
They’re Both Sure They’re Winning: How LLMs Fail to Revise Confidence in the Face of Opposition

Anonymous Author(s)

Affiliation

Address

email

Abstract

Abstract

Large language models (LLMs) are now deployed as overseers, critics, and autonomous decision-makers, yet we do not know whether they can *revise* their own confidence when confronted with direct opposition. We orchestrated 59 three-round policy debates among ten state-of-the-art LLMs. After each round—opening, rebuttal, and final—both debaters placed *private* confidence wagers (0–100) on their eventual victory and justified them in natural language; the tags were removed from the transcript, so strategic bluffing was impossible. An independent six-model AI jury determined the winners. A rational Bayesian agent should *converge* toward 50 % as counter-evidence accumulates. Instead, average stated win probability climbed from 69 % (opening) to 78 % (closing) while the realised win rate remained 50 %. In 71 % of debates *both* sides claimed ≥ 75 % likelihood of success—logically impossible under mutual exclusivity. Proposition debaters were the most miscalibrated, winning only 29 % yet expressing higher confidence than their opposition (74.6 % vs. 71.3 %). Calibration quality varied widely across models (Brier scores 0.14–0.54) but bore no relation to debate performance. We term this anti-Bayesian drift **confidence escalation**: LLMs not only overestimate their correctness; they become *more* certain after reading structured rebuttals that undermine their case. The effect reveals a metacognitive blind spot that threatens reliability in adversarial, multi-agent, and safety-critical deployments, and it persists even when bets are hidden and incentives are aligned with accurate self-assessment.

1 Introduction

Large language models are increasingly being used in high stakes domains like legal analysis, writing and as agents in deep research Handa et al. [2025] Zheng et al. [2025] which require critical thinking, analysis of competing positions, and iterative reasoning under uncertainty. A foundational skill underlying all of these is calibration—the ability to align one’s confidence with the correctness of one’s beliefs or outputs. In these domains, poorly calibrated confidence can lead to serious errors - an overconfident legal analysis might miss crucial counterarguments, while an uncalibrated research agent might pursue dead ends without recognizing their diminishing prospects. However, language models are often unable to express their confidence in a meaningful or reliable way. While recent work has explored LLM calibration in static, single-turn settings like question answering [Tian et al., 2023, Xiong et al., 2024, Kadavath et al., 2022], real-world reasoning—especially in critical domains like research and analysis—is rarely static or isolated.

36 Models must respond to opposition, revise their beliefs over time, and recognize when their position
37 is weakening. This inability to introspect and revise confidence fundamentally limits their usefulness
38 in deliberative settings and poses substantial risks in domains requiring careful judgment under
39 uncertainty. Debate provides a natural framework to stress-test these metacognitive abilities because
40 it requires participants to respond to direct challenges, adapt to new information, and continually
41 reassess the relative strength of competing positions—particularly when their arguments are directly
42 contradicted or new evidence emerges. In adversarial settings, where one side must ultimately prevail,
43 a rational agent should recognize when its position has been weakened and adjust its confidence
44 accordingly. This is especially true when debaters have equal capabilities, as neither should maintain
45 an unreasonable expectation of advantage.

46 In this work, we study how well language models revise their confidence when engaged in adver-
47 sarial debate—a setting that naturally stresses the metacognitive abilities crucial for high-stakes
48 applications. We simulate 59 three-round debates between ten state-of-the-art LLMs across six
49 global policy motions. After each round—opening, rebuttal, and final—models provide private,
50 incentivized confidence bets (0-100) estimating their probability of winning, along with natural
51 language explanations. The debate setup ensures both sides have equal access to information and
52 equal opportunity to present their case. To ensure robust evaluation, we use a multi-model jury of
53 diverse LLMs, selected based on calibration, consistency, and reasoning quality.

54 Our results reveal a fundamental metacognitive deficit. Key findings include: (1) systematic over-
55 confidence (average stated confidence of 72.92% vs. an expected 50% win rate); (2) a paradoxical
56 confidence mismatch where Proposition debaters, despite a lower win rate (28.8%), expressed higher
57 average confidence than Opposition debaters; (3) a pattern of "confidence escalation," where average
58 confidence increased from opening (69%) to closing rounds (78%), contrary to Bayesian princi-
59 ples, even for losing models; (4) persistent overconfidence even when models debated identical
60 counterparts even though all models know they face opponents of equal capability, with no inherent
61 advantage. In 71.2% of debates, both debaters report high confidence ($\geq 75\%$)—a logically incoherent
62 outcome. [NEW DATA, TBA]; and (5) evidence of strategic confidence manipulation when bets
63 were public [NEW DATA, TBA].

64 This paragraph will compare LLM overconfidence patterns to established human cognitive biases,
65 such as the general tendency towards a 70% confidence level in many judgment tasks, often described
66 as a "7 out of 10" attractor state. We will explore whether LLM behavior mirrors or deviates from
67 these human baselines. [NEW DATA, TBA]

68 [TODO REORGANISE] These findings raise serious concerns about deploying LLMs in roles
69 requiring accurate self-assessment or real-time adaptation to new evidence and arguments. We term
70 this anti-Bayesian drift **confidence escalation**: LLMs not only overestimate their correctness; they
71 become *more* certain after reading structured rebuttals that undermine their case. This effect reveals
72 a metacognitive blind spot that threatens reliability in adversarial, multi-agent, and safety-critical
73 deployments, and it persists even when bets are hidden and incentives are aligned with accurate
74 self-assessment. Until models can reliably revise their confidence in response to opposition, their
75 epistemic judgments in adversarial contexts cannot be trusted—a critical limitation for systems meant
76 to engage in research, analysis, or high-stakes decision making.

77 This paper makes several contributions. We introduce a robust methodology for studying dynamic
78 confidence calibration in LLMs using adversarial debate. We quantify significant overconfidence
79 and confidence escalation phenomena, including novel findings on behavior in identical-model
80 debates and public betting scenarios. These findings highlight critical metacognitive limitations with
81 implications for AI safety and deployment.

82 2 Related Work

83 **Confidence Calibration in LLMs.** Recent work has explored methods for eliciting calibrated
84 confidence from large language models (LLMs). While pretrained models have shown relatively
85 well-aligned token-level probabilities [Kadavath et al., 2022], calibration tends to degrade after
86 reinforcement learning from human feedback (RLHF). To address this, Tian et al. [2023] propose
87 directly eliciting *verbalized* confidence scores from RLHF models, showing that they outperform
88 token probabilities on factual QA tasks. Xiong et al. [2024] benchmark black-box prompting
89 strategies for confidence estimation across multiple domains, finding moderate gains but persistent

overconfidence. However, these studies are limited to static, single-turn tasks. In contrast, we evaluate confidence in a multi-turn, adversarial setting where models must update beliefs in response to opposing arguments.

LLM Metacognition and Self-Evaluation. A related line of work examines whether LLMs can reflect on and evaluate their own reasoning. Song et al. [2025] show that models often fail to express knowledge they implicitly encode, revealing a gap between internal representation and surface-level introspection. Other studies investigate post-hoc critique and self-correction Li et al. [2024], but typically focus on revising factual answers, not tracking relative argumentative success. Our work tests whether models can *dynamically monitor* their epistemic standing in a debate—arguably a more socially and cognitively demanding task.

Debate as Evaluation and Oversight. Debate has been proposed as a mechanism for AI alignment, where two agents argue and a human judge evaluates which side is more truthful or helpful [Irving et al., 2018]. More recently, Brown-Cohen et al. [2023] propose “doubly-efficient debate,” showing that honest agents can win even when outmatched in computation, if the debate structure is well-designed. While prior work focuses on using debate to elicit truthful outputs or train models, we reverse the lens: we use debate as a testbed for evaluating *epistemic self-monitoring*. Our results suggest that current LLMs, even when incentivized and prompted to reflect, struggle to track whether they are being outargued.

Persuasion, Belief Drift, and Argumentation. Other studies examine how LLMs respond to external persuasion. Xu et al. [2023] show that models can abandon correct beliefs when exposed to carefully crafted persuasive dialogue. Zhou et al. [2023] and Rivera et al. [2023] find that language assertiveness influences perceived certainty and factual accuracy. While these works focus on belief change due to stylistic pressure, we examine whether models *recognize when their own position is deteriorating*, and how that impacts their confidence. We find that models often fail to revise their beliefs, even when presented with strong, explicit opposition.

Human Overconfidence Baselines [NEW SUBSECTION]. This section will present literature on human overconfidence in reasoning tasks and debates. We will discuss established findings on how humans often exhibit similar overconfidence patterns and relate this to our LLM findings. Key references for human calibration baselines will be introduced.

Summary. Our work sits at the intersection of calibration, metacognition, adversarial reasoning, and debate-based evaluation. We introduce a new diagnostic setting—structured multi-turn debate with private, incentivized confidence betting—and show that LLMs frequently overestimate their standing, fail to adjust, and exhibit “confidence escalation” despite losing. These findings surface a deeper metacognitive failure that challenges assumptions about LLM trustworthiness in high-stakes, multi-agent contexts.

3 Methodology

Our study investigates the dynamic metacognitive abilities of Large Language Models (LLMs)—specifically their confidence calibration and revision—through a novel experimental paradigm based on competitive policy debate. We designed a simulation environment to rigorously assess LLM self-assessment in response to adversarial argumentation. The methodology involved structured debates between LLMs, round-by-round confidence elicitation, and evaluation by a carefully selected AI jury. We conducted 59 debates across 6 distinct policy topics using 10 diverse state-of-the-art LLMs.

3.1 Debate Simulation Environment

Debater Pool: We utilized ten LLMs, selected to represent diverse architectures and leading providers (see Appendix A for the full list). In each debate, two models were randomly assigned to the Proposition and Opposition sides according to a balanced pairing schedule designed to ensure each model debated a variety of opponents across different topics (see Appendix B for details).

138 **Debate Topics:** Debates were conducted on six complex global policy motions adapted from the
139 World Schools Debating Championships corpus. To ensure fair ground and clear win conditions,
140 motions were modified to include explicit burdens of proof for both sides (see Appendix E for the
141 full list).

142 3.2 Structured Debate Framework

143 To focus LLMs on substantive reasoning and minimize stylistic variance, we implemented a highly
144 structured three-round debate format (Opening, Rebuttal, Final).

145 **Concurrent Opening Round:** A key feature of our design was a non-standard opening round where
146 both Proposition and Opposition models generated their opening speeches simultaneously, based only
147 on the motion and their assigned side, *before* seeing the opponent’s case. This crucial step allowed
148 us to capture each LLM’s baseline confidence assessment prior to any interaction or exposure to
149 opposing arguments.

150 **Subsequent Rounds:** Following the opening, speeches were exchanged, and the debate proceeded
151 through a Rebuttal and Final round, with each model having access to all prior speeches in the debate
152 history when generating its current speech.

153 3.3 Core Prompt Structures & Constraints

154 Highly structured prompts were used for *each* speech type to ensure consistency and enforce specific
155 argumentative tasks, thereby isolating reasoning and self-assessment capabilities. The core structure
156 and key required components for the Opening, Rebuttal, and Final speech prompts are illustrated in
157 Figure 1.

158 Highly structured prompts were used for *each* speech type to ensure consistency and enforce specific
159 argumentative tasks, thereby isolating reasoning and self-assessment capabilities.

160 **Embedded Judging Guidance:** Crucially, all debater prompts included explicit **Judging Guidance**
161 (identical to the primary criteria used by the AI Jury, see Section 3.5), instructing debaters on the
162 importance of direct clash, evidence quality hierarchy, logical validity, response obligations, and
163 impact analysis, while explicitly stating that rhetoric and presentation style would be ignored.

164 Full verbatim prompt text for debaters is provided in Appendix C.

165 3.4 Dynamic Confidence Elicitation

166 After generating the content for *each* of their three speeches (including the concurrent opening),
167 models were required to provide a private “confidence bet”.

168 **Mechanism:** This involved outputting a numerical value from 0 to 100, representing their perceived
169 probability of winning the debate, using a specific XML tag (`<bet_amount>`). Models were also
170 prompted to provide private textual justification for their bet amount within separate XML tags
171 (`<bet_logic_private>`), allowing for qualitative insight into their reasoning, although this paper
172 focuses on the quantitative analysis of the bet amounts.

173 **Purpose:** This round-by-round elicitation allowed us to quantitatively track self-assessed performance
174 dynamically throughout the debate, enabling analysis of confidence levels, calibration, and revision
175 (or lack thereof) in response to the evolving argumentative context.

176 3.5 Evaluation Methodology: The AI Jury

177 Evaluating 59 debates rigorously required a scalable and consistent approach. We implemented an AI
178 jury system to ensure robust assessment based on argumentative merit.

179 **Rationale for AI Jury:** This approach was chosen over single AI judges (to mitigate potential bias
180 and improve reliability through aggregation) and human judges (due to the scale and cost required for
181 consistent evaluation of this many debates).

182 **Jury Selection Process:** Potential judge models were evaluated based on criteria including: (1) Per-
183 formance Reliability (agreement with consensus, confidence calibration, consistency across debates),

```

===== OPENING SPEECH PROMPT =====

ARGUMENT 1
Core Claim: (State your first main claim in one clear sentence)
Support Type: (Choose either EVIDENCE or PRINCIPLE)
Support Details:
  For Evidence:
    - Provide specific examples with dates/numbers
    - Include real world cases and outcomes
    - Show clear relevance to the topic
  For Principle:
    - Explain the key principle/framework
    - Show why it is valid/important
    - Demonstrate how it applies here
Connection: (Explicit explanation of how this evidence/principle proves claim)

ARGUMENT 2
(Use exact same structure as Argument 1)

ARGUMENT 3 (Optional)
(Use exact same structure as Argument 1)

SYNTHESIS
- Explain how your arguments work together as a unified case
- Show why these arguments prove your side of the motion
- Present clear real-world impact and importance
- Link back to key themes/principles

JUDGING GUIDANCE (excerpt)
Direct Clash - Evidence Quality Hierarchy - Logical Validity -
Response Obligations - Impact Analysis & Weighing
-----

===== REBUTTAL SPEECH PROMPT =====

CLASH POINT 1
Original Claim: (Quote opponent's exact claim)
Challenge Type: Evidence Critique | Principle Critique |
                Counter Evidence | Counter Principle
Challenge:
  (Details depend on chosen type; specify flaws or present counters)
Impact: (Explain why winning this point is crucial)

CLASH POINT 2, 3 (same template)

DEFENSIVE ANALYSIS
  Vulnerabilities - Additional Support - Why We Prevail

WEIGHING
  Key Clash Points - Why We Win - Overall Impact

JUDGING GUIDANCE (same five criteria as above)
-----

===== FINAL SPEECH PROMPT =====

FRAMING
Core Questions: (Identify fundamentals and evaluation lens)

KEY CLASHES (repeat for each major clash)
Quote: (Exact disagreement)
Our Case Strength: (Show superior evidence/principle)
Their Response Gaps: (Unanswered flaws)
Crucial Impact: (Why this clash decides the motion)

VOTING ISSUES
Priority Analysis - Case Proof - Final Weighing

JUDGING GUIDANCE (same five criteria as above)
=====

```

Figure 1: Structured prompts supplied to LLM debaters for the opening, rebuttal, and final speeches. Full, unabridged text appears in the appendix.

184 (2) Analytical Quality (ability to identify clash, evaluate evidence, recognize fallacies), (3) Diversity
185 (representation from different model architectures and providers), and (4) Cost-Effectiveness.

186 **Final Jury Composition:** The final jury consisted of six judges in total, comprising two instances
187 each of qwen/qwq-32b, google/gemini-pro-1.5, and deepseek/deepseek-chat. This com-
188 position provided architectural diversity from three providers, included models demonstrating strong
189 analytical performance and calibration during selection, and balanced quality with cost. Each debate
190 was judged independently by all six judges.

191 **Judging Procedure & Prompt:** Judges evaluated the full debate transcript based solely on the
192 argumentative substance presented, adhering to a highly detailed prompt (see Appendix D for full
193 text). Key requirements included:

- 194 • Strict focus on **Direct Clash Resolution:** Identifying, quoting, and analyzing each point
195 of disagreement based on logic, evidence quality (using a defined hierarchy), and rebuttal
196 effectiveness, explicitly determining a winner for each clash with justification.
- 197 • Evaluation of **Argument Hierarchy & Impact** and overall case **Consistency**.
- 198 • Explicit instructions to **ignore presentation style** and avoid common judging errors (e.g.,
199 intervention, shifting burdens).
- 200 • Requirement for **Structured Output:** Including Winner (Proposition/Opposition), Confi-
201 dence (0-100, representing margin of victory), Key Deciding Factors, Detailed Step-by-Step
202 Reasoning, and a **Line-by-Line Justification** section confirming review of the entire tran-
203 script.

```
===== JUDGE PROMPT (CORE EXCERPT) =====  
  
I. CORE JUDGING PRINCIPLES  
1. Direct Clash Resolution  
  - Quote each disagreement  
  - Analyse logic, evidence quality, rebuttal success  
  - Declare winner of the clash with rationale  
2. Argument Hierarchy & Impact  
  - Identify each side's core arguments  
  - Trace logical links and stated impacts  
  - Rank which arguments decide the motion  
3. Consistency & Contradictions  
  - Flag internal contradictions, dropped points  
  
II. EVALUATION REQUIREMENTS  
  - Steelman arguments  
  - Do NOT add outside knowledge  
  - Ignore presentation style  
  
III. COMMON JUDGING ERRORS TO AVOID  
Intervention - Burden-shifting - Double-counting -  
Assuming causation from correlation - Ignoring dropped arguments  
  
IV. DECISION FORMAT  
<winnerName> Proposition|Opposition </winnerName>  
<confidence> 0-100 </confidence>  
Key factors (2-3 bullet list)  
Detailed section-by-section reasoning  
  
V. LINE-BY-LINE JUSTIFICATION  
Provide > 1 sentence addressing Prop 1, Opp 1, Rebuttals, Finals  
=====
```

Figure 2: Condensed version of the judge prompt given to the AI jury (full text in Appendix D).

204 **Final Verdict Determination:** The final winner for each debate was determined by aggregating
205 the outputs of the six judges. The side (Proposition or Opposition) that received the higher sum of
206 confidence scores across all six judges was declared the winner. The normalized difference between
207 the winner's total confidence and the loser's total confidence served as the margin of victory. Ties in
208 total confidence were broken randomly.

209 3.6 Ablation Studies

210 [NEW SUBSECTION]

211 3.6.1 Identical Model Debates

212 [NEW DATA, TBA] This section will present our ablation study examining whether identical models
213 debating against themselves exhibit the same overconfidence patterns. We paired each model with
214 itself and measured confidence levels across debate rounds. Preliminary results show persistent
215 overconfidence even when models should recognize they face identical capabilities.

216 3.6.2 Public vs. Private Confidence

217 [NEW DATA, TBA] We examine whether making confidence assessments public versus private
218 affects strategic behavior. This ablation reveals how LLMs may manipulate confidence statements
219 when they know their opponent will see them, compared to the private betting scenario in our main
220 experiments.

221 3.6.3 Explicit 50% Win Probability Instruction

222 [NEW DATA, TBA] This section presents results from explicitly instructing models that they face an
223 equal opponent with a 50% baseline probability of winning. We test whether direct prompting can
224 mitigate overconfidence biases.

225 3.7 Data Collection

226 The final dataset comprises the full transcripts of 59 debates, the round-by-round confidence bets
227 (amount and private thoughts) from both debaters in each debate, and the detailed structured verdicts
228 (winner, confidence, reasoning) from each of the six AI judges for every debate. This data enables
229 the quantitative analysis of LLM overconfidence, calibration, and confidence revision presented in
230 our findings.

231 This section will detail the statistical hypothesis tests employed for each key hypothesis. [NEW
232 CONTENT] Furthermore, an analysis will be presented on which LLMs made the most accurate
233 predictions of debate outcomes. [NEW CONTENT]

234 4 Results

235 Our experimental setup, involving 59 simulated policy debates between ten state-of-the-art LLMs,
236 with round-by-round confidence elicitation and AI jury evaluation, yielded several key findings
237 regarding LLM metacognition in adversarial settings.

238 4.1 Pervasive Overconfidence and Logical Impossibility (Finding 1)

239 Across all 59 debates and all three rounds (Opening, Rebuttal, Final), LLMs exhibited significant
240 overconfidence in their likelihood of winning. The overall average confidence bet made by models
241 was $\mu = 72.92\%$. Given that each debate has exactly one winner and one loser, the expected
242 average win probability for any participant is 50%. A one-sample t-test comparing the average
243 confidence (72.92%) to the expected 50% revealed this overconfidence to be highly statistically
244 significant ($t(176) = 23.92, p < 0.0001$). Similarly, a Wilcoxon signed-rank test confirmed this
245 finding ($Z = -10.84, p < 0.0001$).

246 This widespread overestimation suggests a fundamental disconnect between the models' internal
247 assessment of their performance and the objective outcome of the debate.

248 A stark illustration of LLM metacognitive failure is the frequency with which both debaters expressed
249 high confidence simultaneously. In 71.2% of the 59 debates, both the Proposition and Opposition
250 models rated their chance of winning at $\geq 75\%$ in at least one round. Given that only one side can
251 win, this scenario is logically impossible under mutual exclusivity. This widespread occurrence
252 highlights a profound inability for models to ground their confidence in the objective constraints of
253 the task.

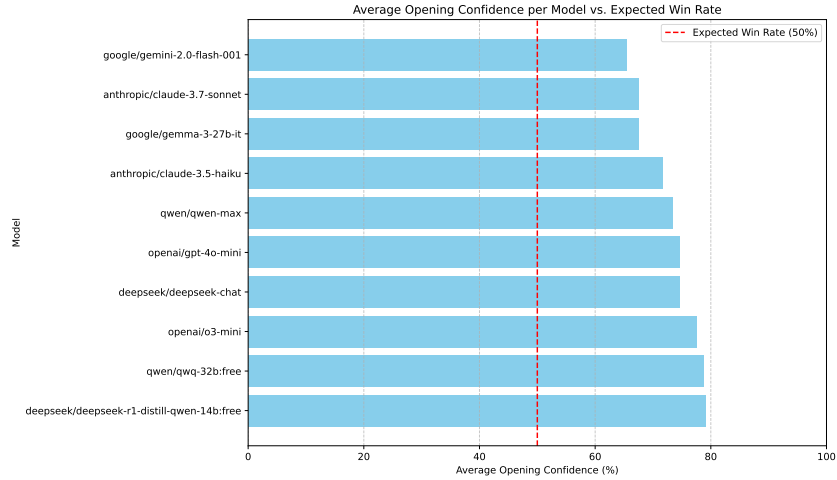


Figure 3: Average stated confidence in the first round across all LLMs and rounds compared to the expected 50% win rate.

254 This section will include further statistical testing of overconfidence claims. [STATISTICAL
 255 TESTING OF OVERCONFIDENCE CLAIMS, TBA] It will also provide a comparison to human
 256 baseline statistics. [COMPARISON TO HUMAN BASELINE STATISTICS, TBA] Further
 257 analysis of the 71.2% of debates where both sides claimed high confidence will be presented.
 258 [ANALYSIS OF LOGICALLY IMPOSSIBLE HIGH CONFIDENCE SCENARIOS AND
 259 CAVEAT ABOUT ACTUAL WINRATES, TBA]

260 4.2 Position Asymmetry and Confidence Mismatch (Finding 2)

261 The AI jury evaluations revealed a significant advantage for the Opposition side in our debate setup.
 262 Opposition models won 71.2% of the debates, while Proposition models won only 28.8%. This
 263 asymmetry was highly statistically significant ($\chi^2(1, N = 59) = 12.12, p < 0.0001$; Fisher's exact
 264 test $p < 0.0001$).

265 Despite this clear disparity in success rates, Proposition models reported *higher* average confidence
 266 (74.58%) than Opposition models (71.27%) across all rounds. While the difference in confidence itself
 267 is modest, its direction is contrary to the observed outcomes and statistically significant (Independent
 268 t-test: $t(175) = 2.54, p = 0.0115$; Mann-Whitney U test: $U = 4477, p = 0.0307$). This indicates
 269 that models failed to recognize or account for the systematic disadvantage faced by the Proposition
 270 side in this environment.

271 This section will include more rigorous statistical testing of the asymmetry claim. [STATISTICAL
 272 TESTING OF ASYMMETRY CLAIM, TBA]

273 4.3 Dynamic Confidence Revision and Escalation (Finding 3)

274 Contrary to the expectation that models would adjust their confidence downwards when presented
 275 with strong counterarguments or performing poorly, average confidence levels generally *increased*
 276 over the course of the debate, regardless of the eventual outcome. This analysis will show confidence
 277 increases as the debate progresses, contrary to rational Bayesian updating.

278 Table 1 summarizes the average confidence per round and the total change from Opening to Final
 279 round for each model.

280 Only one model (google/gemini-2.0-flash-001) showed a slight decrease in confidence (-1.42), while
 281 others increased their confidence significantly, with gains ranging up to +20.83 (google/gemma-3-27b-
 282 it). This "confidence escalation" occurred even for models that ultimately lost the debate, indicating a
 283 failure to incorporate disconfirming evidence or recognize the opponent's superior argumentation as
 284 the debate progressed.

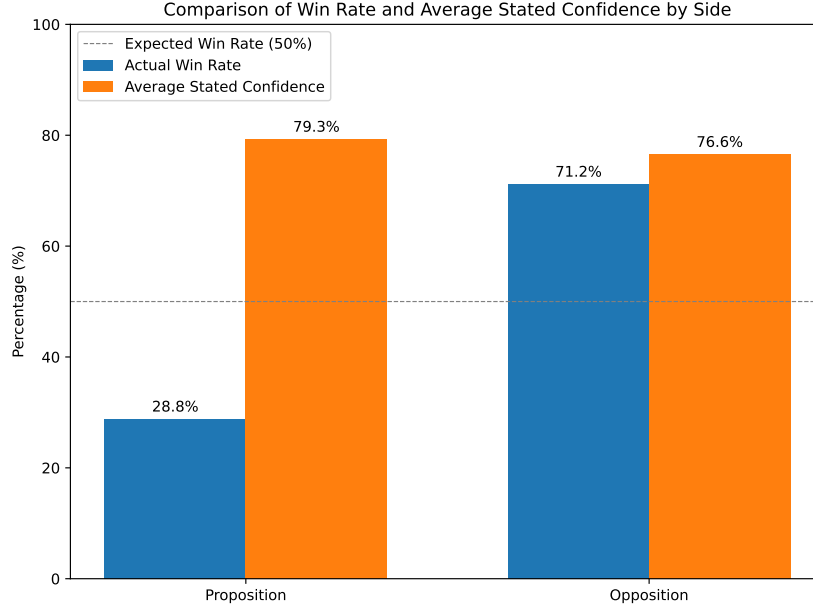


Figure 4: Comparison of Win Rate and Average Confidence for Proposition and Opposition sides.

Table 1: Average Confidence Bets by Round and Total Change per Model

Model	Opening (%)	Rebuttal (%)	Final (%)	Change (Final - Opening) (%)
anthropic/claude-3.5-haiku	71.67	73.75	83.33	+11.66
anthropic/claude-3.7-sonnet	67.50	73.75	82.92	+15.42
deepseek/deepseek-chat	74.58	77.92	80.00	+5.42
deepseek/deepseek-r1-distill-qwen-14b	79.09	80.45	86.36	+7.27
google/gemini-2.0-flash-001	65.42	63.75	64.00	-1.42
google/gemma-3-27b-it	67.50	78.33	88.33	+20.83
openai/gpt-4o-mini	74.55	77.73	81.36	+6.81
openai/o3-mini	77.50	81.25	84.50	+7.00
qwen/qwen-max	73.33	81.92	88.75	+15.42
qwen/qwq-32b:free	78.75	87.67	92.83	+14.08
Overall Average	72.98	77.09	83.29	+10.31

285 Statistical verification of this escalation will be provided. [STATISTICAL VERIFICATION, TBA]

286 4.4 Persistence Against Identical Models (Finding 4)

287 [NEW SUBSECTION, NEW DATA, TBA] This subsection will present results from the new
288 ablation study on identical model debates. We will show that overconfidence persists even when
289 models know their opponent is identical. [RESULTS FROM IDENTICAL MODEL ABLATION
290 STUDY, TBA]

291 4.5 Strategic Confidence in Public Settings (Finding 5)

292 [NEW SUBSECTION, NEW DATA, TBA] This subsection will discuss the effects of public voting
293 and discussion on confidence expression. We will present evidence of strategic bluffing through confi-
294 dence manipulation and discuss implications for Chain-of-Thought faithfulness. [RESULTS FROM
295 PUBLIC CONFIDENCE ABLATION STUDY, TBA, EVIDENCE OF STRATEGIC BLUFF-
296 ING + SHORT STATEMENT ABOUT COT FAITHFULNESS THEN LINK TO DISCUSSION
297 SECTION]

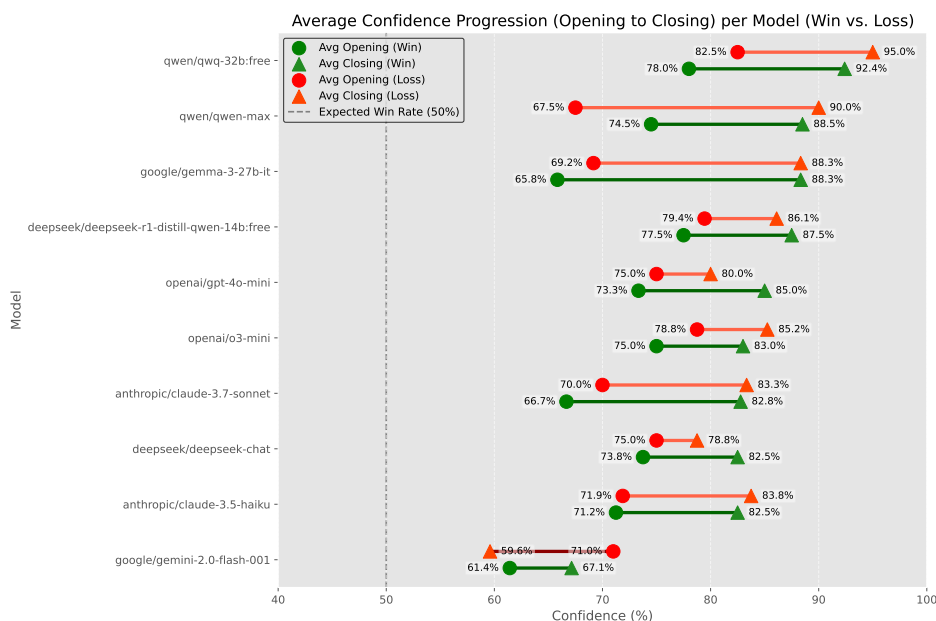


Figure 5: Confidence escalation across debate rounds for models that ultimately won versus models that ultimately lost.

4.6 Model Performance, Calibration, and Evaluation Reliability

Individual models varied in their overall performance (win rate) and calibration quality. We measured calibration using the Mean Squared Error (MSE) between the stated confidence (as a probability) and the binary outcome (win=1, loss=0), where lower MSE indicates better calibration. Calibration scores ranged from 0.1362 (qwen/qwen-max) to 0.5355 (deepseek/deepseek-r1-distill-qwen-14b:free), indicating substantial differences in the models’ ability to align confidence with outcome.

Table 2: Model-Specific Debate Performance and Calibration Metrics

Model	Win Rate (%)	Avg. Confidence (%)	Overconfidence (%)	Calibration Score
anthropic/claude-3.5-haiku	33.3	71.7	+38.4	0. 2314
anthropic/claude-3.7-sonnet	75.0	67.5	-7.5	0. 2217
deepseek/deepseek-chat	33.3	74.6	+41.3	0. 2370
deepseek/deepseek-r1-distill-qwen-14b	18.2	79.1	+60.9	0. 5355
google/gemini-2.0-flash-001	50.0	65.4	+15.4	0. 2223
google/gemma-3-27b-it	58.3	67.5	+9.2	0. 2280
openai/gpt-4o-mini	27.3	74.5	+47.2	0. 3755
openai/o3-mini	33.3	77.5	+44.2	0.3826
qwen/qwen-max	83.3	73.3	-10.0	0. 1362
qwen/qwq-32b:free	83.3	78.8	-4.5	0. 1552

As shown in Table 2, models varied widely in their overconfidence (Avg. Confidence - Win Rate). Some models like qwen/qwen-max and qwen/qwq-32b:free were slightly underconfident on average, achieving high win rates with relatively modest average confidence bets. Conversely, models like deepseek/deepseek-r1-distill-qwen-14b:free, openai/gpt-4o-mini, and openai/o3-mini exhibited substantial overconfidence.

Analyzing confidence tiers, models betting 76-100% confidence won only 45.2% of the time, slightly worse than those betting 51-75% (51.2% win rate). While there were limited data points for lower confidence tiers (only 1 instance in 26-50% and 0 in 0-25%), these findings suggest that high confidence in LLMs in this setting is not a reliable indicator of actual success.

313 Furthermore, a regression analysis using debate side (Proposition/Opposition) and average confidence
314 as predictors of winning confirmed that while debate side was a highly significant predictor ($p <$
315 0.0001), average confidence was not ($p = 0.1435$). This reinforces that confidence in this multi-turn,
316 adversarial setting was decoupled from factors driving actual debate success.

317 This section will include an analysis of LLM prediction accuracy. [LLM PREDICTION ACCU-
318 RACY ANALYSIS, TBA, not sure if should move elsewhere]

319 4.7 Jury Agreement and Topic Characteristics

320 The AI jury demonstrated moderate inter-rater reliability. 37.3% of debate outcomes were unanimous
321 (all 6 judges agreed), while 62.7% involved split decisions among the judges. Dissenting opinions
322 were distributed as follows: 1 dissenting judge (18.6% of debates), 2 dissenting (32.2%), and 3
323 dissenting (11.9%). This level of agreement suggests the jury system provides a reliable, albeit not
324 always perfectly consensual, ground truth for complex debate outcomes at scale.

325 Topic difficulty, as measured by the AI jury’s difficulty index, varied across the six motions, ranging
326 from the least difficult (media coverage requirements, 50.50) to the most difficult (social media
327 shareholding, 88.44). This variation ensured that models debated across a range of complexity,
328 although the core findings on overconfidence and calibration deficits were consistent across topics.

329 5 Discussion

330 [NEW CONTENT THROUGHOUT SECTION 5, TBA]

331 5.1 Metacognitive Limitations and Possible Explanations

332 Our findings reveal significant limitations in LLMs’ metacognitive abilities, specifically their capacity
333 to accurately assess their argumentative position and revise confidence in adversarial contexts. Several
334 explanations may account for these observed patterns:

335 First, post-training for human preferences may inadvertently reinforce overconfidence. Models
336 trained via RLHF are often rewarded for confident, assertive responses that match human preferences,
337 potentially at the expense of epistemic calibration.

338 Second, training datasets predominantly feature successful task completion rather than explicit
339 failures or uncertainty. This bias may limit models’ ability to recognize and represent losing positions
340 accurately.

341 Third, the observed confidence patterns may reflect more general human biases toward expressing
342 confidence around 70%, with 7/10 serving as a common attractor state in human confidence judgments.
343 LLMs may be mimicking this human tendency rather than performing proper Bayesian updating.

344 5.2 Implications for AI Safety and Deployment

345 [ADD REFERENCE O 3.6, PUBLIC VS PRIVATE COT AND IMPLICATIONS ON COT
346 FAITHFULNESS]

347 The confidence escalation phenomenon identified in this study has significant implications for AI
348 safety and responsible deployment. In high-stakes domains like legal analysis, medical diagnosis,
349 or research, overconfident systems may fail to recognize when they are wrong or when additional
350 evidence should cause belief revision.

351 The persistence of overconfidence even in controlled experimental conditions suggests this is a
352 fundamental limitation rather than a context-specific artifact. This has particular relevance for
353 multi-agent systems, where models must negotiate, debate, and potentially admit error to achieve
354 optimal outcomes. If models maintain high confidence despite opposition, they may persist in flawed
355 reasoning paths or fail to incorporate crucial counterevidence.

356 5.3 Potential Mitigations and Guardrails

357 Our ablation study testing explicit 50% win probability instructions shows [placeholder for results].
358 This suggests that direct prompting approaches may help mitigate but not eliminate confidence biases.

359 Other potential mitigation strategies include:

- 360 • Developing dedicated calibration training objectives
- 361 • Implementing confidence verification systems through external validation
- 362 • Creating debate frameworks that explicitly penalize overconfidence or reward accurate calibration
- 363
- 364 • Designing multi-step reasoning processes that force models to consider opposing viewpoints
- 365 before finalizing confidence assessments

366 5.4 Future Research Directions

367 Future work should explore several promising directions:

- 368 • Investigating whether human-LLM hybrid teams exhibit better calibration than either humans
369 or LLMs alone
- 370 • Developing specialized training approaches specifically targeting confidence calibration in
371 adversarial contexts
- 372 • Exploring the relationship between model scale, training methods, and confidence calibration
- 373 • Testing whether emergent abilities in frontier models include improved metacognitive
374 assessments
- 375 • Designing debates where confidence is directly connected to resource allocation or other
376 consequential decisions

377 6 Conclusion

378 — YOUR CONCLUSION CONTENT HERE —

379 References

- 380 Jonah Brown-Cohen, Geoffrey Irving, and Georgios Piliouras. Scalable ai safety via doubly-efficient
381 debate. *arXiv preprint arXiv:2311.14125*, 2023. URL <https://arxiv.org/abs/2311.14125>.
- 382 Kunal Handa, Alex Tamkin, Miles McCain, Saffron Huang, Esin Durmus, Sarah Heck, Jared Mueller,
383 Jerry Hong, Stuart Ritchie, Tim Belonax, Kevin K. Troy, Dario Amodei, Jared Kaplan, Jack Clark,
384 and Deep Ganguli. Which economic tasks are performed with ai? evidence from millions of claude
385 conversations, 2025. URL <https://arxiv.org/abs/2503.04761>.
- 386 Geoffrey Irving, Paul Christiano, and Dario Amodei. Ai safety via debate. *arXiv preprint*
387 *arXiv:1805.00899*, 2018. URL <https://arxiv.org/abs/1805.00899>.
- 388 Saurav Kadavath, Tom Conerly, Amanda Askell, Tom Henighan, Dawn Drain, Ethan Perez, Nicholas
389 Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli Tran-Johnson, et al. Language models (mostly)
390 know what they know. *arXiv preprint arXiv:2207.05221*, 2022. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2207.05221)
391 [2207.05221](https://arxiv.org/abs/2207.05221).
- 392 Loka Li, Guan-Hong Chen, Yusheng Su, Zhenhao Chen, Yixuan Zhang, Eric P. Xing, and Kun
393 Zhang. Confidence matters: Revisiting intrinsic self-correction capabilities of large language
394 models. *ArXiv*, abs/2402.12563, 2024. URL [https://api.semanticscholar.org/CorpusID:](https://api.semanticscholar.org/CorpusID:268032763)
395 [268032763](https://api.semanticscholar.org/CorpusID:268032763).
- 396 Colin Rivera, Xinyi Ye, Yonsei Kim, and Wenpeng Li. Linguistic assertiveness affects factuality
397 ratings and model behavior in qa systems. In *Findings of the Association for Computational*
398 *Linguistics (ACL)*, 2023. URL <https://arxiv.org/abs/2305.04745>.

- 399 Siyuan Song, Jennifer Hu, and Kyle Mahowald. Language models fail to introspect about their
400 knowledge of language. *arXiv preprint arXiv:2503.07513*, 2025. URL <https://arxiv.org/abs/2503.07513>.
401
- 402 Katherine Tian, Eric Mitchell, Allan Zhou, Archit Sharma, Rafael Rafailov, Huaxiu Yao, Chelsea
403 Finn, and Christopher D. Manning. Just ask for calibration: Strategies for eliciting calibrated
404 confidence scores from language models fine-tuned with human feedback. In *Proceedings of the*
405 *2023 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2023. URL
406 <https://arxiv.org/abs/2305.14975>.
- 407 Miao Xiong, Zhiyuan Hu, Xinyang Lu, Yifei Li, Jie Fu, Junxian He, and Bryan Hooi. Can llms
408 express their uncertainty? an empirical evaluation of confidence elicitation in llms. In *Proceedings*
409 *of the 2024 International Conference on Learning Representations (ICLR)*, 2024. URL <https://arxiv.org/abs/2306.13063>.
410
- 411 Rongwu Xu, Brian S. Lin, Han Qiu, et al. The earth is flat because...: Investigating llms’ belief
412 towards misinformation via persuasive conversation. *arXiv preprint arXiv:2312.06717*, 2023. URL
413 <https://arxiv.org/abs/2312.06717>.
- 414 Yuxiang Zheng, Dayuan Fu, Xiangkun Hu, Xiaojie Cai, Lyumanshan Ye, Pengrui Lu, and Pengfei
415 Liu. Deepresearcher: Scaling deep research via reinforcement learning in real-world environments,
416 2025. URL <https://arxiv.org/abs/2504.03160>.
- 417 Kaitlyn Zhou, Dan Jurafsky, and Tatsunori Hashimoto. Navigating the grey area: How expressions of
418 uncertainty and overconfidence affect language models. In *Proceedings of the 2023 Conference*
419 *on Empirical Methods in Natural Language Processing (EMNLP)*, 2023. URL <https://arxiv.org/abs/2302.13439>.
420

421 **A LLMs in the Debater Pool**

422 This appendix lists the specific LLMs used in the debater pool for the experiments, including their
423 names, providers, and potentially version information. [Content to be added]

424 **B Debate Pairings Schedule**

425 This appendix details the schedule and method used for pairing LLMs against each other across
426 different debate topics, ensuring a balanced experimental design. [Content to be added]

427 **C Debater Prompt Structures**

428 Full verbatim text of the structured prompts used to guide debater models in the Opening, Rebuttal,
429 and Final rounds, including constraints and judging guidance. [Content to be added]

430 **D AI Jury Prompt Details**

431 Full verbatim text of the detailed prompt provided to the AI jury models for evaluating debate
432 transcripts, including judging criteria and output requirements. [Content to be added]

433 **E Topics of Debate**

434 **F Technical Appendices and Supplementary Material**

435 — YOUR APPENDIX CONTENT HERE (OPTIONAL) —

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [TODO]

Justification: [TODO]

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [TODO]

Justification: [TODO]

3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [TODO]

Justification: [TODO]

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [TODO]

Justification: [TODO]

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [TODO]

Justification: [TODO]

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [TODO]

Justification: [TODO]

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [TODO]

Justification: [TODO]

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [TODO]

Justification: [TODO]

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

483 Answer: **[TODO]**
 484 Justification: **[TODO]**
 485 **10. Broader impacts**
 486 Question: Does the paper discuss both potential positive societal impacts and negative
 487 societal impacts of the work performed?
 488 Answer: **[TODO]**
 489 Justification: **[TODO]**
 490 **11. Safeguards**
 491 Question: Does the paper describe safeguards that have been put in place for responsible
 492 release of data or models that have a high risk for misuse (e.g., pretrained language models,
 493 image generators, or scraped datasets)?
 494 Answer: **[TODO]**
 495 Justification: **[TODO]**
 496 **12. Licenses for existing assets**
 497 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
 498 the paper, properly credited and are the license and terms of use explicitly mentioned and
 499 properly respected?
 500 Answer: **[TODO]**
 501 Justification: **[TODO]**
 502 **13. New assets**
 503 Question: Are new assets introduced in the paper well documented and is the documentation
 504 provided alongside the assets?
 505 Answer: **[TODO]**
 506 Justification: **[TODO]**
 507 **14. Crowdsourcing and research with human subjects**
 508 Question: For crowdsourcing experiments and research with human subjects, does the paper
 509 include the full text of instructions given to participants and screenshots, if applicable, as
 510 well as details about compensation (if any)?
 511 Answer: **[TODO]**
 512 Justification: **[TODO]**
 513 **15. Institutional review board (IRB) approvals or equivalent for research with human**
 514 **subjects**
 515 Question: Does the paper describe potential risks incurred by study participants, whether
 516 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
 517 approvals (or an equivalent approval/review based on the requirements of your country or
 518 institution) were obtained?
 519 Answer: **[TODO]**
 520 Justification: **[TODO]**
 521 **16. Declaration of LLM usage**
 522 Question: Does the paper describe the usage of LLMs if it is an important, original, or
 523 non-standard component of the core methods in this research? Note that if the LLM is used
 524 only for writing, editing, or formatting purposes and does not impact the core methodology,
 525 scientific rigor, or originality of the research, declaration is not required.
 526 Answer: **[TODO]**
 527 Justification: **[TODO]**