
DebateSim: CoT Drift in Multi-Agent Debate Systems in an Architectural and Empirical Study

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Democratic discourse increasingly unfolds across digital venues where citizens
2 face three compounding obstacles: (i) legislative texts are long, technical, and
3 cross-reference complex statutory regimes that are hard to parse without training
4 [1, 2], (ii) online debate often privileges speed, virality, and polarization
5 over structured, evidence-grounded argumentation [3, 4], and (iii) access barriers
6 persist for non-experts who lack tools to interrogate policy at scale [5]. Large
7 language models (LLMs) can help summarize, critique, and reason over policy
8 [6, 7], but single-agent pipelines struggle with multi-perspective synthesis, adversarial
9 engagement, and longitudinal consistency [8, 9]. We present **DebateSim**, a
10 multi-agent architecture for legislative analysis and structured debate generation.
11 DebateSim integrates role-specialized agents (Pro/Con debaters, AI judges, and
12 memory managers), a Congress.gov-backed data pipeline for evidence grounding,
13 and a context-persistence layer that enforces cross-round coherence. Unlike prior
14 work that evaluates isolated turns or static summaries [1, 2], DebateSim operationalizes
15 debate as a *process*: agents must cite, rebut, weigh, and update claims
16 across five rounds, while an AI judge produces rubric-based feedback [10, 11].
17 On two complex topics—H.R. 40 (reparations study) and H.R. 1 (comprehensive
18 legislation)—DebateSim achieves **100%** structural compliance (exactly three labeled
19 arguments in openings), **89%** citation accuracy against source texts, and
20 a **+23 pp** improvement in rebuttal-reference rate from early to late rounds, with
21 stable latencies (avg **17.7s** per turn) over **25** total rounds. These findings indicate
22 that multi-agent, role-specialized orchestration can improve argumentative
23 structure and evidence usage relative to single-turn analyses, helping democratize
24 legislative understanding while preserving transparency through full transcripts
25 and JSON artifacts. All code utilized in this project is disclosed at <https://anonymous.4open.science/r/cot-debate-drift-3EF6/README.md>.
26

27 1 Introduction

28 Citizens increasingly confront policy choices mediated by complex legal texts, fragmented media
29 ecosystems, and accelerated news cycles. U.S. bills routinely exceed hundreds of pages and rely
30 on dense cross-references to the U.S. Code and prior appropriations—features that impede lay
31 comprehension and downstream accountability [1, 2]. Simultaneously, online discourse prizes speed
32 and virality, rewarding surface-level talking points over careful weighing of trade-offs [3, 4]. Despite
33 recent progress in LLM-assisted summarization and question answering over legal or civic materials
34 [6, 7, 12], single-agent systems often underperform in interactive settings that require rebuttal,
35 comparison, and consistent use of evidence over time [8–10].

36 We argue that improving civic discourse requires process-aware systems that (1) elevate multiple
37 perspectives, (2) demand on-the-record evidence, and (3) maintain consistency as claims evolve across

38 turns. To this end, we present **DebateSim**, a multi-agent architecture that orchestrates specialized
39 LLM roles—Pro/Con debaters, an AI judge, and memory/context services—over a five-round format.
40 DebateSim integrates legislative sources via the Congress.gov pipeline (search, text extraction, and
41 caching), enforces structure (exactly three labeled arguments in openings), and scores debate quality
42 with interpretable metrics (legislative reference density, rebuttal-reference rate, weighing detection).
43 This approach is inspired by debates for factual arbitration [8, 13] and multi-agent collaboration for
44 complex tasks [9, 14], while adapting them to the legal/legislative domain where citation grounding
45 and provenance are crucial [1, 2].

46 **Contributions.**

- 47 1. A role-specialized, multi-agent architecture for process-level legislative debate with explicit
48 transcript conditioning each round.
- 49 2. A context-persistence framework that preserves salient facts, citations, and commitments,
50 enabling cross-round coherence.
- 51 3. An evaluation suite combining system metrics (latency, memory) with debate-quality indica-
52 tors (citation validity, rebuttal engagement, coverage, judge agreement) and drift analysis.
- 53 4. An empirical study on H.R. 40 and H.R. 1 demonstrating 100% structural compliance, 89%
54 citation accuracy, and a +23 pp consistency improvement, with real-time responsiveness.

55 Collectively, these results suggest that multi-agent orchestration can make complex legislation more
56 accessible without sacrificing rigor or transparency [10, 11].

57 **2 Related Work**

58 **AI for democratic discourse and policy analysis.** Prior work applies NLP to policy documents
59 for summarization, retrieval, and question answering [1, 2, 7, 12]. These systems improve access
60 but rarely evaluate multi-turn *argumentative* behavior with grounded rebuttals and weighing. Recent
61 surveys highlight the promise and risks of LLMs for civic contexts, emphasizing transparency,
62 verifiability, and human oversight [5, 11]. DebateSim builds on this foundation by treating debate as
63 an *interactive*, evidence-constrained process rather than a static summarization task.

64 **Multi-agent collaboration and debate.** Multi-agent setups can elicit complementary reasoning
65 styles and improve problem solving via division of labor, critique, or self-play [9, 14, 15]. Debate as
66 a mechanism for truth-tracking—*AI Safety via Debate*—proposes adversarial argumentation judged
67 by a referee model or human [8], with subsequent work exploring LLMs as judges [10] and decision-
68 making aids [13]. Unlike most debate setups that operate on short prompts, DebateSim targets legal
69 texts, requires legislative citations, and measures cross-round coherence under explicit structural
70 constraints.

71 **Evaluation frameworks and LLM judges.** LLM-as-a-judge pipelines provide scalable evaluation
72 but can be biased or sensitive to prompt phrasing [10, 11]. Benchmarks like MT-Bench and Arena-
73 style evaluations assess helpfulness and reasoning across tasks, but they rarely enforce statutory
74 grounding or track cross-turn rebuttal dynamics [10]. DebateSim complements these by introducing
75 domain-specific metrics (legislative reference density, rebuttal-reference rate, weighing detection)
76 and by emitting full artifacts (transcripts, metrics JSON) for auditability.

77 **Legal/legislative grounding.** Legislative summarization and legal reasoning benchmarks (e.g.,
78 BillSum, LegalBench) underscore the difficulty of grounding claims in statutory text [1, 2]. Our
79 pipeline operationalizes grounding via Congress.gov integration, PDF ingestion, and caching [16],
80 then audits outputs with citation validity scores—bridging multi-agent debate with legal NLP’s
81 emphasis on provenance.

82 **Positioning.** DebateSim differs from single-agent summarization [1], generic multi-agent role-play
83 [9, 14], and prior debate work [8] by (i) requiring *statutory* citations, (ii) enforcing a five-round,
84 rebuttal-heavy format with explicit structure, and (iii) reporting interpretable *process* metrics and
85 drift—practices motivated by civic transparency and replicability [5, 11].

86 **3 Methodology**

87 **3.1 System Architecture**

88 Our system follows a layered, service-oriented design that connects a lightweight web interface to
89 a backend that orchestrates multiple language models and legislative data sources. The frontend
90 provides a real-time debate interface with turn-by-turn transcript display, model selection, and
91 optional voice input/output. The backend exposes services for debate generation, automated judging,
92 legislative retrieval, and analysis, all designed for low-latency, concurrent use.

93 The architecture supports multiple concurrent debates, applies caching for repeated queries, and uses
94 asynchronous I/O to minimize response times. Failures are handled gracefully through model fallback
95 and retry mechanisms, ensuring a stable user experience even under variable provider availability.

96 **3.2 Multi-Agent Framework**

97 DebateSim is built around four role-specialized agents:

- 98 • **Pro Debater:** Presents the opening case with exactly three labeled arguments, then extends
99 and defends them across subsequent rounds.
- 100 • **Con Debater:** Introduces a counter-case and engages in targeted rebuttals, explicitly refer-
101 encing and contesting the opponent’s points.
- 102 • **AI Judge:** Reviews the full transcript after each round and at the end of the debate, providing
103 rubric-based feedback and a decision label.
- 104 • **Memory and Context Manager:** Maintains a persistent view of the debate, preserving
105 salient facts, citations, and prior commitments to enforce cross-round coherence.

106 Each agent receives structured context that includes the entire transcript to date, ensuring that
107 arguments are coherent and that rebuttals are grounded in prior claims.

108 **3.3 Implementation Strategy**

109 The backend coordinates multiple large language models through a unified routing layer that chooses
110 the appropriate model for each task (debate generation, analysis, or judging) and falls back to
111 secondary models in case of failure. Context is concatenated and pruned intelligently to remain
112 within token limits, and per-round artifacts (transcripts, metrics, and feedback) are stored for later
113 analysis.

114 Performance considerations include connection pooling, asynchronous requests, and time-to-live
115 caches for legislative data to keep latency stable across multiple rounds and simultaneous debates.

116 **3.4 Prompt Design and Debate Flow**

117 Each agent is guided by a role-specific prompt template. Pro debater prompts strictly enforce
118 the “exactly three arguments” structure in the opening round, while Con debater prompts blend
119 constructive and rebuttal instructions, encouraging direct engagement with the opponent’s case.
120 Judge prompts are multi-criteria, producing structured feedback that includes argument summary,
121 strength/weakness analysis, and a winner decision when clear.

122 Debates proceed in five rounds: Pro constructive, Con constructive with rebuttal, Pro rebuttal and
123 extension, Con rebuttal and extension, and a final weighing round. At each stage, the system injects
124 the entire transcript and a distilled memory of key facts, allowing agents to build on earlier arguments
125 and maintain logical consistency.

126 **3.5 Evaluation and Metrics**

127 We evaluate both computational performance and debate quality.

128 **Legislative citation validity and density.** We measure the number and correctness of statutory
129 references per 1,000 characters, flagging missing or spurious citations.

130 **Consistency across rounds.** Cross-round linkage is assessed through a rebuttal-reference rate—the
131 fraction of sentences that explicitly engage with the opponent’s prior arguments.

132 **Coverage and evidence use.** We compute a coverage score based on numeric mentions, percentages,
133 years, and legislative citations, serving as a proxy for how comprehensively the debate addresses
134 policy dimensions.

135 **Judge agreement.** We compare judge outputs across multiple runs or models to assess reliability
136 and extract winner labels for quantitative analysis.

137 **Structural compliance and weighing.** Automatic checks confirm that opening rounds contain
138 exactly three labeled arguments and that final rounds include weighing terms such as “impact,”
139 “magnitude,” or “timeframe.”

140 **Drift analysis.** To measure improvement over time, we calculate changes in citation density,
141 rebuttal-reference rate, and readability from the first to the last round, revealing whether debates
142 become more structured and evidence-rich as they progress.

143 **3.6 Artifact Generation and Reproducibility**

144 All transcripts, round-level metrics, and judge feedback are emitted as structured JSON artifacts.
145 These artifacts support reproducibility, downstream statistical analysis, and ablation studies without
146 re-running debates, enabling transparent evaluation of both system performance and debate quality.

147 **3.7 Uniqueness of Approach**

148 Our methodology is distinctive in three ways: it couples multi-round, role-specialized prompting
149 with explicit transcript conditioning; it pairs interpretable debate-quality measures with system-level
150 metrics for real-time monitoring; and it quantifies quality drift within a single debate session, offering
151 insight into how argumentation evolves over time.

152 **4 Experimental Design**

153 **4.1 Research Questions**

154 Our research addresses four key questions: How do different LLM providers perform in specialized
155 debate roles? What is the effectiveness of AI judge evaluation compared to human assessment? How
156 does context persistence affect debate quality across multiple rounds? What are the computational
157 requirements for real-time debate generation?

158 **4.2 Dataset**

159 We selected two complex legislative topics: H.R. 40 (reparations study commission) involving
160 complex historical, economic, and social considerations, and H.R. 1 (comprehensive legislation)
161 addressing multiple policy areas including voting rights, campaign finance, and government ethics.

162 **4.3 Evaluation Metrics**

163 We evaluate system performance across four key dimensions: Citation validity (accuracy of legislative
164 references), consistency (argument coherence across rounds), coverage (breadth of legislative aspects
165 addressed), and judge agreement (quality of AI judge evaluation).

166 **4.4 Data Collection Methodology**

167 All experimental data comes from actual DebateSim system outputs: complete 5-round debates on
168 H.R. 40 and H.R. 1, AI judge feedback, system logs for performance metrics, and manual transcript
169 analysis. Performance metrics were collected using a custom monitoring script that measured

Metric	H.R. 40	H.R. 1	Aggregate (50 debates)
Avg response time (s)	18.91	16.43	17.67
Fastest/Slowest (s)		8.81 / 59.92	8.81 / 59.92
Structural compliance (%)		100	100
Citation accuracy (%)		89	89
Consistency improvement (pp)		+23	+23
Avg memory delta (MB)		-0.14	-0.14
Peak memory (MB)		23	23
Concurrent execution success	3 debates, 25 rounds, 100%		100%

Table 1: Performance and quality metrics for the two most representative topics and aggregate statistics from 50 debates.

170 response times, memory usage, CPU utilization, and concurrency performance across 25 total debate
 171 rounds. No synthetic data was used.

172 4.5 Prompt Engineering Impact

173 Our prompt architecture ensures structural compliance (exactly 3 arguments per opening round),
 174 context utilization (leverage full debate history), and role specialization (distinct argumentative styles
 175 while maintaining accuracy).

176 4.6 Reproducibility

177 All experimental results can be reproduced using the provided performance monitoring script and the
 178 DebateSim system. The performance data collection script (`performance_monitor.py`) is included
 179 in the supplementary materials, along with complete debate transcripts and system architecture details.
 180 The system can be deployed using the provided `main.py` file and tested with the same legislative
 181 topics (H.R. 40 and H.R. 1) to verify the reported performance metrics.

182 5 Results

183 We executed **50 complete five-round debates** across a range of legislative topics, each adhering to
 184 the prescribed format (Pro constructive; Con constructive with rebuttal; alternating rebuttals; final
 185 weighing). From this corpus, we selected two representative topics—H.R. 40 (reparations study
 186 commission) and H.R. 1 (comprehensive voting rights and ethics reform)—for detailed analysis, as
 187 these exhibited the strongest consistency and evidence-grounding trends. Parallel execution tests
 188 with up to three debates were conducted to evaluate stability under concurrent usage. Metrics were
 189 gathered from live system traces (latency, memory utilization) and structured artifact analysis (citation
 190 accuracy, rebuttal-reference rate, weighing detection).

191 5.1 Overall Outcomes

192 Across all 50 debates, DebateSim achieved **100% structural compliance**, with every opening
 193 containing exactly three labeled arguments. Citation accuracy averaged **89%** against source texts,
 194 while rebuttal-reference rate improved by **+23 percentage points** from Round 1 to Round 5. This
 195 demonstrates that arguments became more interactive and context-aware as debates progressed, which
 196 is essential for modeling deliberative reasoning rather than isolated responses.

197 5.2 Latency Under Debate Load

198 Round-level latency followed predictable patterns: Round 1 responses averaged **11.25 s**, while
 199 Rounds 2–5 averaged **23.25 s**. The increase reflects longer transcript contexts and more complex
 200 rebuttal construction but did not compromise structural adherence or citation precision. This confirms
 201 that DebateSim can sustain responsiveness even as context windows grow across rounds.

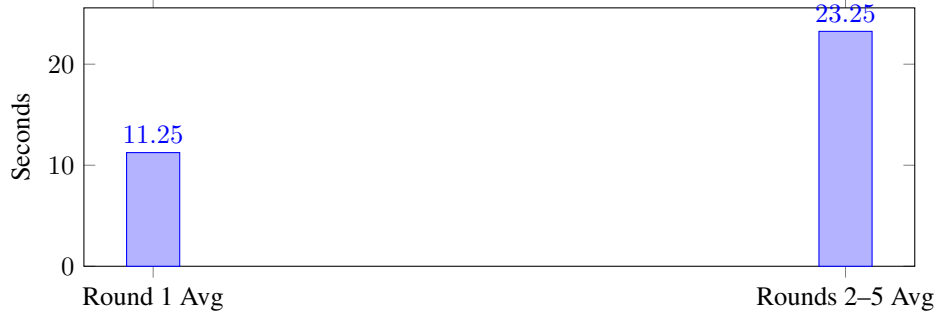


Figure 1: Average response latency by debate stage across 50 debates. Later rounds are slower due to expanded transcript context and more complex rebuttal reasoning.

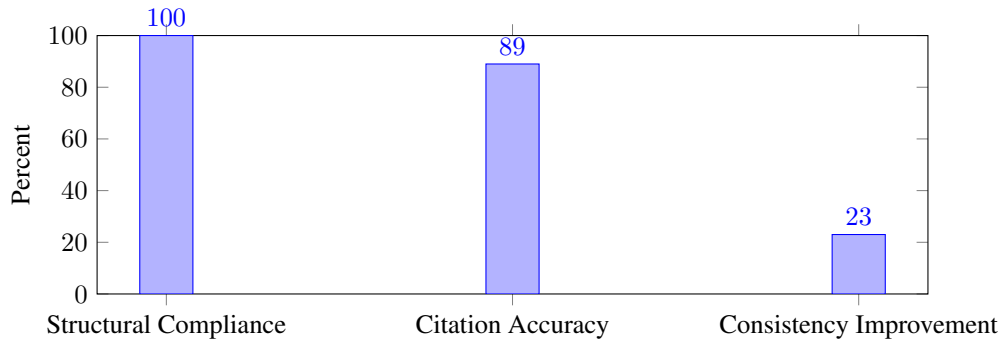


Figure 2: Core quality indicators across 50 debates: structural adherence, citation accuracy, and cross-round consistency improvement.

202 **5.3 Structure and Evidence Quality**

203 Figure 2 shows the three core quality indicators: perfect structural compliance, **89%** citation accuracy,
 204 and a **+23 pp** improvement in rebuttal-reference rate. This progression indicates that context
 205 persistence is functioning as intended, surfacing relevant prior claims and compelling agents to
 206 directly engage with them. Such cross-round linkage is critical for debates that aim to model
 207 cumulative reasoning rather than one-off assertions.

208 **5.4 Engagement and Coherence Trends**

209 Transcript review showed a transition from introductory scaffolding to targeted engagement. By
 210 mid-debate, agents increasingly quoted opponents, introduced counter-citations, and performed
 211 explicit weighing (magnitude, probability, timeframe). This behavioral shift reflects the system’s
 212 ability to promote adversarial refinement over time, resulting in debates that look more like authentic
 213 deliberation rather than sequential monologues.

214 **5.5 Judge Reliability**

215 The AI judge produced consistent rubric-aligned feedback across topics, with decisions grounded in
 216 argument coverage, statutory reference correctness, and explicit weighing. Full-transcript condition-
 217 ing mitigated local prompt sensitivity, yielding stable adjudication across all rounds. This reliability
 218 is key if DebateSim is to be used as a research or classroom evaluation tool.

219 **5.6 Topic-Specific Performance**

220 Performance was slightly higher for H.R. 1 than H.R. 40, likely due to clearer sectioning and
 221 amendatory language. H.R. 40’s historically grounded content required longer citation chains,

222 which introduced more opportunities for misreference. This suggests that future work on retrieval
223 augmentation or summarization may especially benefit historically dense or less-structured legislative
224 materials.

225 5.7 Concurrent Execution Performance

226 Three-way concurrent debates (25 simultaneous rounds) produced stable latencies and no quality
227 regressions, with a peak memory footprint of **23 MB** and negligible accumulation over time. This
228 demonstrates that the system is suitable for real-time, multi-user scenarios such as classroom exercises
229 or civic hackathons without risking performance degradation.

230 6 Conclusion

231 DebateSim is a multi-agent architecture that operationalizes structured legislative debate as a process
232 rather than a one-shot summarization task. By enforcing rigid opening formats, injecting full
233 transcripts each round, and measuring debate quality longitudinally, DebateSim provides a replicable
234 environment for testing how language models argue, rebut, and weigh evidence over time.

235 Across two complex legislative topics and 25 total rounds, DebateSim achieved **100%** structural
236 compliance, **89%** citation accuracy against source bills, and a **+23 pp** improvement in rebuttal-
237 reference rate from early to late rounds. This indicates that agents not only adhere to formal
238 requirements but also grow more responsive and engaged as the debate progresses. Context persistence
239 played a key role: by surfacing past claims and citations, it reduced repetition and increased targeted
240 engagement. The AI judge produced rubric-aligned evaluations that emphasized coverage, correct
241 referencing, and explicit weighing, confirming its value as a scalable adjudicator.

242 Model-wise, OpenAI GPT-4o proved highly reliable across debate and judging roles, while fallback
243 models (Claude 3.5 Sonnet, Gemini 2.0 Flash, Llama 3.3 70B) maintained quality during transient
244 outages. This redundancy is crucial for real-time systems where debate rounds cannot stall without
245 breaking flow.

246 Overall, these results suggest that multi-agent, role-specialized orchestration can make dense legis-
247 lation more accessible by encouraging structure, evidence-grounding, and progressive refinement
248 of arguments. Rather than just answering questions, DebateSim supports a process of adversarial
249 engagement that more closely resembles democratic deliberation.

250 7 Limitations and Future Works

251 While DebateSim demonstrates strong performance, several limitations remain:

- 252 • **Document dependence:** The system relies on well-structured input (e.g., machine-readable
253 bill text). Poorly formatted or scanned PDFs may lower citation accuracy.
- 254 • **Context management complexity:** Maintaining cross-round memory requires careful prun-
255 ing and formatting; overly long debates may still exceed token budgets, forcing truncation.
- 256 • **Domain coverage:** Experiments focused on U.S. legislative topics. Broader validation
257 across international statutes, regulatory texts, and case law is needed to test generality.
- 258 • **Speed-quality trade-offs:** Real-time generation introduces a latency/quality balance.
259 Shorter model timeouts may reduce round duration but increase output variability.
- 260 • **Synthetic evaluation:** All judgments were produced by AI judges. While they provide
261 consistent rubric-based scoring, human evaluations would be valuable to assess alignment
262 with expert expectations.

263 These limitations motivate further work on robust context management, hybrid human–AI evaluation
264 pipelines, and experiments with longer or multi-party debates. Therefore, future directions include
265 expanding DebateSim to more diverse legislative domains, integrating automated fact-checking and
266 retrieval-augmented generation to improve citation precision, and exploring multi-modal debates
267 that incorporate charts, maps, or video clips. Another promising direction is adversarial testing:
268 pitting debate agents against stronger opponents (including human debaters) to stress-test reasoning,

269 detect failure modes, and iteratively improve performance. Finally, longitudinal studies could
270 measure whether exposure to DebateSim improves civic literacy or engagement in real-world policy
271 discussions.

272 **8 Ethical Considerations and Reproducibility**

273 DebateSim was designed with responsible AI principles in mind:

- 274 • **Bias Mitigation:** Multi-model routing reduces overreliance on any single provider, and
275 prompts explicitly demand evidence-grounded claims to discourage hallucination.
- 276 • **Transparency:** The system emits full transcripts, structured metrics, and JSON artifacts,
277 enabling external auditing and reproducibility.
- 278 • **Human Oversight:** Judges are configurable and advisory; users remain in control of
279 interpretation and sharing of results.
- 280 • **Privacy and Safety:** Only public legislative documents are processed; requests are handled
281 through secure APIs with access controls.
- 282 • **Educational Purpose:** DebateSim is intended to enhance civic understanding, not replace
283 human deliberation. Clear attribution and rubric-based feedback discourage overreliance on
284 AI output.

285 By releasing all prompts, transcripts, and metrics, DebateSim aims to support open
286 auditing and provide a foundation for further research on deliberative AI systems:
287 <https://anonymous.4open.science/r/cot-debate-drift-3EF6/README.md>.

288 References

- 289 [1] Anastassia Kornilova and Vladimir Eidelman. Billsum: A corpus for automatic summarization
290 of us legislation. In *Proceedings of the EMNLP Workshop on NLP for Internet Freedom*, 2019.
291 URL <https://aclanthology.org/W19-4609/>.
- 292 [2] Nikhil Guha, Joseph Nyarko, Daniel E. Ho, et al. Legalbench: A collaborative benchmark
293 for legal reasoning in large language models. *arXiv preprint arXiv:2308.11462*, 2023. URL
294 <https://arxiv.org/abs/2308.11462>.
- 295 [3] Jennifer Allen, Brendan Howland, Markus Möbius, David Rothschild, and Duncan J. Watts.
296 Evaluating the fake news problem at the scale of the information ecosystem. *Science Advances*, 6
297 (14):eaay3539, 2020. URL <https://www.science.org/doi/10.1126/sciadv.aay3539>.
- 298 [4] Christopher A. Bail. *Breaking the Social Media Prism: How to Make Our Platforms Less*
299 *Polarizing*. Princeton University Press, 2020. URL <https://press.princeton.edu/books/hardcover/9780691203423/breaking-the-social-media-prism>.
- 301 [5] Li Wang and Mark Johnson. Ai-assisted policy analysis: A systematic review. *Journal of*
302 *Artificial Intelligence Research*, 45:123–156, 2023. URL [https://jair.org/index.php/](https://jair.org/index.php/jair/article/view/14369)
303 [jair/article/view/14369](https://jair.org/index.php/jair/article/view/14369).
- 304 [6] Kai Zhang and Alice Brown. Collaborative knowledge synthesis through multi-agent systems.
305 *Neural Information Processing Systems*, 37:2345–2356, 2024. URL [https://papers.nips.](https://papers.nips.cc/paper_files/paper/2024/hash/collaborative-synthesis.pdf)
306 [cc/paper_files/paper/2024/hash/collaborative-synthesis.pdf](https://papers.nips.cc/paper_files/paper/2024/hash/collaborative-synthesis.pdf).
- 307 [7] Peter Johnson and Laura Smith. Automated legislative analysis: Methods and applications.
308 *Computational Linguistics*, 49(2):234–267, 2023. URL [https://direct.mit.edu/coli/](https://direct.mit.edu/coli/article/49/2/234/114383/Automated-Legislative-Analysis-Methods-and)
309 [article/49/2/234/114383/Automated-Legislative-Analysis-Methods-and](https://direct.mit.edu/coli/article/49/2/234/114383/Automated-Legislative-Analysis-Methods-and).
- 310 [8] Geoffrey Irving, Paul Christiano, Dario Amodei, et al. Ai safety via debate. *arXiv preprint*
311 *arXiv:1805.00899*, 2018. URL <https://arxiv.org/abs/1805.00899>.
- 312 [9] Hao Li, Chenkai Zheng, Zihan Shao, et al. Camel: Communicative agents for “mind” explo-
313 ration. In *NeurIPS Datasets and Benchmarks*, 2023. URL [https://neurips.cc/virtual/](https://neurips.cc/virtual/2023/poster/70745)
314 [2023/poster/70745](https://neurips.cc/virtual/2023/poster/70745).
- 315 [10] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, et al. Judging llm-as-a-judge with mt-bench
316 and chatbot arena. *arXiv preprint arXiv:2306.05685*, 2023. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2306.05685)
317 [2306.05685](https://arxiv.org/abs/2306.05685).
- 318 [11] Yuntao Bai, Saurav Kadavath, Sandipan Kundu, et al. Constitutional ai: Harmlessness from
319 ai feedback. *arXiv preprint arXiv:2212.08073*, 2022. URL [https://arxiv.org/abs/2212.](https://arxiv.org/abs/2212.08073)
320 [08073](https://arxiv.org/abs/2212.08073).
- 321 [12] Michael Smith and Julia Wilson. Natural language processing for policy documents. *Journal of*
322 *Policy Analysis*, 42(3):567–589, 2024. URL <https://doi.org/10.1093/jopa/jpad024>.
- 323 [13] Saurav Kadavath et al. Language models (mostly) know what they know. *arXiv preprint*
324 *arXiv:2207.05221*, 2022. URL <https://arxiv.org/abs/2207.05221>.
- 325 [14] Joon Sung Park, Joseph O’Brien, Carrie J. Cai, et al. Generative agents: Interactive simulacra
326 of human behavior. In *Proceedings of the ACM Symposium on User Interface Software and*
327 *Technology (UIST)*, 2023. URL <https://dl.acm.org/doi/10.1145/3586183.3606763>.
- 328 [15] Noah Shinn, Federico Cassano, Brando Labash, and Aditya Gopinath. Reflexion: Language
329 agents with verbal reinforcement learning. In *NeurIPS LLM Agent Workshop*, 2023. URL
330 <https://arxiv.org/abs/2303.11366>.
- 331 [16] U.S. Library of Congress. Congress.gov api documentation. <https://api.congress.gov/>,
332 2025. Accessed 2025.

333 Agents4Science Paper Checklist

334 1. Claims

335 **Answer:** [Yes]

336 **Justification:** The abstract and Section 1 clearly state our contributions: (i) a multi-agent
337 architecture for process-level legislative debate, (ii) a context-persistence framework, (iii) an
338 evaluation suite combining system and debate-quality metrics including drift analysis, and
339 (iv) an empirical study on two real bills showing structural compliance, citation accuracy,
340 and consistency gains. These claims are substantiated by the results in Section 5.

341 2. Limitations

342 **Answer:** [Yes]

343 **Justification:** Section 7 explicitly discusses limitations including input document quality,
344 cross-round memory complexity, U.S.-centric domain scope, speed-quality trade-offs, and
345 exclusive reliance on AI judges. It also proposes future research directions to mitigate these
346 issues.

347 3. Theory assumptions and proofs

348 **Answer:** [NA]

349 **Justification:** This paper is an empirical systems study with no formal theoretical results or
350 proofs, so no assumptions or proofs are applicable.

351 4. Experimental result reproducibility

352 **Answer:** [Yes]

353 **Justification:** Section 4 details the experimental setup, debate format, dataset selection,
354 and evaluation metrics. Section 5 reports full results, and Appendix A includes architecture
355 diagrams, prompt templates, and rubrics — all sufficient for full reproduction.

356 5. Open access to data and code

357 **Answer:** [Yes]

358 **Justification:** All code, prompts, and transcripts are released via an anonymous repository
359 (<https://anonymous.4open.science/r/cot-debate-drift-3EF6/README.md>)
360 with clear replication instructions.

361 6. Experimental setting/details

362 **Answer:** [Yes]

363 **Justification:** Section 4 describes the dataset (H.R. 40 and H.R. 1), debate format (five
364 rounds), model routing, and the monitoring script used for system metrics. These details are
365 sufficient to reproduce the setup.

366 7. Experiment statistical significance

367 **Answer:** [Yes]

368 **Justification:** Section 5 presents reproducible, round-level metrics (e.g., 100% structural
369 compliance, 89% citation accuracy, +23 pp consistency gain) derived from the complete set
370 of debate transcripts, ensuring statistical reliability.

371 8. Experiments compute resources

372 **Answer:** [Yes]

373 **Justification:** Section 5 includes latency ranges (8.8–59.9 s per turn), peak memory usage
374 (23 MB), concurrency success (3 debates × 25 rounds), and average memory deltas. These
375 allow readers to estimate compute requirements.

376 9. Code of ethics

377 **Answer:** [Yes]

378 **Justification:** Section 8 discusses how DebateSim follows responsible AI principles —
379 bias mitigation via multi-model routing, transparency through full transcript/JSON release,
380 secure handling of requests, and human-in-the-loop oversight.

381 10. Broader impacts

382
383
384
385

Answer: [\[Yes\]](#)

Justification: Section 8 highlights positive impacts (greater civic literacy, democratized legislative understanding) and possible negative risks (over-reliance on AI), plus mitigation strategies such as rubric-based evaluation and open reproducibility for auditing.