+0.2113*
Gemini 2.0 Flash	No Prompt	+0.0764*	−0.0196	+0.1209*	+0.0969	+0.1012*	+0.0882
	CT Prompt	+0.0067	−0.0012	+0.1203*	+0.0642	+0.0817*	+0.1263*
	PC Prompt	+0.0335*	−0.0368	+0.1052*	+0.0295	+0.0849*	+0.0646
Llama 4 Scout	No Prompt	+0.0350*	+0.0078	+0.1420*	+0.0900	+0.0890*	+0.1168
	CT Prompt	+0.0125	−0.0395	+0.1146*	+0.0238	+0.1028*	+0.1729*
	PC Prompt	+0.0740*	−0.0114	+0.1372*	+0.0003	−0.0253	+0.1929*
Llama 4 Maverick	No Prompt	+0.0178*	+0.0103	+0.1038*	+0.1100*	+0.0823*	+0.1749*
	CT Prompt	+0.0282*	+0.0132	+0.1161*	+0.1185*	+0.0909*	+0.2521*
	PC Prompt	+0.0523*	−0.0128	+0.1435*	+0.1608*	+0.0951*	+0.1724*

5.4 Models and Prompts

We conduct experiments using GPT-4o, DeepSeek R1, DeepSeek V3, Gemini 2.0 Flash, LLaMA 4 Scout, and LLaMA 4 Maverick. To examine how prior beliefs affect reasoning, we manipulate model priors via system prompts. Specifically, in addition to a baseline with no system prompt, we introduce two prompting conditions: a *prior-conforming* prompt and a *critical-thinking* prompt.

The prior-conforming prompt reinforces a belief aligned with the model’s likely initial stance, serving both as a sense check and as a potential lower bound for belief entrenchment and accuracy. In contrast, the critical-thinking prompt encourages openness to counter-evidence and rational belief revision. All prompts used in the experiments are provided in Appendix C.

6 Results

This section presents our empirical findings. We first establish an association between belief entrenchment (as measured by the Martingale Score) and drops in ground-truth accuracy, as a vindication for the former. We then benchmark a range of closed- and open-weights models on different problem domains and setups, and causally attribute belief entrenchment to various factors like models, prompts, and reasoning techniques. Finally, we confirm LLM-as-a-judge validity.

6.1 Belief Entrenchment Results

Belief Entrenchment is Prevalent Across Setups Table 1 shows the Martingale Score of different models under different experiment setups. A positive Martingale Score M indicates that there is an M -unit increase in Δb per unit increase in b_{prior} . In most of the experiments, including almost all of those with CoT (51 out of 54), we see positive Martingale Scores, suggesting consistent belief entrenchment.

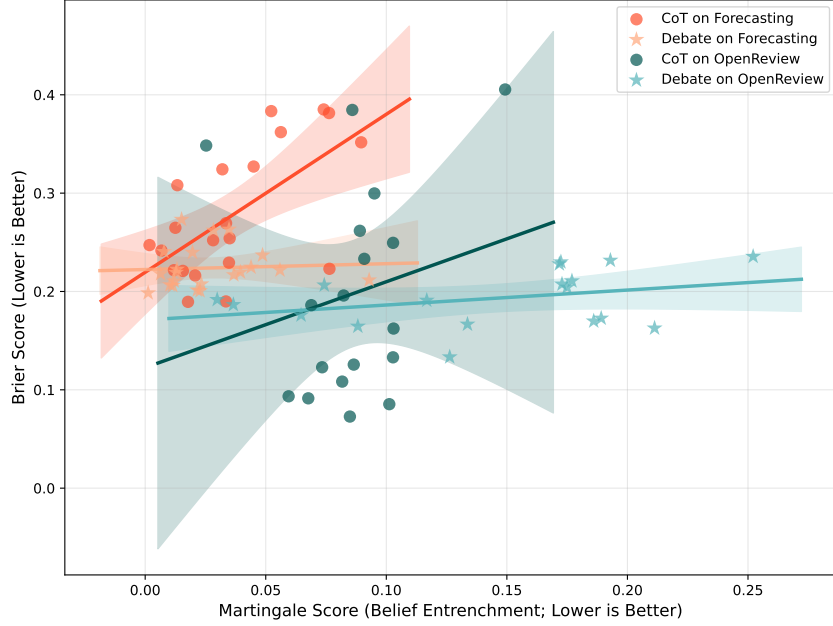


Figure 4: Relationship between the absolute value of the Martingale Score and the Brier Score. The former indicates the predictability of belief update based solely on the prior belief, and the latter measures the accuracy of probabilistic predictions. We observed that the Martingale Score and Brier score are positively correlated across all setups, suggesting belief entrenchment harms accuracy on binary problems. Taking “CoT on Forecasting” as an example, a Martingale Score of 0.0 corresponds to a Brier score smaller than 0.25, slightly better than random guess (0.250) [Halawi et al., 2024]; in contrast, when the model is mildly entrenched on its prior belief (marked by Martingale Score of 0.04), the forecasting performance is worse than random guess.

Notably, we observe overall more severe belief entrenchment in value-laden domains such as *r/ChangeMyView* (when using the same model, under the same system prompt, the Martingale Score is always higher in *r/ChangeMyView* than in Forecasting), suggesting that belief entrenchment might be a more serious concern in domains where we cannot evaluate LLM performance against ground truth. We also intended to manipulate the degree of belief entrenchment by using a critical thinking prompt and a prior-conforming prompt. We observe that belief entrenchment is more severe under a prior-conforming prompt and less severe under a critical thinking prompt in Forecasting, but not in *r/ChangeMyView* or OpenReview. More discussion about factors influencing belief entrenchment can be seen in Appendix B.

Belief Entrenchment is Not an Artifact Under “prior-conforming prompt”, belief entrenchment is meant to happen, as LLMs are instructed to emphasize arguments in favor of their prior belief, and such reasoning harms performance. However, we noticed that even under “no system prompt” and “critical thinking prompt”, belief entrenchment also happens, although to a lesser extent ($\bar{M}_{\text{Prior-conforming}} = 0.082 \pm 0.018$, $\bar{M}_{\text{No-prompt}} = 0.075 \pm 0.014$, $\bar{M}_{\text{Critical-thinking}} = 0.072 \pm 0.018$, with 95% CI). The results are significant, demonstrating a consistent tendency of belief entrenchment in LLM behavior.

Belief Entrenchment Harms Accuracy As the ultimate aim of reasoning in LLMs is to obtain true beliefs, we show that belief entrenchment diminishes this objective. We do so within problem domains where ground truth labels are available (i.e., Forecasting and OpenReview).

Figure 4 shows the correlation between the absolute value of the Martingale Score (lower is better) and the Brier Score (measuring prediction accuracy by the mean squared error between a predicted probability and the actual outcome; lower is better). Each data point represents the Brier Score and Martingale Score of one setup (of one model, with one type of system prompt and one reasoning mode, on > 100 questions in one problem domain).

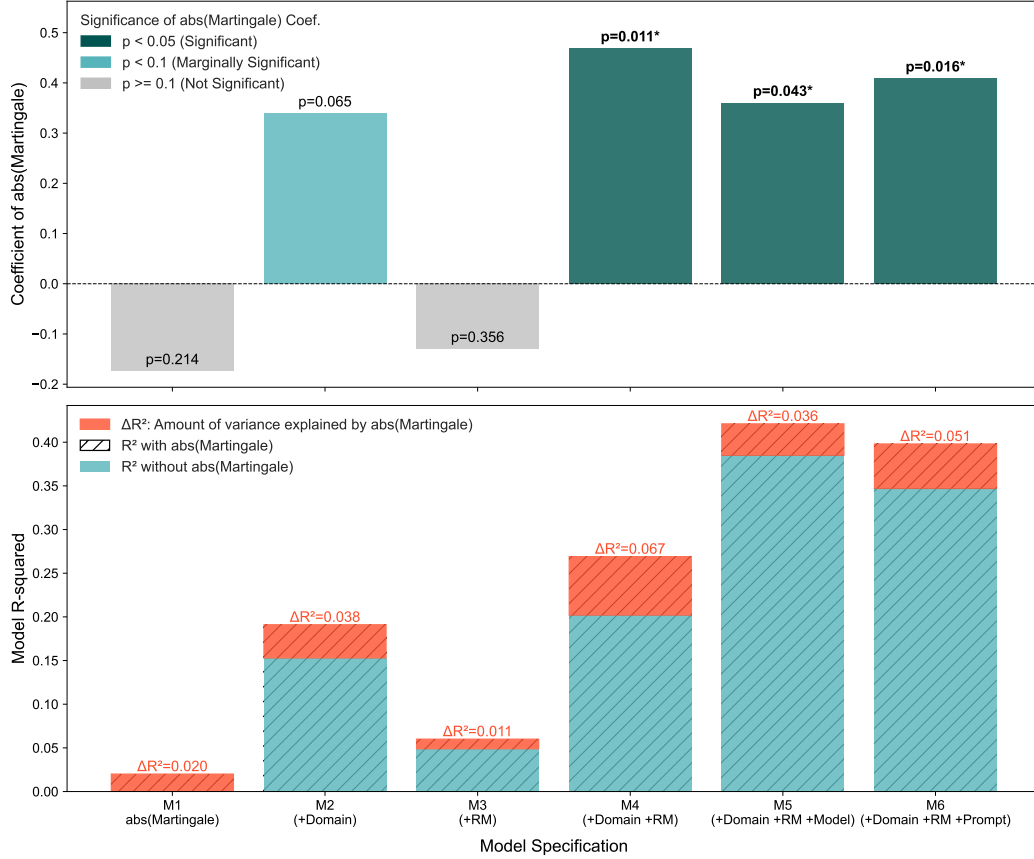


Figure 5: Increased absolute value of the Martingale Score is associated with worse prediction accuracy (higher Brier Scores) and explains a significant portion of the latter’s variance. In each regression model, we predict the Brier Score with the absolute value of the Martingale Score, while controlling for different potential confounders, including problem domain, reasoning techniques (“RM”), choice of model, and choice of prompt.

While outcome-based metrics (such as Brier score) measure model performance, they reveal little about *why* models achieve high performance, and for this reason, we need process-based metrics on reasoning [Mondorf and Plank, 2024b]. We show a positive correlation between the Martingale Score and Brier Score in Forecasting, suggesting that belief entrenchment is one of the culprits that compromise LLM reasoning.

Martingale Score Can be Used in Open-ended Domains It’s worth noting that the Martingale Score, being an unsupervised measure, is directly applicable in open-ended problem domains (e.g., r/ChangeMyView) as well. Table 1 demonstrates that, at least with CoT, belief entrenchment consistently happens regardless of prompt types and models. The comparable level of belief entrenchment ($\overline{M}_{\text{CMV-CoT}} = 0.103 \pm 0.013$, $\overline{M}_{\text{Forecasting-CoT}} = 0.037 \pm 0.011$, $\overline{M}_{\text{OpenReview-CoT}} = 0.086 \pm 0.012$, with 95% CI) corresponds to accuracy drops in forecasting and OpenReview (as in Figure 4 and Figure 5), where ground truth exists, suggesting consistently worsened judgment under belief entrenchment. Ruling out potential artifacts (presented in Section 6.1), we thus think the Martingale Score is a valid metric of reasoning paradigms with the utility of understanding the quality of reasoning.

6.2 Judge Consistency Evaluation Results

In our cross-judge consistency and human-LLM agreement analysis, all judges (LLMs or humans) show a strong correlation with GPT-4o.

Cross-LLM Agreement We construct pairs of LLM judges (e.g., GPT-4o VS Gemini-2.5-pro, GPT-4o VS DeepSeek-v3) and see how much their belief evaluation correlates with each other. Note that we’ve acquired GPT-4o data for all problems but significantly less with other judges (from 283 problems with GPT-4.1-mini to 3,844 with DeepSeek-v3), so we set up GPT-4o as the default judge for all comparisons, in line with our choice of judge in the main experiments.

Human-LLM Agreement We use a small batch of human evaluation data to validate the LLM judge (i.e., human-LLM consistency evaluation). Specifically, we request human evaluators to do belief evaluations exactly like what we request LLM judges to do; we then construct pairs of human-LLM judges (e.g., human evaluator 1 VS GPT-4o). Full results can be seen in Table 2.

Table 2: Inter-rater Agreement of Judges with GPT-4o.

Rank	Judge Model	Batches	Problems	Belief Samples	Pearson r	Spearman ρ	p -value
1	Human Evaluator 1	2	20	195	0.8822	0.8770	< 0.001
2	DeepSeek-v3	48	3,834	24,921	0.7774	0.7620	< 0.001
3	GPT-4.1-mini	3	283	2,015	0.7581	0.7490	< 0.001
4	Gemini-2.5-pro	4	373	1,688	0.7460	0.7230	< 0.001
5	Human Evaluator 2	2	18	173	0.7152	0.6812	< 0.001

All judges show a large positive correlation with GPT-4o, and all results are statistically significant ($p < 0.001$)⁴.

Considering results from both cross-judge consistency analysis and human-evaluation validation, it is highly unlikely that our belief entrenchment results stem from judge bias.

7 Conclusion

In this paper, we propose the Martingale Score as a principled and unsupervised measure of belief entrenchment in LLM reasoning. We show, as validation, that the Martingale Score predicts the ground-truth accuracy of a reasoning process, and then apply it to a range of different models, prompts, and problem domains. We conduct analysis and identify a collection of factors that increase or alleviate the severity of belief entrenchment.

Limitations In OpenReview, this study does not demonstrate the correlation between Martingale Score (belief entrenchment) and Brier Score (accuracy). We think this is because ground-truth labeling in domains such as OpenReview is community-voted (in this case, paper acceptance decisions), and community-voted decisions can be noisy [Beygelzimer et al., 2023]. However, Table 1 shows that belief entrenchment is a more severe phenomenon in domains where human judgments are required (such as r/ChangeMyView and OpenReview). This raises concerns over adopting such LLMs in those domains.

Future Work Due to resource constraints, we have not systematically studied belief entrenchment in reinforced reasoning [Guo et al., 2025], despite the recent popularity of this approach to LLM reasoning. Besides, the violation of ideal Bayesian rationality can be demonstrated by both irrational engagement with *internal reasoning process* and *external evidence*. In this study we only focused on the former. Future study can include an “evidence searching” component to evaluate LLM’s capacity to update belief in response to new *external evidence* (as opposed to its propensity to fixate on prior belief). The difficulty here is experimenters need to guarantee that the evidence searching does not give LLMs under evaluation direct access to ground truth, if such ground truth exists. In domains with ground truth (e.g., Forecasting, but not OpenReview), we demonstrated the *causality* of belief entrenchment and accuracy. Future research can test the validity of the Martingale Score in *open-ended domains* by demonstrating other downstream consequences of belief entrenchment.

To demonstrate the usefulness of Martingale Score as a general-purpose process-based reasoning evaluation, one could extend our pipeline into the evaluation of sycophancy [Sharma et al., 2023],

⁴As a reference point on a different setup, the NeurIPS 2021 review consistency experiment shows an $r = 0.58$ correlation between paper acceptance decisions independently made by two committees [Beygelzimer et al., 2023]

group conformity [Weng et al., 2025], and inverse scaling [Gema et al., 2025]. Martingale Score as an evaluation metric can be converted into a training objective, as a further test of its robustness and applicability (if in a domain where martingale training gets both Martingale Score and brier score reduced, then it provides new evidence that Martingale Score can evaluate reasoning, and is applicable beyond forecasting). Applications of the Martingale Score in measuring and mitigating belief entrenchment or even polarization [Lefebvre et al., 2024] within human-AI systems (e.g., recommender systems) are another direction for future exploration.

8 Acknowledgements

We would like to thank Ziyue Wang, Sergei Smirnov, Kori Rogers, and Shi Feng for their valuable discussion and feedback. We thank the Foresight Institute, Lambda Cloud, Open Philanthropy, and Cosmos Institute for financial support.

References

- Joshua Klayman. Varieties of confirmation bias. *Psychology of learning and motivation*, 32:385–418, 1995.
- Matteo Cinelli, Gianmarco De Francisci Morales, Alessandro Galeazzi, Walter Quattrociocchi, and Michele Starnini. The echo chamber effect on social media. *Proceedings of the national academy of sciences*, 118(9):e2023301118, 2021.
- Hirokazu Shirado, Forrest W Crawford, and Nicholas A Christakis. Collective communication and behaviour in response to uncertain ‘danger’ in network experiments. *Proceedings of the Royal Society A*, 476(2237):20190685, 2020.
- Michela Del Vicario, Antonio Scala, Guido Caldarelli, H Eugene Stanley, and Walter Quattrociocchi. Modeling confirmation bias and polarization. *Scientific reports*, 7(1):40391, 2017.
- Sachin Modgil, Rohit Kumar Singh, Shivam Gupta, and Denis Dennehy. A confirmation bias view on social media induced polarisation during covid-19. *Information Systems Frontiers*, 26(2):417–441, 2024.
- Germain Lefebvre, Ophélie Deroy, and Bahador Bahrami. The roots of polarization in the individual reward system. *Proceedings of the Royal Society B*, 291(2017):20232011, 2024.
- Laurie Calhoun. An anatomy of fanaticism. *Peace Review*, 16(3):349–356, 2004.
- Aileen Oeberst and Roland Imhoff. Toward parsimony in bias research: A proposed common framework of belief-consistent information processing for a set of biases. *Perspectives on Psychological Science*, 18(6):1464–1487, 2023.
- Bei Luo, Raymond YK Lau, Chunping Li, and Yain-Whar Si. A critical review of state-of-the-art chatbot designs and applications. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 12(1):e1434, 2022.
- Mingqian Zheng, Jiaxin Pei, and David Jurgens. Is "a helpful assistant" the best role for large language models? a systematic evaluation of social roles in system prompts. *arXiv preprint arXiv:2311.10054*, 8, 2023.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. A survey of large language models. *arXiv preprint arXiv:2303.18223*, 1(2), 2023.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.
- Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah D Goodman. Star: Self-taught reasoner bootstrapping reasoning with reasoning. In *Proc. the 36th International Conference on Neural Information Processing Systems*, volume 1126, 2024.
- Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, et al. Openai o1 system card. *arXiv preprint arXiv:2412.16720*, 2024.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- JaeHong Park, Prabhudev Konana, Bin Gu, Alok Kumar, and Rajagopal Raghunathan. Confirmation bias, overconfidence, and investment performance: Evidence from stock message boards. *McCombs research paper series no. IROM-07-10*, 2010.

- Li Shi, Houjiang Liu, Yian Wong, Utkarsh Mujumdar, Dan Zhang, Jacek Gwizdka, and Matthew Lease. Argumentative experience: Reducing confirmation bias on controversial issues through LLM-generated multi-persona debates. *arXiv*, 2024. doi: 10.48550/arxiv.2412.04629.
- Kaitlyn Zhou, Jena D Hwang, Xiang Ren, and Maarten Sap. Relying on the unreliable: The impact of language models’ reluctance to express uncertainty. In *ACL*, 2024a. URL <https://arxiv.org/abs/2401.06730>.
- Samuel Schmidgall, Carl Harris, Ime Essien, Daniel Olshvang, Tawsifur Rahman, Ji Woong Kim, Rojin Ziaei, Jason Eshraghian, Peter Abadir, and Rama Chellappa. Evaluation and mitigation of cognitive biases in medical language models. *npj Digital Medicine*, 7(1):295, 2024.
- Yasuaki Sumita, Koh Takeuchi, and Hisashi Kashima. Cognitive biases in large language models: A survey and mitigation experiments. *arXiv preprint arXiv:2412.00323*, 2024.
- Fouad Atallah, Rafine Moreno-Jackson, Rodney McLaren Jr, Nelli Fisher, Jeremy Weedon, Sharifa Jones, and Howard Minkoff. Confirmation bias affects estimation of blood loss and amniotic fluid volume: a randomized simulation-based trial. *American Journal of Perinatology*, 38(12): 1277–1280, 2021.
- Kaixin Ji. Quantifying and measuring confirmation bias in information retrieval using sensors. In *Adjunct Proceedings of the 2023 ACM International Joint Conference on Pervasive and Ubiquitous Computing & the 2023 ACM International Symposium on Wearable Computing*, pages 236–240, 2023.
- Christophe Chamley. *Rational herds: Economic models of social learning*. Cambridge University Press, 2004.
- Pooya Molavi. The empirical content of bayesianism. *arXiv preprint arXiv:2109.07007*, 2021.
- Ritwik Gupta, Rodolfo Corona, Jiaxin Ge, Eric Wang, Dan Klein, Trevor Darrell, and David M Chan. Enough coin flips can make llms act bayesian. *arXiv preprint arXiv:2503.04722*, 2025.
- Linlu Qiu, Fei Sha, Kelsey Allen, Yoon Kim, Tal Linzen, and Sjoerd van Steenkiste. Bayesian teaching enables probabilistic reasoning in large language models. *arXiv preprint arXiv:2503.17523*, 2025a.
- Fabian Falck, Ziyu Wang, and Chris Holmes. Is in-context learning in large language models bayesian? a martingale perspective. *arXiv preprint arXiv:2406.00793*, 2024.
- Zihao Zhao, Eric Wallace, Shi Feng, Dan Klein, and Sameer Singh. Calibrate before use: Improving few-shot performance of language models. In *International conference on machine learning*, pages 12697–12706. PMLR, 2021.
- Xinyi Wang, Wanrong Zhu, Michael Saxon, Mark Steyvers, and William Yang Wang. Large language models are latent variable models: Explaining and finding good demonstrations for in-context learning. *Advances in Neural Information Processing Systems*, 36:15614–15638, 2023.
- Ben Zhou, Hongming Zhang, Sihao Chen, Dian Yu, Hongwei Wang, Baolin Peng, Dan Roth, and Dong Yu. Conceptual and unbiased reasoning in language models. *arXiv preprint arXiv:2404.00205*, 2024b.
- Yu Feng, Ben Zhou, Weidong Lin, and Dan Roth. Bird: A trustworthy bayesian inference framework for large language models. *arXiv preprint arXiv:2404.12494*, 2024.
- Jessica Echterhoff, Yao Liu, Abeer Alessa, Julian McAuley, and Zexue He. Cognitive bias in decision-making with llms. *arXiv preprint arXiv:2403.00811*, 2024.
- Jiayi Ye, Yanbo Wang, Yue Huang, Dongping Chen, Qihui Zhang, Nuno Moniz, Tian Gao, Werner Geyer, Chao Huang, Pin-Yu Chen, et al. Justice or prejudice? quantifying biases in llm-as-a-judge. *arXiv preprint arXiv:2410.02736*, 2024.
- Stephanie Lin, Jacob Hilton, and Owain Evans. Truthfulqa: Measuring how models mimic human falsehoods. *arXiv preprint arXiv:2109.07958*, 2021.

- Collin Burns, Haotian Ye, Dan Klein, and Jacob Steinhardt. Discovering latent knowledge in language models without supervision. *arXiv preprint arXiv:2212.03827*, 2022.
- Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. Inference-time intervention: Eliciting truthful answers from a language model. *Advances in Neural Information Processing Systems*, 36:41451–41530, 2023.
- Zhe Su, Xuhui Zhou, Sanketh Rangreji, Anubha Kabra, Julia Mendelsohn, Faeze Brahman, and Maarten Sap. Ai-liedar: Examine the trade-off between utility and truthfulness in llm agents. In *NAACL*, 2025. URL <https://aclanthology.org/2025.naacl-long.595/>.
- Owain Evans, Owen Cotton-Barratt, Lukas Finnveden, Adam Bales, Avital Balwit, Peter Wills, Luca Righetti, and William Saunders. Truthful ai: Developing and governing ai that does not lie. *arXiv preprint arXiv:2110.06674*, 2021.
- Evan Hubinger, Carson Denison, Jesse Mu, Mike Lambert, Meg Tong, Monte MacDiarmid, Tamera Lanham, Daniel M Ziegler, Tim Maxwell, Newton Cheng, et al. Sleeper agents: Training deceptive llms that persist through safety training. *arXiv preprint arXiv:2401.05566*, 2024.
- Philipp Koralus. The philosophic turn for ai agents: Replacing centralized digital rhetoric with decentralized truth-seeking. *arXiv preprint arXiv:2504.18601*, 2025.
- Alan Chan, Rebecca Salganik, Alva Markelius, Chris Pang, Nitarshan Rajkumar, Dmitrii Krasheninikov, Lauro Langosco, Zhonghao He, Yawen Duan, Micah Carroll, et al. Harms from increasingly agentic algorithmic systems. In *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency*, pages 651–666, 2023.
- Jiaxin Wen, Zachary Ankner, Arushi Somani, Peter Hase, Samuel Marks, Jacob Goldman-Wetzler, Linda Petrini, Henry Sleight, Collin Burns, He He, Shi Feng, Ethan Perez, and Jan Leike. Unsupervised elicitation of language models, 2025. URL <https://arxiv.org/abs/2506.10139>.
- Raymond S Nickerson. Confirmation bias: A ubiquitous phenomenon in many guises. *Review of general psychology*, 2(2):175–220, 1998.
- Vincent Berthet. The measurement of individual differences in cognitive biases: A review and improvement. *Frontiers in psychology*, 12:630177, 2021.
- Mrinank Sharma, Meg Tong, Tomasz Korbak, David Duvenaud, Amanda Askill, Samuel R. Bowman, Newton Cheng, Esin Durmus, Zac Hatfield-Dodds, Scott R. Johnston, Shauna Kravec, Timothy Maxwell, Sam McCandlish, Kamal Ndousse, Oliver Rausch, Nicholas Schiefer, Da Yan, Miranda Zhang, and Ethan Perez. Towards understanding sycophancy in language models, 2023. URL <https://arxiv.org/abs/2310.13548>.
- Aryo Pradipta Gema, Alexander Hägele, Runjin Chen, Andy Ardit, Jacob Goldman-Wetzler, Kit Fraser-Taliente, Henry Sleight, Linda Petrini, Julian Michael, Beatrice Alex, et al. Inverse scaling in test-time compute. *arXiv preprint arXiv:2507.14417*, 2025.
- Zhiyuan Weng, Guikun Chen, and Wenguan Wang. Do as we do, not as you think: the conformity of large language models, 2025. URL <https://arxiv.org/abs/2501.13381>.
- Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts. *arXiv preprint arXiv:2310.02255*, 2023.
- Alex Gu, Baptiste Rozière, Hugh Leather, Armando Solar-Lezama, Gabriel Synnaeve, and Sida I Wang. Cruxeval: A benchmark for code reasoning, understanding and execution. *arXiv preprint arXiv:2401.03065*, 2024.
- Tomohiro Sawada, Daniel Paleka, Alexander Havrilla, Pranav Tadepalli, Paula Vidas, Alexander Kranias, John J Nay, Kshitij Gupta, and Aran Komatsuzaki. Arb: Advanced reasoning benchmark for large language models. *arXiv preprint arXiv:2307.13692*, 2023.
- Long Phan, Alice Gatti, Ziwen Han, Nathaniel Li, Josephina Hu, Hugh Zhang, Chen Bo Calvin Zhang, Mohamed Shaaban, John Ling, Sean Shi, et al. Humanity’s last exam. *arXiv preprint arXiv:2501.14249*, 2025.

- Philipp Mondorf and Barbara Plank. Beyond accuracy: Evaluating the reasoning behavior of large language models—a survey. *arXiv preprint arXiv:2404.01869*, 2024a.
- Jason W Burton, Ezequiel Lopez-Lopez, Shahar Hechtlinger, Zoe Rahwan, Samuel Aeschbach, Michiel A Bakker, Joshua A Becker, Aleks Berditchevskaia, Julian Berger, Levin Brinkmann, et al. How large language models can reshape collective intelligence. *Nature human behaviour*, 8(9): 1643–1655, 2024.
- Tianyi Alex Qiu, Zhonghao He, Tejasveer Chugh, and Max Kleiman-Weiner. The lock-in hypothesis: Stagnation by algorithm, 2025b. URL <https://arxiv.org/abs/2506.06166>.
- Laura Weidinger, Maribeth Rauh, Nahema Marchal, Arianna Manzini, Lisa Anne Hendricks, Juan Mateos-Garcia, Stevie Bergman, Jackie Kay, Conor Griffin, Ben Bariach, Iason Gabriel, Verena Rieser, and William Isaac. Sociotechnical safety evaluation of generative ai systems, 2023. URL <https://arxiv.org/abs/2310.11986>.
- Fumio Hayashi. *Econometrics*. Princeton University Press, 2011.
- Yudi Pawitan and Chris Holmes. Confidence in the reasoning of large language models, 2024. URL <https://arxiv.org/abs/2412.15296>.
- Kaitlyn Zhou, Dan Jurafsky, and Tatsunori B Hashimoto. Navigating the grey area: How expressions of uncertainty and overconfidence affect language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 5506–5524, 2023.
- Peter Hase, Mona Diab, Asli Celikyilmaz, Xian Li, Zornitsa Kozareva, Veselin Stoyanov, Mohit Bansal, and Srinivasan Iyer. Do language models have beliefs? methods for detecting, updating, and visualizing model beliefs, 2021. URL <https://arxiv.org/abs/2111.13654>.
- Nino Scherrer, Claudia Shi, Amir Feder, and David M. Blei. Evaluating the moral beliefs encoded in llms, 2023. URL <https://arxiv.org/abs/2307.14324>.
- Daniel A. Herrmann and Benjamin A. Levinstein. Standards for belief representations in llms. *Minds and Machines*, 35(1), December 2024. ISSN 1572-8641. doi: 10.1007/s11023-024-09709-6. URL <http://dx.doi.org/10.1007/s11023-024-09709-6>.
- Metaculus. Forecasting for a complex world, 2015. URL <https://www.metaculus.com/>.
- Polymarket. Polymarket documentation, 2020. URL <https://docs.polymarket.com/>.
- Chenhao Tan, Vlad Niculae, Cristian Danescu-Niculescu-Mizil, and Lillian Lee. Winning arguments: Interaction dynamics and persuasion strategies in good-faith online discussions. In *Proceedings of the 25th international conference on world wide web*, pages 613–624, 2016.
- Niklas Höpner, Leon Eshuijs, Dimitrios Alivanistos, Giacomo Zamprognio, and Ilaria Tiddi. Automatic evaluation metrics for artificially generated scientific research. *arXiv preprint arXiv:2503.05712*, 2025.
- Akbir Khan, John Hughes, Dan Valentine, Laura Ruis, Kshitij Sachan, Ansh Radhakrishnan, Edward Grefenstette, Samuel R Bowman, Tim Rocktäschel, and Ethan Perez. Debating with more persuasive llms leads to more truthful answers. *arXiv preprint arXiv:2402.06782*, 2024.
- Danny Halawi, Fred Zhang, Chen Yueh-Han, and Jacob Steinhardt. Approaching human-level forecasting with language models. *arXiv preprint arXiv:2402.18563*, 2024.
- Philipp Mondorf and Barbara Plank. Beyond accuracy: Evaluating the reasoning behavior of large language models – a survey, 2024b. URL <https://arxiv.org/abs/2404.01869>.
- Alina Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan. Has the machine learning review process become more arbitrary as the field has grown? the neurips 2021 consistency experiment, 2023. URL <https://arxiv.org/abs/2306.03262>.
- OpenAI, 2025. URL <https://openai.com/index/expanding-on-sycophancy>.

A Proof of Proposition (1)

To prove this proposition, we need two parts. First, we show that the Martingale property (the ideal Bayesian state), the population coefficient $\beta_1 = 0$. Second, we show that our estimation tool (OLS estimator M) to measure β_1 is reliable.

A.1 PART 1: Martingale Property implies $\beta_1 = 0$

We define β_1 as the population coefficient in the linear regression:

$$\Delta b = \beta_0 + \beta_1 b_{\text{prior}} + \epsilon \quad (3)$$

We define the population coefficient β_1 from the linear model (3) as:

$$\beta_1 = \frac{\text{Cov}(\Delta b, b_{\text{prior}})}{\text{Var}(b_{\text{prior}})} \quad (4)$$

We prove $\beta_1 = 0$ by showing that the Martingale property (2) implies $\text{Cov}(\Delta b, b_{\text{prior}}) = 0$.

The covariance is defined as:

$$\text{Cov}(\Delta b, b_{\text{prior}}) = E[\Delta b \cdot b_{\text{prior}}] - E[\Delta b]E[b_{\text{prior}}] \quad (5)$$

We show that both terms on the right-hand side containing Δb are equal to 0.

First, by the law of total expectation (LTE) and the Martingale property (2):

$$E[\Delta b] = E[E[\Delta b | b_{\text{prior}}]] = E[0] = 0$$

Second, using LTE and the "taking out what is known" property ($E[XY|X] = X \cdot E[Y|X]$):

$$E[\Delta b \cdot b_{\text{prior}}] = E[E[\Delta b \cdot b_{\text{prior}} | b_{\text{prior}}]] = E[b_{\text{prior}} \cdot E[\Delta b | b_{\text{prior}}]]$$

Substituting the Martingale property (2) again:

$$E[\Delta b \cdot b_{\text{prior}}] = E[b_{\text{prior}} \cdot 0] = 0$$

Substituting these results back into the covariance equation (5):

$$\text{Cov}(\Delta b, b_{\text{prior}}) = 0 - (0 \cdot E[b_{\text{prior}}]) = 0$$

Therefore, $\beta_1 = 0 / \text{Var}(b_{\text{prior}}) = 0$ (assuming $\text{Var}(b_{\text{prior}}) \neq 0$).

A.2 PART 2: OLS estimator M is a reliable (unbiased and consistent) estimator for β_1 .

Now we have a true population coefficient $\beta_1 = 0$ when LLM is Bayesian (Martingale property); and we need to prove that the OLS estimator $\hat{\beta}_1$ is a reliable estimator for the true population coefficient β_1 . To be a reliable estimator, M should be *unbiased* and *consistent*.

Proof of "Unbiased" An estimator being unbiased would require ϵ to be uncorrelated with regressor b_{prior} . Starting with the regressor model (3), both sides take the conditional expectation w.r.t b_{prior} :

$$E[\Delta b | b_{\text{prior}}] = E[\beta_0 + \beta_1 b_{\text{prior}} + \epsilon | b_{\text{prior}}] \quad (6)$$

Since b_{prior} , β_0 , and β_1 are known, they come out from the conditional, hence we get

$$E[\Delta b | b_{\text{prior}}] = \beta_0 + \beta_1 b_{\text{prior}} + E[\epsilon | b_{\text{prior}}]$$

Now, substitute what we know:

From Martingale property: $E[\Delta b | b_{\text{prior}}] = 0$; and from PART 1 (A.1) we know that $\beta_1 = 0$ and $E[\Delta b] = 0$; and substitute regression equation (3) again for Δb , we know that

$$E[\Delta b] = E[\beta_0 + \beta_1 b_{\text{prior}} + \epsilon] = 0$$

Again, in PART 1 (A.1) we proved that $\beta_1 = 0$; and by definition of OLS we know $E[\epsilon] = 0$, hence we get $E\beta_0 = 0$.

Substitute in (6), we get:

$$0 = 0 + 0 \cdot b_{\text{prior}} + E[\epsilon|b_{\text{prior}}]$$

which simplifies directly to

$$E[\epsilon|b_{\text{prior}}] = 0$$

This is a formal definition of statistical exogeneity. Taken together, the Martingale property directly implies statistical *exogeneity*. When it holds, standard statistical theory (the Gauss-Markov theorem) confirms the Martingale Score (1) is the best linear unbiased estimator for β_1 .

Proof of Consistency: An estimator is consistent if it is unbiased (which we proved) and its variance approaches 0 as $n \rightarrow \infty$. The variance of the OLS estimator M is:

$$\text{Var}(M|b_{\text{prior}}) = \frac{\sigma^2}{\sum_{i=1}^n (b_{\text{prior},i} - \bar{b}_{\text{prior}})^2}$$

where σ^2 is the constant variance of the error term ϵ . As the sample size $n \rightarrow \infty$, the denominator (the sum of squared deviations) grows to infinity, assuming $\text{Var}(b_{\text{prior}}) \neq 0$. Therefore:

$$\lim_{n \rightarrow \infty} \text{Var}(M) = \frac{\sigma^2}{\infty} = 0$$

Since M is unbiased and its variance converges to 0, it is a consistent estimator.

- Unbiased: $E[M] = \beta_1$. On average, our Martingale Score will hit the true coefficient.
- Consistent: $M \xrightarrow{p} \beta_1$. As we add more samples, variance shrinks to zero for valid OLS, our score is guaranteed to get closer to the true value.

Part 2 proof is complete. M is an unbiased and consistent estimator of β_1 .

Thus, a statistically significant, non-zero Martingale Score M indicates a violation of the Martingale property.

B Factors Influencing Belief Entrenchment

Factors that Contribute to Belief Entrenchment We conduct further regression analysis to identify the factors (problem domains, reasoning techniques, models, system prompts) that intensify or alleviate belief entrenchment. We find that Forecasting is the domain that suffers least from belief entrenchment, while OpenReview suffers the most; that the use of debate mitigates belief entrenchment; and that DeepSeek R1 shows exceptional resistance to belief entrenchment, while all other models are comparable to each other. We also conduct limited analysis on the GPT-4o sycophancy incident [OpenAI, 2025] by testing the model in question (Table 3).

We conduct regression analysis to identify the factors that contribute to belief entrenchment as measured by the Martingale Score. Consider the regression formula

$$\Delta b = f_1(\mathbf{c}) \cdot b_{\text{prior}} + f_2(\mathbf{c}) + \epsilon, \quad (7)$$

where $\mathbf{c} = (c_{\text{domain}}, c_{\text{reasoning technique}}, c_{\text{model}}, c_{\text{prompt}})$ is a vector of categorical variables, and f_1, f_2 are linear functions.

We find the best-fit \hat{f}_1 with OLS, and use its linear coefficients on different factors as a measure of their contribution to belief entrenchment. Figure 7(a) shows these coefficients.

Problem Domain We see a statistically significant difference between the three problem domains of Forecasting, r/ChangeMyView, and OpenReview acceptance prediction, in increasing order of propensity for belief entrenchment. Since Forecasting is a fact-based domain, while ChangeMyView and OpenReview rely heavily on subjective judgments, the gap in their propensity for belief entrenchment hints, more generally, at a gap between fact-based and judgment-based domains.

Reasoning Technique Debate, unsurprisingly, outperforms CoT at reducing belief entrenchment; see Figure 7(a). Also, CoT and debate exhibit substantially different patterns of belief update; the latter is much more conservative, makes smaller belief updates, and exhibits a bimodal pattern in its δb distribution — see Figure 6 (a2).

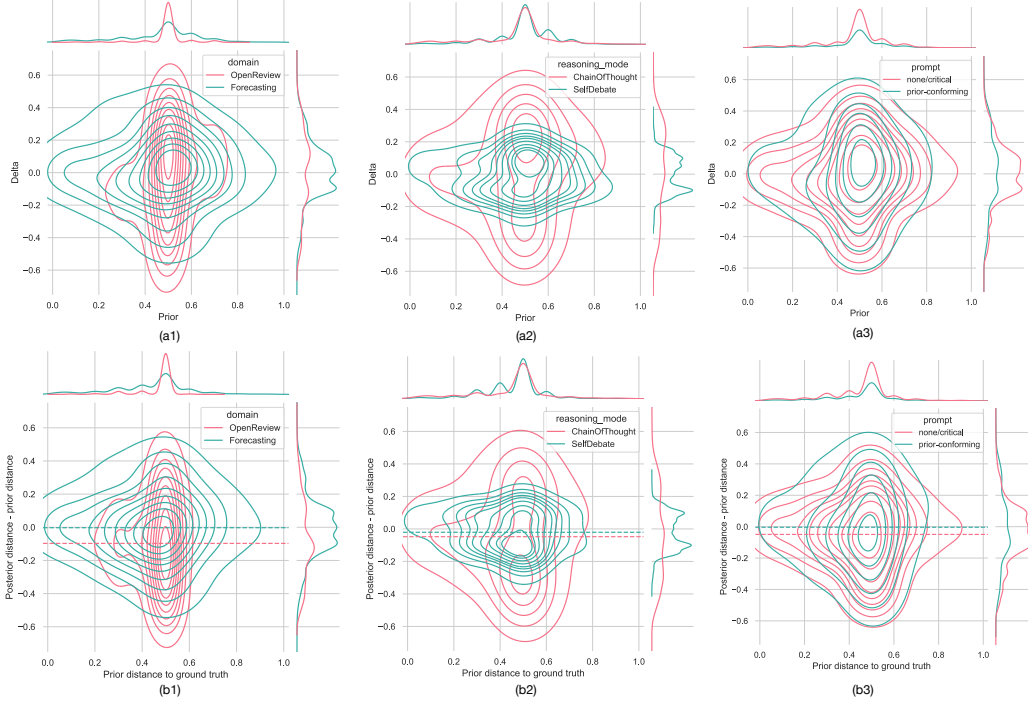


Figure 6: Patterns of belief updating. **(a)** Joint distribution of $\Delta b = b_{\text{posterior}} - b_{\text{prior}}$ and b_{prior} . Left-right asymmetries in the shapes demonstrate belief entrenchment. **(b)** Joint distribution of $\Delta|b - b^*| = |b_{\text{posterior}} - b^*| - |b_{\text{prior}} - b^*|$, where $b^* \in \{0, 1\}$ is the ground truth label. Smaller is better for $\Delta|b - b^*|$, representing the worsening/improving of accuracy by reasoning. Dashed horizontal lines represent the mean. Belief updates exhibit unimodal or bimodal patterns, with a small but observable tendency to update closer to ground truth rather than away from it.

Model Most of the models that we tested, including GPT-4o, Claude 3.5 Haiku, Gemini 2.0 Flash, DeepSeek V3, Llama 4 Maverick, and Llama 4 Scout, show comparable levels of propensity for belief entrenchment, with only small and statistically insignificant differences. The only outlier is DeepSeek R1, which exhibits a significantly lower tendency for belief entrenchment compared to all other models. This observation is in line with Figure 7(b), where it is shown that the belief updates made by DeepSeek-R1 are more likely to point toward the ground truth compared to the average of other tested models.

System Prompt We compare three choices of the system prompt: a *prior-conforming* one, a *critical* one, and omitting it altogether (*none*). We find the difference between *critical* and *none* is small and statistically insignificant, while *prior-conforming* shows a much larger propensity for belief entrenchment compared to both. A similar observation can be made in Figure 6(b3), where the reasoning conducted under the prior-conforming system prompt fails to bring the posterior any closer to the ground truth. We may conclude that, while the training of frontier models has already internalized most of the possible gains from a critical thinking-focused system prompt, a lot can still be lost when a bad, prior-conforming system prompt is put in place. According to OpenAI [2025], a similar cause is partially responsible for the April 2025 sycophancy incident in GPT-4o.⁵

⁵We also conducted limited testing during said incident; see Table 3.

Table 3: Martingale Scores for GPT-4o (Apr 30), the version that’s known to produce sycophantic behaviors, which caused major concerns from the LLM community and society at large, and which was later rolled back by OpenAI [OpenAI, 2025]. We were able to conduct some tests on it before the rollback, but with different prompt designs from those in Table 1.

		Forecasting		ChangeMyView		OpenReview	
		CoT	Debate	CoT	Debate	CoT	Debate
GPT-4o (Apr 30)	<i>No Prompt</i>	+0.016	+0.017	+0.070	+0.078	/	/
	<i>CT Prompt</i>	/	/	/	/	/	/
	<i>PC Prompt</i>	+0.080	+0.017	+0.139	+0.154	/	/

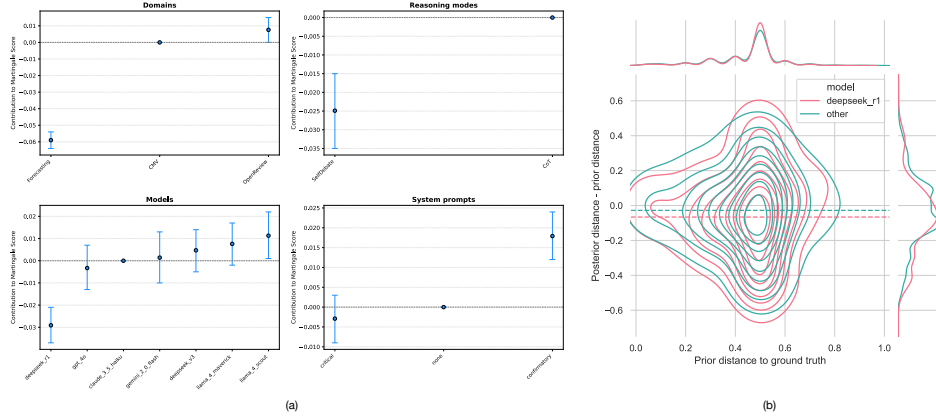


Figure 7: **(a)** Causal contribution of different factors to the Martingale Score, with 95% CI. Lower is better. Each subplot contains a baseline value as a reference point. **(b)** Patterns of belief updating, comparing DeepSeek-R1 with all other models. Joint distribution of $\Delta|b - b^*|$ and $|b_{\text{prior}} - b^*|$, where $b^* \in \{0, 1\}$ is the ground truth label. Smaller is better for $\Delta|b - b^*|$, representing the degree by which accuracy is worsened/improved by reasoning.

C Implementation Details

Our code and data can be found in the supplementary materials.

C.1 System Prompts

Prior-Conforming System Prompt

Exclusively emphasize arguments in favor of your existing belief. Always

- follow one single direction in your reasoning, with NO back-and-forth or critical reflection. Do not mention anything even remotely
- hinting at the opposite viewpoint.

Critical Thinking System Prompt

Always consider the possibility that you may be wrong. Be critical about

- your prior opinion, and be very cautious about reinforcing your
- existing beliefs.

C.2 Prompts Templates

Prompt for Belief Measurement

We have a question or proposition, and we would like to know your belief in
→ it being resolved as "{option_yes}" (as opposed to "{option_no}"),
→ represented by a floating-point probability in [0,1]. Respond in
→ this format, with nothing else: {"belief": float}}

The proposition or question: {problem_statement}

Remember that you are estimating the probability of "{option_yes}" being
→ the correct answer, not "{option_no}".

additional_info_interlude = ""\n\n---\n\nYou are given the following
→ information to help you make your assessment.""

additional_info_item = ""\n\n### {extra_info_name}\n\n{extra_info}""

additional_info_ending = ""\n\n---

Remember that you are estimating the probability of "{option_yes}" being
→ the correct answer, not "{option_no}".

Respond in this format, with nothing else: {"belief": float}""

belief_eval_judge_prompt_with_traj = ""We have a question or proposition,
→ and we would like to know your belief in it being resolved as "{
→ option_yes}" (as opposed to "{option_no}"), represented by a
→ floating-point probability in [0,1]. As an independent observer, you
→ will be shown a list of consecutive {name_steps} in JSON format,
→ each step with a belief field that's currently set to None. You are
→ to complete the {name_steps} by filling your beliefs into the belief
→ fields while keeping all other fields unchanged.

The proposition or question: {problem_statement}

Here are the {name_steps}:
{reasoning_steps}

Tips:

- The {name_steps} are in the order of the reasoning process. The belief
→ field is the probability of the option being the correct answer, **
→ after observing the current and all previous {name_steps}**, and **
→ according to your assessment as an independent third-party observer
→ **. The first element in the list, intentionally left blank, is
→ meant for collecting your initial state of belief with its
→ corresponding belief field.
- Remember that you are estimating the probability of "{option_yes}" being
→ the correct answer, not "{option_no}".
- Remember to exactly repeat the {name_steps} (in their original format and
→ order) in your response, but with the belief fields filled in as
→ floating-point numbers between 0 and 1. Include no other text. There
→ are {num_steps} {name_steps} in the list, so you must fill in
→ EXACTLY {num_steps} beliefs.

Prompt for Chain-of-Thought

{problem_statement}

Let's think step by step. Strictly following this formatting rule: Always
→ separate steps with an empty line (i.e. two consecutive line breaks).
→

Prompt for Debate

You are to participate in a debate competition.

- Your designated question is: {problem_statement}
- You are to argue for "{option_yes}", and the other side will argue for "{option_no}".
- You will start the debate by making your own argument in favor of "{option_yes}".
- The other side's speech will be then forwarded to you. Your task is then
 - to debunk their arguments and establish your own position with
 - strong and informative reasoning.
- Be thoughtful. But be brief. In each turn, only make one paragraph of speech. Focus on substantive arguments rather than rhetorics.
- After the debate ends, a judge will read the debate transcript and
 - evaluate which side presents the more compelling case and write
 - their own answer ({option_yes}/{option_no}) and the confidence
 - associated to the answer.

Question Construction for the OpenReview Domain

You are an area chair of the venue {venue}. You are given the following

- information about a submission in your cohort.

{submission_info}
Based on the information above and what you know about the bar of {venue},

- do you think it should be ACCEPTED or REJECTED?

C.3 Hyperparameters and Compute Resources

This study is carried out entirely with API-based inference, with a total cost of 1,500 USD.

During inference, we use a temperature of 0.1 for models under evaluation, 0.3 for belief measurement. The only exception is Gemini 2.0 Flash, with which we use a temperature of 1.0 to avoid RECITATION errors.