

---

# When Researchers Say Mental Model/Theory of Mind of AI, What Are They Really Talking About?

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1           When researchers claim AI systems possess ToM or mental models, they are funda-  
2           mentally discussing behavioral predictions and bias corrections rather than genuine  
3           mental states. This position paper argues that the current discourse conflates so-  
4           phisticated pattern matching with authentic cognition, missing a crucial distinction  
5           between simulation and experience. While recent studies show LLMs achieving  
6           human-level performance on ToM laboratory tasks, these results are based only on  
7           behavioral mimicry. More importantly, the entire testing paradigm may be flawed  
8           in applying individual human cognitive tests to AI systems, but assessing human  
9           cognition directly in the moment of human-AI interaction. I suggest shifting focus  
10          toward mutual ToM frameworks that acknowledge the simultaneous contributions  
11          of human cognition and AI algorithms, emphasizing the interaction dynamics,  
12          instead of testing AI in isolation.

## 13   1   Introduction

14       Humans develop theories to explain each other's behaviors (Sellars, 1956). This ability to infer  
15       each other's mental states (Premack & Woodruff, 1978) has been called the theory of mind (ToM).  
16       However, we are in direct continuous contact with our own minds but not the mental states of others;  
17       therefore, studying human mental states is difficult. This abstract concept is not related to any specific  
18       parameters or rubric; it's the theory of an experienced mental state that we can't observe in any  
19       specific event or attribute to different, distinct events.

20       Kosinski (2023) argues that individual artificial neurons in LLMs function like "Chinese rooms",  
21       which follow mathematical instructions without genuine understanding. The author also suggests  
22       that complex cognitive abilities may emerge at the network level, similar to how human cognition  
23       is assumed to emerge from networks of neurons. Mechanically speaking, this comparison makes  
24       sense. However, debating whether models are merely Chinese rooms may be unproductive. The key  
25       is to study how they interact, as humans typically don't question others' cognition during normal  
26       interaction. Meanwhile, Strachan et al. (2024) claim that LLM behavior is "indistinguishable from  
27       human behavior" on ToM tasks on various dimensions. However, if you listen to an authentic  
28       interaction between two humans and an interaction between a human and an AI agent, you can tell a  
29       difference. Under the assumption that a ToM is fundamental to human social interaction, rather than  
30       a byproduct of it, this finding has motivated researchers to develop benchmarks that offer a practical  
31       and cost-effective approach to assessing ToM capabilities in LLMs.

32       Researchers like Gu et al. (2024) tested GPT-4 using their SimpleToM dataset. While GPT-4  
33       achieves approximately 90% accuracy on ToM questions, performance drops to approximately 50%  
34       on behavior prediction and 15% on behavioral judgment. The test results are not ideal, and the  
35       test itself is also questionable, unless being able to answer a battery of ToM questions proves that  
36       GPT-4 has a ToM. The limitations of this kind of approaches has also been pointed out by Wang

37 et al. (2025), that current ToM tasks have three major limitations: 1) Theoretically, ToM should be  
38 multidimensional but has only been tested in one dimension; 2) Current ToM tests lack construct  
39 validity; and 3) Evaluations always use third-person static scenarios rather than spontaneous dynamic  
40 interactions. They also report that 75.5% of ToM measures focus only on beliefs rather than other  
41 considerations such as emotion, desires, and intentions. In a narrow technical sense, AI systems that  
42 track beliefs and predict behaviors could be said to have "ToM". However, this technical capability  
43 differs fundamentally from human ToM grounded in embodied experience.

## 44 **2 The Fundamental Flaw in Cognitive Testing for LLMs**

45 The trend in evaluating LLMs using cognitive tasks may show a fundamental misunderstanding  
46 about what these systems actually are and what we need to know about them. Researchers designed  
47 ToM tests, reasoning challenges, and planning problems for LLMs by assuming that humans process  
48 information like LLMs do. But as Kambhampati (2024) argues, LLMs are fundamentally "n-gram  
49 models on steroids" performing "universal approximate retrieval" rather than reasoning focused on  
50 survival that humans do.

51 Consider how humans may develop ToM through embodied experience. Children learn that others  
52 have different perspectives through physical interactions such as hiding objects, playing peek-a-boo,  
53 and observing emotional responses. Psychologists study how children build up their abilities using  
54 developmental scales. For example, Wellman and Liu (2004) described scales of a sequence of tasks  
55 to test children's diverse desires, diverse beliefs, knowledge access, contents false belief, explicit  
56 false belief, belief-emotion, and real-apparent emotion. While this is representative work, it doesn't  
57 mean we should test LLMs in the same way or even compare them to humans (such as the SimToM  
58 prompting framework by Wilf et al. (2024)) by assuming these researchers are correct in reducing  
59 human interaction to their identified dimensions of a ToM.

60 Some researchers have designed benchmarks to measure ToM abilities specifically for LLMs, such as  
61 Kanishk et al. (2023) with their social reasoning benchmark (BigToM) and Chen et al. (2024) with  
62 their ToMBench. But why should we accept that LLMs even have a belief that can be tested? More  
63 importantly, what does passing these tests tell us about how AI will function in dynamic human-AI  
64 interactions?

65 Many benchmarks developed for testing LLMs' ToM are derived from psychological tests such as the  
66 Sally-Anne test (Scassellati, 2001) but lack validity and reliability. As ToM is proposed as an ability  
67 to understand and respond to ongoing feedback from the social environment, the evaluation of ToM  
68 should not be static (Wang et al., 2025). The tests are measuring something about ToM, certainly,  
69 but not meaningful interaction in a valid sense. They're measuring how well statistical patterns from  
70 human cognitive behavior in the training data can be reproduced in response to prompts. In other  
71 words, it is learning patterns that we intend it to, and we are asked to accept this as having a ToM that  
72 some researchers want it to.

73 This becomes clearer when we examine how LLMs fail. When planning problems are presented with  
74 obfuscated names, removing the statistical patterns from training data, LLM performance plummets  
75 dramatically. They're not reasoning inductively through the structure of problems; they're deductively  
76 matching surface patterns to similar examples in their training corpus. The same limitation appears  
77 in self-critique tasks. Despite claims that LLMs can verify their own reasoning, studies show  
78 performance actually worsens with self-verification as models "hallucinate" both false positives and  
79 false negatives.

80 The recent phenomenon that Cuadron et al. (2025) termed the "Reasoning-Action Dilemma" further  
81 exposes this limitation. Large Reasoning Models exhibit what appears to be "overthinking", where  
82 they bias internal reasoning over environmental feedback. But this isn't genuine overthinking in  
83 the human sense, such as in anxiety-driven rumination. It's the model getting stuck in loops of text  
84 generation that statistically resemble reasoning without genuine deliberation or decision-making  
85 authority.

86 The test is legitimate for measuring behavioral reproduction accuracy, but not valid for inferring  
87 genuine social cognitive processes or, more importantly, for predicting that AI functions as a human  
88 does in human contexts.

### 89 **3 Why Testing Misses the Point**

90 The distinction between simulation and authentic mental processes matters, but not for reasons  
91 typically discussed. LLMs are good at simulating ToM responses because they've been trained on  
92 billions of examples of human discourse about mental states aimed at generating those responses.  
93 They can produce outputs that seem to demonstrate understanding of beliefs, desires, and intentions.  
94 However, this cannot be called "ToM" in an adaptive human sense.

95 This distinction becomes clearer when examining motivated reasoning. Humans don't always reason  
96 correctly, even when they possess sufficient cognitive resources. We engage in "wishful thinking",  
97 "confirmation bias", and emotionally driven reasoning. As Kunda (1990) notes, human reasoning  
98 involves "hot motivation to reach desired goals" alongside "cold motivation favoring accuracy". AI  
99 systems, following optimization objectives, will lack these intrinsic motivational states. When AI  
100 systems make errors, it's not because desire overrides logic but because they've encountered edge  
101 cases that their training and algorithm do not cover.

102 Proponents of AI cognition often invoke "emergence", which is a mysterious way of claiming that  
103 sufficiently complex neural networks spontaneously develop genuine understanding. They point to  
104 studies like Zhu et al. (2024) showing LLMs internally represent different agents' beliefs in their  
105 neural activations. However, finding structured representations doesn't prove consciousness exists any  
106 more than finding face-selective neurons proves the fusiform face area of the human brain experiences  
107 faces.

108 Debate about genuine versus simulated cognition distracts from the central issue: isolated ToM tests  
109 do not tell us what actually matters, nor how AI systems should function within human social and  
110 collaborative contexts.

### 111 **4 Mutual ToM**

112 Rather than asking whether AI "has" a ToM, we should examine how humans and AI systems can  
113 mutually develop understanding. Instead of validating AI performance on scales of ToM using isolated  
114 cognitive tests, we should test the system dynamics when AI operates within human environments.  
115 The critical question isn't whether AI can pass a false-belief task in isolation, but how human-AI  
116 interaction changes both behavior and outcomes when the human-AI system is working on a project  
117 rather than taking a test.

118 Wang and Goel's (2022) Mutual ToM framework offers one approach that preserves the ToM concept  
119 by focusing on three iteratively shaping elements: interpretation, feedback, and mutuality. In this  
120 framework, humans and AI agents construct representations of each other, but these representations  
121 serve different functions. Humans apply ToM to AI systems through anthropomorphism, attributing  
122 mental states that the AI doesn't possess. Meanwhile, AI systems build statistical models predicting  
123 human behavior without assuming any genuine understanding of mental states.

124 The key insight is that effective human-AI collaboration doesn't require AI to have a ToM but rather  
125 requires systems that support mutual adaptation and understanding. Zhang et al.'s (2024) findings  
126 support this approach. They found AI's independent ToM capability didn't significantly impact  
127 team performance. What mattered was enhancing human understanding of the agent. Surprisingly,  
128 bidirectional communication sometimes decreased performance, with participants reporting increased  
129 cognitive workload. This suggests that pretending AI has ToM might actually hinder rather than help  
130 human-AI collaboration. Considering the dynamic nature of mutual theory of mind, human-autonomy  
131 team (HAT) resilience studies employing dynamic measurements (e.g., Grimm et al., 2023) can serve  
132 as a valuable reference for future ToM studies, particularly those focusing on team-level states in  
133 human-AI collaboration.

134 This distinction matters for practical purposes. We should stop asking "Does this AI understand  
135 false beliefs?" and start asking "How does this AI system change human behavior in collaborative  
136 settings?" Then stop endless philosophical debates about "emergent" consciousness and ToM test  
137 scores.

## 138 5 Conclusion

139 When researchers test AI’s ToM, they’re observing sophisticated behavioral prediction, not genuine  
140 mental state attribution. But more importantly, they’re asking the wrong question. The current ap-  
141 proach of administering human cognitive tests to AI systems, whether adapted from child psychology  
142 or specially designed for LLMs, fundamentally misunderstands both what these systems are from a  
143 psychological perspective and what we need to know to use them effectively.

144 LLMs achieve remarkable performance on ToM tasks through statistical pattern matching across  
145 massive datasets, but this differs fundamentally from human experience, grounded in embodied  
146 experience, motivated reasoning, and genuine understanding, which led to the development of the  
147 ToM concept in the first place. However, even if AI does replicate ToM behaviors in testing conditions,  
148 these isolated assessments do not tell us about what actually matters: how human-AI systems function  
149 together as a cognitive system. Specifically, we should abandon attempts to improve AI’s scores on  
150 ToM tests and instead study mutual adaptation and understanding between humans and AI systems.  
151 The mutual approach acknowledges both human tendencies to anthropomorphize and AI’s statistical  
152 modeling capabilities, designing for effective collaboration within these constraints. The question  
153 isn’t whether AI thinks like or simulates the human mind. It’s how humans and AI can work together  
154 in the service of human survival and flourishing. After all, what matters isn’t how well models score  
155 on ToM tests, but rather how good the models can serve us, humans.

## 156 References

- 157 [1] Chen, Z., Wu, J., Zhou, J., Wen, B., Bi, G., Jiang, G., Cao, Y., Hu, M., Lai, Y., Xiong, Z., & Huang, M.  
158 (2024). ToMBench: Benchmarking theory of mind in large language models. *arXiv preprint arXiv:2402.15052*.
- 159 [2] Cuadron, A., Li, D., Ma, W., Wang, X., Wang, Y., Zhuang, S., Liu, S., Schroeder, L. G., Xia, T., Mao,  
160 H., Thumiger, N., Desai, A., Stoica, I., Klimovic, A., Neubig, G., & Gonzalez, J. E. (2025). The danger of  
161 overthinking: Examining the reasoning-action dilemma in agentic tasks. *arXiv preprint arXiv:2502.08235*.
- 162 [3] Gandhi, K., Fränken, J. P., Gerstenberg, T., & Goodman, N. (2023). Understanding social reasoning in  
163 language models with language models. In *Advances in Neural Information Processing Systems* (Vol. 36, pp.  
164 13518-13529).
- 165 [4] Grimm, D. A., Gorman, J. C., Cooke, N. J., Demir, M., & McNeese, N. J. (2023). Dynamical measurement  
166 of team resilience. *Journal of Cognitive Engineering and Decision Making*, 17(4), 351-382.
- 167 [5] Gu, Y., Tafjord, O., Kim, H., Moore, J., Bras, R. L., Clark, P., & Choi, Y. (2024). SimpleTom: Exposing the  
168 gap between explicit tom inference and implicit tom application in llms. *arXiv preprint arXiv:2410.13648*.
- 169 [6] Kambhampati, S. (2024). Can large language models reason and plan? *Annals of the New York Academy of*  
170 *Sciences*, 1534(1), 15-18.
- 171 [7] Kosinski, M. (2023). Theory of mind may have spontaneously emerged in large language models. *arXiv*  
172 *preprint arXiv:2302.02083*.
- 173 [8] Kunda, Z. (1990). The case for motivated reasoning. *Psychological Bulletin*, 108(3), 480-498.
- 174 [9] Premack, D., & Woodruff, G. (1978). Does the chimpanzee have a theory of mind? *Behavioral and Brain*  
175 *Sciences*, 1(4), 515-526.
- 176 [10] Scassellati, B. M. (2001). *Foundations for a Theory of Mind for a Humanoid Robot* (Doctoral dissertation,  
177 Massachusetts Institute of Technology).
- 178 [11] Sellars, W. (1956). Empiricism and the philosophy of mind. *Minnesota Studies in the Philosophy of Science*,  
179 1(19), 253-329.
- 180 [12] Strachan, J. W. A., Albergo, D., Borghini, G., Pansardi, O., Scaliti, E., Gupta, S., Saxena, K., Rufo, A.,  
181 Panzeri, S., Manzi, G., Graziano, M. S. A., & Becchio, C. (2024). Testing theory of mind in large language  
182 models and humans. *Nature Human Behaviour*, 8(7), 1285-1295.
- 183 [13] Wang, Q., & Goel, A. K. (2022). Mutual theory of mind for human-AI communication. *arXiv preprint*  
184 *arXiv:2210.03842*.
- 185 [14] Wang, Q., Zhou, X., Sap, M., Forlizzi, J., & Shen, H. (2025). Rethinking theory of mind benchmarks for  
186 llms: Towards a user-centered perspective. *arXiv preprint arXiv:2504.10839*.
- 187 [15] Wellman, H. M., & Liu, D. (2004). Scaling of theory-of-mind tasks. *Child Development*, 75(2), 523-541.

- 188 [16] Wilf, A., Lee, S. S., Liang, P. P., & Morency, L. P. (2023). Think twice: Perspective-taking improves large  
189 language models' theory-of-mind capabilities. *arXiv preprint arXiv:2311.10227*.
- 190 [17] Zhang, S., Wang, X., Zhang, W., Chen, Y., Gao, L., Wang, D., Zhang, W., Wang, X., & Wen, Y. (2024).  
191 Mutual theory of mind in human-ai collaboration: An empirical study with llm-driven ai agents in a real-time  
192 shared workspace task. *arXiv preprint arXiv:2409.08811*.
- 193 [18] Zhu, W., Zhang, Z., & Wang, Y. (2024). Language models represent beliefs of self and others. *arXiv*  
194 *preprint arXiv:2402.18496*.

195 **NeurIPS Paper Checklist**

196 The checklist is designed to encourage best practices for responsible machine learning research,  
197 addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove  
198 the checklist: **The papers not including the checklist will be desk rejected.** The checklist should  
199 follow the references and follow the (optional) supplemental material. The checklist does NOT count  
200 towards the page limit.

201 Please read the checklist guidelines carefully for information on how to answer these questions. For  
202 each question in the checklist:

- 203 • You should answer [Yes] , [No] , or [NA] .
- 204 • [NA] means either that the question is Not Applicable for that particular paper or the  
205 relevant information is Not Available.
- 206 • Please provide a short (1–2 sentence) justification right after your answer (even for NA).

207 **The checklist answers are an integral part of your paper submission.** They are visible to the  
208 reviewers, area chairs, senior area chairs, and ethics reviewers. You will be asked to also include it  
209 (after eventual revisions) with the final version of your paper, and its final version will be published  
210 with the paper.

211 The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation.  
212 While "[Yes]" is generally preferable to "[No]", it is perfectly acceptable to answer "[No]" provided a  
213 proper justification is given (e.g., "error bars are not reported because it would be too computationally  
214 expensive" or "we were unable to find the license for the dataset we used"). In general, answering  
215 "[No]" or "[NA]" is not grounds for rejection. While the questions are phrased in a binary way, we  
216 acknowledge that the true answer is often more nuanced, so please just use your best judgment and  
217 write a justification to elaborate. All supporting evidence can appear either in the main paper or the  
218 supplemental material, provided in appendix. If you answer [Yes] to a question, in the justification  
219 please point to the section(s) where related material for the question can be found.

220 IMPORTANT, please:

- 221 • **Delete this instruction block, but keep the section heading “NeurIPS Paper Checklist”,**
- 222 • **Keep the checklist subsection headings, questions/answers and guidelines below.**
- 223 • **Do not modify the questions and only use the provided macros for your answers.**

224 **1. Claims**

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

227 Answer: [Yes]

228 Justification: The main claims made in the abstract and introduction accurately reflect the  
229 paper’s contributions and scope.

230 Guidelines:

- 231 • The answer NA means that the abstract and introduction do not include the claims  
232 made in the paper.
- 233 • The abstract and/or introduction should clearly state the claims made, including the  
234 contributions made in the paper and important assumptions and limitations. A No or  
235 NA answer to this question will not be perceived well by the reviewers.
- 236 • The claims made should match theoretical and experimental results, and reflect how  
237 much the results can be expected to generalize to other settings.
- 238 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
239 are not attained by the paper.

240 **2. Limitations**

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

242 Answer: [NA]

243 Justification: It's a theoretical position paper; the whole paper was purposed to discuss the  
244 limitations of the current approach.

245 Guidelines:

- 246 • The answer NA means that the paper has no limitation while the answer No means that  
247 the paper has limitations, but those are not discussed in the paper.
- 248 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 249 • The paper should point out any strong assumptions and how robust the results are to  
250 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
251 model well-specification, asymptotic approximations only holding locally). The authors  
252 should reflect on how these assumptions might be violated in practice and what the  
253 implications would be.
- 254 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
255 only tested on a few datasets or with a few runs. In general, empirical results often  
256 depend on implicit assumptions, which should be articulated.
- 257 • The authors should reflect on the factors that influence the performance of the approach.  
258 For example, a facial recognition algorithm may perform poorly when image resolution  
259 is low or images are taken in low lighting. Or a speech-to-text system might not be  
260 used reliably to provide closed captions for online lectures because it fails to handle  
261 technical jargon.
- 262 • The authors should discuss the computational efficiency of the proposed algorithms  
263 and how they scale with dataset size.
- 264 • If applicable, the authors should discuss possible limitations of their approach to  
265 address problems of privacy and fairness.
- 266 • While the authors might fear that complete honesty about limitations might be used by  
267 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
268 limitations that aren't acknowledged in the paper. The authors should use their best  
269 judgment and recognize that individual actions in favor of transparency play an impor-  
270 tant role in developing norms that preserve the integrity of the community. Reviewers  
271 will be specifically instructed to not penalize honesty concerning limitations.

### 272 3. Theory assumptions and proofs

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

275 Answer: [NA]

276 Justification: The paper does not include theoretical results.

277 Guidelines:

- 278 • The answer NA means that the paper does not include theoretical results.
- 279 • All the theorems, formulas, and proofs in the paper should be numbered and cross-  
280 referenced.
- 281 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 282 • The proofs can either appear in the main paper or the supplemental material, but if  
283 they appear in the supplemental material, the authors are encouraged to provide a short  
284 proof sketch to provide intuition.
- 285 • Inversely, any informal proof provided in the core of the paper should be complemented  
286 by formal proofs provided in appendix or supplemental material.
- 287 • Theorems and Lemmas that the proof relies upon should be properly referenced.

### 288 4. Experimental result reproducibility

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

292 Answer: [NA]

293 Justification: No experiments.

294 Guidelines:

- 295 • The answer NA means that the paper does not include experiments.
- 296 • If the paper includes experiments, a No answer to this question will not be perceived
- 297 well by the reviewers: Making the paper reproducible is important, regardless of
- 298 whether the code and data are provided or not.
- 299 • If the contribution is a dataset and/or model, the authors should describe the steps taken
- 300 to make their results reproducible or verifiable.
- 301 • Depending on the contribution, reproducibility can be accomplished in various ways.
- 302 For example, if the contribution is a novel architecture, describing the architecture fully
- 303 might suffice, or if the contribution is a specific model and empirical evaluation, it may
- 304 be necessary to either make it possible for others to replicate the model with the same
- 305 dataset, or provide access to the model. In general, releasing code and data is often
- 306 one good way to accomplish this, but reproducibility can also be provided via detailed
- 307 instructions for how to replicate the results, access to a hosted model (e.g., in the case
- 308 of a large language model), releasing of a model checkpoint, or other means that are
- 309 appropriate to the research performed.
- 310 • While NeurIPS does not require releasing code, the conference does require all submissions
- 311 to provide some reasonable avenue for reproducibility, which may depend on the
- 312 nature of the contribution. For example
  - 313 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
  - 314 to reproduce that algorithm.
  - 315 (b) If the contribution is primarily a new model architecture, the paper should describe
  - 316 the architecture clearly and fully.
  - 317 (c) If the contribution is a new model (e.g., a large language model), then there should
  - 318 either be a way to access this model for reproducing the results or a way to reproduce
  - 319 the model (e.g., with an open-source dataset or instructions for how to construct
  - 320 the dataset).
  - 321 (d) We recognize that reproducibility may be tricky in some cases, in which case
  - 322 authors are welcome to describe the particular way they provide for reproducibility.
  - 323 In the case of closed-source models, it may be that access to the model is limited in
  - 324 some way (e.g., to registered users), but it should be possible for other researchers
  - 325 to have some path to reproducing or verifying the results.

## 326 5. Open access to data and code

327 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
 328 tions to faithfully reproduce the main experimental results, as described in supplemental  
 329 material?

330 Answer: [NA]

331 Justification: No code.

332 Guidelines:

- 333 • The answer NA means that paper does not include experiments requiring code.
- 334 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/  
 335 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 336 • While we encourage the release of code and data, we understand that this might not be
- 337 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not
- 338 including code, unless this is central to the contribution (e.g., for a new open-source
- 339 benchmark).
- 340 • The instructions should contain the exact command and environment needed to run to
- 341 reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 342 • The authors should provide instructions on data access and preparation, including how
- 343 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 344 • The authors should provide scripts to reproduce all experimental results for the new
- 345 proposed method and baselines. If only a subset of experiments are reproducible, they
- 346 should state which ones are omitted from the script and why.
- 347 • At submission time, to preserve anonymity, the authors should release anonymized
- 348 versions (if applicable).
- 349

- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

## 6. Experimental setting/details

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

Answer: [NA]

Justification: No experiments.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

## 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: [NA]

Justification: No experiments.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

## 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: [NA]

Justification: No experiments.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.

- 400 • The paper should provide the amount of compute required for each of the individual  
401 experimental runs as well as estimate the total compute.  
402 • The paper should disclose whether the full research project required more compute  
403 than the experiments reported in the paper (e.g., preliminary or failed experiments that  
404 didn't make it into the paper).

#### 405 9. Code of ethics

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

408 Answer: [Yes]

409 Justification: No harm with the paper.

410 Guidelines:

- 411 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.  
412 • If the authors answer No, they should explain the special circumstances that require a  
413 deviation from the Code of Ethics.  
414 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-  
415 eration due to laws or regulations in their jurisdiction).

#### 416 10. Broader impacts

417 Question: Does the paper discuss both potential positive societal impacts and negative  
418 societal impacts of the work performed?

419 Answer: [Yes]

420 Justification: The paper discuss both potential positive societal impacts and negative societal  
421 impacts of the work performed

422 Guidelines:

- 423 • The answer NA means that there is no societal impact of the work performed.  
424 • If the authors answer NA or No, they should explain why their work has no societal  
425 impact or why the paper does not address societal impact.  
426 • Examples of negative societal impacts include potential malicious or unintended uses  
427 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations  
428 (e.g., deployment of technologies that could make decisions that unfairly impact specific  
429 groups), privacy considerations, and security considerations.  
430 • The conference expects that many papers will be foundational research and not tied  
431 to particular applications, let alone deployments. However, if there is a direct path to  
432 any negative applications, the authors should point it out. For example, it is legitimate  
433 to point out that an improvement in the quality of generative models could be used to  
434 generate deepfakes for disinformation. On the other hand, it is not needed to point out  
435 that a generic algorithm for optimizing neural networks could enable people to train  
436 models that generate Deepfakes faster.  
437 • The authors should consider possible harms that could arise when the technology is  
438 being used as intended and functioning correctly, harms that could arise when the  
439 technology is being used as intended but gives incorrect results, and harms following  
440 from (intentional or unintentional) misuse of the technology.  
441 • If there are negative societal impacts, the authors could also discuss possible mitigation  
442 strategies (e.g., gated release of models, providing defenses in addition to attacks,  
443 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from  
444 feedback over time, improving the efficiency and accessibility of ML).

#### 445 11. Safeguards

446 Question: Does the paper describe safeguards that have been put in place for responsible  
447 release of data or models that have a high risk for misuse (e.g., pretrained language models,  
448 image generators, or scraped datasets)?

449 Answer: [NA]

450 Justification: No data no model.

451 Guidelines:

- 452 • The answer NA means that the paper poses no such risks.
- 453 • Released models that have a high risk for misuse or dual-use should be released with
- 454 necessary safeguards to allow for controlled use of the model, for example by requiring
- 455 that users adhere to usage guidelines or restrictions to access the model or implementing
- 456 safety filters.
- 457 • Datasets that have been scraped from the Internet could pose safety risks. The authors
- 458 should describe how they avoided releasing unsafe images.
- 459 • We recognize that providing effective safeguards is challenging, and many papers do
- 460 not require this, but we encourage authors to take this into account and make a best
- 461 faith effort.

## 462 12. Licenses for existing assets

463 Question: Are the creators or original owners of assets (e.g., code, data, models), used in  
464 the paper, properly credited and are the license and terms of use explicitly mentioned and  
465 properly respected?

466 Answer: [NA]

467 Justification: No code no data.

468 Guidelines:

- 469 • The answer NA means that the paper does not use existing assets.
- 470 • The authors should cite the original paper that produced the code package or dataset.
- 471 • The authors should state which version of the asset is used and, if possible, include a
- 472 URL.
- 473 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 474 • For scraped data from a particular source (e.g., website), the copyright and terms of
- 475 service of that source should be provided.
- 476 • If assets are released, the license, copyright information, and terms of use in the
- 477 package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets)
- 478 has curated licenses for some datasets. Their licensing guide can help determine the
- 479 license of a dataset.
- 480 • For existing datasets that are re-packaged, both the original license and the license of
- 481 the derived asset (if it has changed) should be provided.
- 482 • If this information is not available online, the authors are encouraged to reach out to
- 483 the asset's creators.

## 484 13. New assets

485 Question: Are new assets introduced in the paper well documented and is the documentation  
486 provided alongside the assets?

487 Answer: [NA]

488 Justification: No new assets.

489 Guidelines:

- 490 • The answer NA means that the paper does not release new assets.
- 491 • Researchers should communicate the details of the dataset/code/model as part of their
- 492 submissions via structured templates. This includes details about training, license,
- 493 limitations, etc.
- 494 • The paper should discuss whether and how consent was obtained from people whose
- 495 asset is used.
- 496 • At submission time, remember to anonymize your assets (if applicable). You can either
- 497 create an anonymized URL or include an anonymized zip file.

## 498 14. Crowdsourcing and research with human subjects

499 Question: For crowdsourcing experiments and research with human subjects, does the paper  
500 include the full text of instructions given to participants and screenshots, if applicable, as  
501 well as details about compensation (if any)?

502 Answer: [NA]

503 Justification: Not an experimental paper.

504 Guidelines:

- 505 • The answer NA means that the paper does not involve crowdsourcing nor research with  
506 human subjects.
- 507 • Including this information in the supplemental material is fine, but if the main contribu-  
508 tion of the paper involves human subjects, then as much detail as possible should be  
509 included in the main paper.
- 510 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,  
511 or other labor should be paid at least the minimum wage in the country of the data  
512 collector.

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: [NA]

520 Justification: No participants.

521 Guidelines:

- 522 • The answer NA means that the paper does not involve crowdsourcing nor research with  
523 human subjects.
- 524 • Depending on the country in which research is conducted, IRB approval (or equivalent)  
525 may be required for any human subjects research. If you obtained IRB approval, you  
526 should clearly state this in the paper.
- 527 • We recognize that the procedures for this may vary significantly between institutions  
528 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the  
529 guidelines for their institution.
- 530 • For initial submissions, do not include any information that would break anonymity (if  
531 applicable), such as the institution conducting the review.

532 **16. Declaration of LLM usage**

533 Question: Does the paper describe the usage of LLMs if it is an important, original, or  
534 non-standard component of the core methods in this research? Note that if the LLM is used  
535 only for writing, editing, or formatting purposes and does not impact the core methodology,  
536 scientific rigor, or originality of the research, declaration is not required.

537 Answer: [NA]

538 Justification: The core method development in this research does not involve LLMs as any  
539 important, original, or non-standard components.

540 Guidelines:

- 541 • The answer NA means that the core method development in this research does not  
542 involve LLMs as any important, original, or non-standard components.
- 543 • Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>)  
544 for what should or should not be described.