
CogniBias: A Benchmark for Cognitive Biases in AI–Human Dialogue

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Cognitive biases are predictable departures from rational judgment that impact
2 people’s decisions and communication. As large language models (LLMs) are
3 increasingly interfaced with everyday interactions, it is valuable to understand how
4 biases arise and spread through AI–human conversations. We present CogniBias,
5 the first benchmark dataset to study and assess cognitive biases in conversational
6 settings. CogniBias is composed of 30+ established types of biases, like anchoring,
7 framing, confirmation bias, and optimism bias, distilled from a variety of authentic
8 question and answer scenarios. Each dialogue sample includes an LLM suggestion,
9 a human-like response, and expert-informed annotations with notes, bias labels, and
10 confidence scores. To establish the benchmark, we describe a generation pipeline
11 incorporating multiple LLMs and include baseline results with pretrained and
12 fine-tuned models on classification and detection tasks. Our analysis highlights the
13 identified challenges in detecting subtle biases or overlapping biases and identified
14 each model’s frequent failures. We hope that by releasing CogniBias, it will align
15 divergent perspectives on cognitive bias assessment and establish a baseline dataset
16 for fairer, more trustworthy conversational AI systems.

17 1 Introduction

18 Cognitive biases are systematic deviations from rational judgment that shape perception, decision-
19 making, and interaction, including well-studied cases such as anchoring, confirmation bias, framing
20 effects, and overconfidence (6). While these insights from psychology are valuable for understanding
21 human reasoning, they also pose risks in AI–human interaction, where conversational systems may
22 amplify existing user biases or introduce new ones. Although recent NLP research has advanced
23 fairness, bias detection, and debiasing, most benchmarks (e.g., WinoBias, StereoSet, CrowS-Pairs)
24 focus on demographic or stereotypical biases (10), overlooking cognitive biases that emerge in
25 decision-making and dialogue. Prior datasets on framing or persuasion capture isolated phenomena
26 within narrow domains, but they lack the breadth needed to generalize across multiple categories of
27 cognitive bias.

28 To address this gap, we present CogniBias, a benchmark dataset for evaluating cognitive biases
29 in AI–human dialogue. CogniBias integrates over 30 bias types from established psychological
30 taxonomies (e.g., anchoring, availability heuristic, sunk cost fallacy), instantiated in decision-making
31 dialogues. Each dialogue consists of a decision-oriented question, an AI-generated suggestion, and
32 a human-like response. Expert annotations provide bias labels and confidence scores, enabling
33 systematic study of how biases emerge and interact in machine-mediated dialogue (8).

34 Our contributions are threefold:

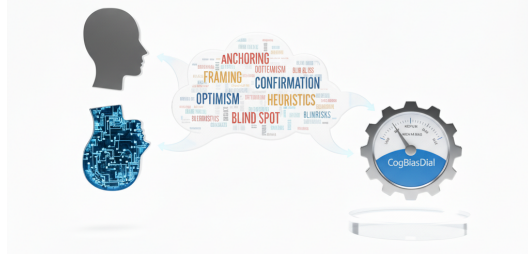


Figure 1: Motivation for CogniBias. Conversational AI may introduce or amplify cognitive biases in user decision-making. CogniBias provides the first benchmark to systematically study such effects.

- 35 • CogniBias dataset: A large-scale, multi-bias benchmark covering 30+ types of cognitive
- 36 biases in AI–human dialogue.
- 37 • Generation pipeline: A reproducible LLM-based framework combining dialogue simulation
- 38 with expert-informed annotations.

39 By releasing CogniBias, we aim to foster cross-disciplinary research at the intersection of cognitive
 40 psychology, NLP, and human–AI interaction. We envision this resource as a foundation for building
 41 conversational AI systems that are accurate, fair, and free from cognitive distortions that undermine
 42 trust and decision quality.

43 2 Related work

44 The NLP field has developed various benchmarks which have aimed to understand and measure
 45 bias in language models. This work has, at least to date, focused on representational-type harms
 46 related to bias based on gender, race, and/or culture. For instance, StereoSet, CrowS-Pairs and
 47 WinoBias (7; 13; 10) showed that pretrained models reinforced systemic stereotypes, and inspired
 48 new research in debiasing. However, these benchmarks and the corpus based on them, only assess
 49 demographic and social stereotypes and do not measure the cognitive processes underlying reasoning
 50 and/or decision-making.

51 Cognitive psychology has long documented systematic deviations in judgment, such as anchoring,
 52 framing, and the availability heuristic (5), yet these cognitive biases are largely absent from computa-
 53 tional evaluation benchmarks. Existing NLP studies on persuasion, framing, or fallacies typically
 54 focus on one or two bias types within narrow domains like news, politics, or social media, limiting
 55 their generalizability across bias taxonomies and failing to capture conversational dynamics.

56 Recent research (11; 12) has indicated that large language models (LLMs) not only inherit social
 57 stereotypes, but can also exacerbate cognitive biases in dialog. For example, an anchoring effect
 58 can occur when users overweight an initial suggestion from a model, and a framing effect can bias
 59 user judgment depending on which way the information is presented. Regardless of these important
 60 considerations, there is not yet a systematic benchmark on measuring cognitive biases in a human
 61 interactive AI setting. There has been a focus on ethical risk or persuasive misuse of interacting with
 62 an LLM but no empirical datasets to use for a thorough examination of those concepts.

63 CogniBias differs from prior bias benchmarks by addressing cognitive rather than demographic biases
 64 within dialogue contexts 1.

65 3 Dataset design

66 3.1 Bias Taxonomy

67 CogniBias is founded on a classification of 30+ cognitive biases situated in cognitive psychology.
 68 These include classic heuristics and judgment errors such as anchoring, confirmation bias, framing,
 69 availability heuristic, gambler’s fallacy, sunk cost fallacy, halo effect, overconfidence, optimism bias,
 70 and spotlight effect, as well as various other biases. The classification was an aggregation from the
 71 psychology literature and was modified slightly for conversational contexts to ensure that biases

Table 1: Comparative analysis of bias datasets. CogniBias uniquely focuses on cognitive biases in AI–human dialogue, unlike prior benchmarks addressing demographic stereotypes.

Aspect	CogniBias	StereoSet	CrowS-Pairs	WinoBias
Primary Focus	Cognitive biases in AI–human dialogues (e.g., anchoring, framing, optimism bias)	Stereotypical social biases (gender, race, religion, profession)	Social biases across demographic categories	Gender bias in coreference resolution
Bias Type	30+ cognitive bias types from psychology	4 social bias domains	9 social bias categories	Gender stereotypes only
Data Format	Three-turn dialogues: <i>decision question</i> → <i>AI suggestion</i> → <i>human-like response</i>	Triplets: <i>stereotype / anti-stereotype / unrelated</i>	Paired sentences: <i>stereotype vs. anti-stereotype</i>	Sentence templates with gendered pronouns
Scale	3,000 dialogues; 30+ bias types; 50 simulated participants	16,995 instances; 321 target terms	1,508 pairs; 9 bias types	396 examples (2 test sets)
Domains Covered	Finance, health, lifestyle, and social decision-making	Gender, race, religion, profession	Race, gender, religion, age, etc.	Occupation–gender stereotypes
Interaction Type	Conversational — includes both AI and human-like turns	Static sentence or discourse	Static sentence pairs	Static coreference sentences
Uniqueness	First dataset to study cognitive biases in AI–human dialogue	Measures demographic stereotypes only	Broad demographic coverage but not conversational	Template-based; focused narrowly on gender
Human–AI Interaction	Present — AI output may influence human bias	Absent	Absent	Absent

72 covered a wide range of reasoning errors. We code each dialogue with exactly one main bias type (or
73 No Bias, when appropriate) (6).

74 3.2 Dialogue Format

75 Each dataset instance follows a three-part dialogue structure:

- 76 • Question of action: A plausible situation that calls for decision-making (e.g., "Is it an
77 opportune moment for me to invest in cryptocurrency?").
- 78 • AI proposal: The model-generated response emulates the kinds of advice, or context framing,
79 that an LLM might use.
- 80 • Response of a human-like persona: A responding utterance replicates how a person would
81 respond directly in a moment of cognitive bias.

82 Annotations include:

- 83 • Bias label (one of the 30+ categories)
- 84 • Confidence score (0–100), reflecting annotator certainty,
- 85 • Optional notes, used during quality control.

86 This design captures interactive dynamics where AI output may influence or reinforce human cognitive
87 distortions.

88 3.3 Generation Pipeline

89 To create CogniBias, we implemented a two-stage generation pipeline:

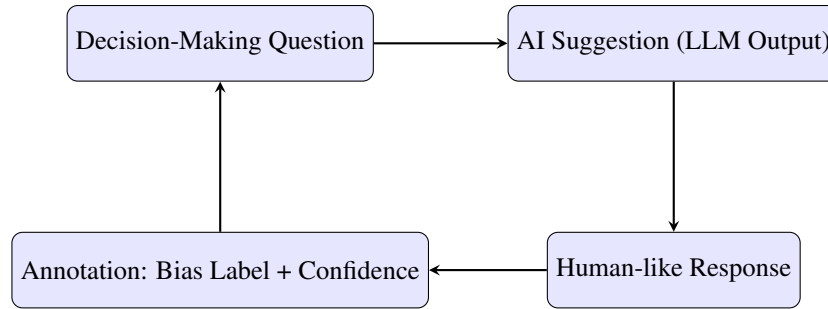


Figure 2: Overview of the CogniBias dataset generation pipeline. Each instance includes a decision-making question, an AI-generated suggestion, and a human-like response annotated with bias type and confidence score.

- 90 • Prompted LLM generation: A number of large language models (Gemini, Groq’s LLaMA
- 91 variants, etc.) were prompted with a pool of 60+ fallback questions covering the domains
- 92 of finance, health, lifestyle, and social decision-making. The models generated both AI-
- 93 generated suggestions and participant-like responses that varied stylistically and with chain
- 94 of reasoning.
- 95 • Expert-informed annotation: Sampled each dialogue for one or more candidate bias from
- 96 the taxonomy. Disagreement was resolved by discussion and there are now a bias label and
- 97 confidence score in the final version of the annotations.

98 This hybrid pipeline has the benefits of leveraging the scale of LLM simulation, while using experts
 99 to maintain fidelity to psychological definitions.

100 3.4 Dataset Statistics

101 The released dataset includes:

- 102 • Participants: 50 simulated participants, each with a dialogue history (e.g., P01, P26)
- 103 • Dialogues: 3000 annotated dialogue samples across diverse domains.
- 104 • Bias coverage: All 30+ bias categories are represented, with frequency distributions available
- 105 in Appendix A.
- 106 • Annotations: Every instance contains a bias type and confidence score (mean confidence 73
- 107 across the dataset).

108 3.5 Quality Control

109 Each annotation in CogniBias was validated against cognitive psychology definitions to ensure
 110 accurate bias categorization. To reduce overfitting, the dataset also includes instances labeled as
 111 No Bias. Furthermore, the questions span diverse domains such as finance, health, relationships,
 112 consumer choices, and ethics, ensuring broad applicability.

113 4 Challenges & Future Use

114 4.1 Challenges and Future Use

115 We highlight several open challenges for future research. Many biases are expressed subtly and
 116 implicitly, making them difficult to identify and model. Categories can overlap, as a single utterance
 117 may fit multiple bias types, and some biases occur more frequently than others, creating class
 118 imbalance. Additionally, context sensitivity complicates modeling, since the same surface form may
 119 signal different biases depending on the conversational setting. Future work could explore bias-
 120 adaptive prompting strategies, where large language models dynamically adjust their communication
 121 style to mitigate bias triggers by reframing questions, providing counter-examples, or requesting

122 supporting evidence. Another promising direction is to treat cognitive bias detection as a multi-
123 task, hierarchical NLP problem that captures both sentence-level and dialogue-level phenomena,
124 accounting for overlapping biases and leveraging relationships among cognitive bias categories.
125 By surfacing these challenges and directions, CogniBias provides not only a benchmark but also a
126 foundation for advancing methods in bias detection, dialogue modeling, and cognitively aware AI
127 evaluation.

128 **5 Discussion & Limitations**

129 CogniBias provides a new benchmark for evaluating cognitive biases in conversations, yet several
130 limitations should be noted. Annotation ambiguity arises since biases often overlap (e.g., framing and
131 optimism), and we chose to label only the dominant bias to simplify the task, which sacrifices some
132 realism in representing the complexity of human reasoning (9). Our taxonomy covers over 30 biases
133 but does not capture the full spectrum of distortions (e.g., moral licensing, hot-cold empathy gap) due
134 to annotation complexity and limited examples (1). Furthermore, all dialogues are LLM-generated,
135 which, while enabling scalability and control, reflect model-specific linguistic tendencies rather
136 than fully naturalistic human behavior (12). The dataset also has uneven domain balance, with
137 certain topics such as consumer decisions being overrepresented compared to others like health or
138 finance. Annotation confidence averaged around 70, highlighting both the difficulty and subjectivity
139 of identifying cognitive biases, which in turn increases the challenge of evaluation. Despite these
140 limitations, CogniBias remains a systematic and scalable resource to study bias in dialogue, offering a
141 foundation for research on bias detection, conversational fairness, and evaluation of language models
142 grounded in psychological theory.

143 **6 Broader Impact & Ethics**

144 CogniBias is released to support research on how cognitive biases emerge in AI–human dialogue, with
145 the aim of improving detection tools, bias-aware training, and the trustworthiness of conversational
146 systems (11). While such a resource can advance transparency, equity, and reliability, it also carries
147 risks of misuse, such as exploiting anchoring or framing for manipulation. To mitigate this, the dataset
148 is framed strictly for evaluation and responsible use (12). Since the dialogues are simulated with
149 LLMs and expert annotation, they do not fully capture cultural or contextual diversity, but future work
150 could expand coverage across languages and settings (2). Importantly, no personal data are included,
151 avoiding privacy concerns associated with scraped corpora (4). With these safeguards, CogniBias is
152 intended to encourage responsible progress in conversational AI research while acknowledging both
153 its potential and limitations. By situating biases in dialogue explicitly, it also invites interdisciplinary
154 collaboration between AI, psychology, and ethics. This, in turn, can help inform design guidelines for
155 safer systems. Ultimately, we hope CogniBias will serve not just as a benchmark but as a foundation
156 for building more equitable and trustworthy AI dialogue technologies.

157 **7 Conclusion**

158 We present CogniBias, the first benchmark dataset designed for studying cognitive biases in AI–
159 human dialogues. It includes over 30 types of biases identified in psychology and about 3,000
160 annotated dialogues, each consisting of a decision-oriented question, an AI recommendation, and a
161 human-like response labeled for bias type and confidence. CogniBias enables systematic investigation
162 of how biases manifest in conversational contexts and supports the development of bias detection or
163 mitigation models. Our analysis highlights challenges such as overlapping categories, annotation
164 ambiguity, and the complexity of conversational cues. By profiling cognitive biases in a structured way,
165 CogniBias aims to inform future research across cognitive psychology, natural language processing,
166 and human–AI interaction, fostering the development of conversational AI systems that are not only
167 fluent and accurate but also fair, trustworthy, and aligned with human reasoning.

168 **References**

169 [1] Haselton, M. G., Nettle, D., & Andrews, P. W. (2015). The evolution of cognitive bias. In D. M.
170 Buss (Ed.), *The handbook of evolutionary psychology* (pp. 724–746).

- 171 [2] R. Segerer. Cultural value alignment in large language models: A prompt-based analysis of
172 Schwartz values in Gemini, ChatGPT, and DeepSeek. *arXiv preprint arXiv:2505.17112*, 2025.
- 173 [3] Singh, S., Singla, Y. K., Si, H., & Krishnamurthy, B. (2024). Measuring and improving
174 persuasiveness of large language models. *arXiv preprint arXiv:2410.02653*.
- 175 [4] Gebru, T., Morgenstern, J., Vecchione, B., Vaughan, J. W., Wallach, H., Daumé III, H., &
176 Crawford, K. (2018). Datasheets for datasets. *Proceedings of the 5th Workshop on Fairness,
177 Accountability, and Transparency in Machine Learning (FAT/ML)*, 1–10.
- 178 [5] Kahneman, D. (2011). *Fast and slow thinking*. Allen Lane and Penguin Books, New York.
- 179 [6] Tversky, A. & Kahneman, D. (1974). Judgment under uncertainty: Heuristics and biases.
180 *Science*, 185(4157), 1124–1131.
- 181 [7] Asad, N., Sahoo, N. R., Murthy, R., Nath, S., & Bhattacharyya, P. (2025, July). “You are
182 Beautiful, Body Image Stereotypes are Ugly!” BISTereo: A Benchmark to Measure Body Image
183 Stereotypes in Language Models. In *Findings of the Association for Computational Linguistics:
184 ACL 2025* (pp. 24471–24496).
- 185 [8] Chen, H., Ghosal, D., Majumder, N., Hussain, A., & Poria, S. (2021). Persuasive dialogue
186 understanding: The baselines and negative results. *Neurocomputing*, 431, 47–56.
- 187 [9] Pavlick, E., & Kwiatkowski, T. (2019). Inherent disagreements in human textual inferences.
188 *Transactions of the Association for Computational Linguistics*, 7, 677–694.
- 189 [10] Zhao, J., Wang, T., Yatskar, M., Ordonez, V., & Chang, K. W. (2018). Gender bias in coreference
190 resolution: Evaluation and mitigation. *EMNLP*, 2018.
- 191 [11] Sheng, E., Chang, K.-W., Natarajan, P., & Peng, N. (2019). The Woman Worked as a Babysitter:
192 On Bias in Language Models. *Proceedings of the 2019 Conference on Empirical Methods in
193 Natural Language Processing (EMNLP)*, 3407–3412.
- 194 [12] Bender, E. M., Gebru, T., McMillan-Major, A., & Shmitchell, S. (2021). On the dangers
195 of stochastic parrots: Can language models be too big? In *Proceedings of the 2021 ACM
196 Conference on Fairness, Accountability, and Transparency (FAccT)*, 610–623.
- 197 [13] Nangia, N., Vania, C., Bhlerao, R., & Bowman, S. R. (2020). CrowS-pairs: A challenge dataset
198 for measuring social biases in masked language models. *arXiv preprint arXiv:2010.00133*.

199 **A Bias Taxonomy**

200 Table 2 lists the cognitive biases covered in COGNIBIAS with concise definitions.

Bias Type	Definition
Anchoring	Relying too much on the first piece of information encountered.
Availability Heuristic	Judging likelihood by easily recalled examples.
Confirmation Bias	Seeking or interpreting information that confirms prior beliefs.
Framing Effect	Making different decisions depending on how information is presented.
Hindsight Bias	Viewing past events as more predictable than they were.
Loss Aversion	Preferring to avoid losses rather than acquire equivalent gains.
Status Quo Bias	Preferring things to stay the same instead of changing.
Optimism Bias	Overestimating the likelihood of positive outcomes.
Pessimism Bias	Overestimating the likelihood of negative outcomes.
Bandwagon Effect	Adopting beliefs or behaviors because others do.
Sunk Cost Fallacy	Continuing due to prior investments despite poor prospects.
Gambler’s Fallacy	Believing past random events affect future probabilities.
Overconfidence	Overestimating one’s knowledge, ability, or predictions.
Halo Effect	Letting one positive trait influence overall judgment.
Self-Serving Bias	Attributing success to oneself, failures to external factors.
Dunning–Kruger Effect	Overestimating competence due to limited knowledge.
Negativity Bias	Giving greater weight to negative experiences over positive ones.
Survivorship Bias	Focusing on successful cases while ignoring failures.
Authority Bias	Placing undue weight on the opinions of authority figures.
Recency Bias	Overemphasizing recent events compared to older information.
Outcome Bias	Judging decisions by their outcomes instead of reasoning quality.
Planning Fallacy	Underestimating the time or resources needed for tasks.
Spotlight Effect	Overestimating how much others notice one’s actions or appearance.
Illusory Correlation	Perceiving a relationship between unrelated events.
Base Rate Fallacy	Ignoring statistical base rates in favor of anecdotal evidence.

Table 2: Concise taxonomy of cognitive biases in COGNIBIAS.

201 **B Dataset Statistics**

202 Table 3 reports the frequency and average annotator confidence for all 30+ cognitive biases in
203 COGNIBIAS.

204 **B.1 Domain Distribution**

205 The dataset currently covers four domains: finance, health, lifestyle, and social decision-making.
206 In this release, domain labels are under-specified (annotated as “unknown” in the JSON files), but
207 prompts were designed to ensure cross-domain diversity. Future expansions will include explicit
208 domain annotations for finer-grained analysis.

209 **C Annotation Guidelines**

210 This section summarizes the annotation instructions provided to expert annotators for labeling
211 dialogues in COGNIBIAS. The goal was to ensure consistent application of cognitive bias categories
212 across 30+ types drawn from psychology.

213 **C.1 General Instructions**

214 Annotators were asked to:

- 215 • Read the dialogue in full (decision-oriented question, AI suggestion, and human-like re-
216 sponse).
- 217 • Identify whether the human-like response exhibits a cognitive bias.

Bias Type	Count	Avg. Confidence
Pessimism Bias	121	60.4
Optimism Bias	110	82.4
Framing Effect	108	74.8
Spotlight Effect	108	66.0
No Bias	104	70.8
Hindsight Bias	103	64.4
Recency Bias	100	65.0
Planning Fallacy	98	53.8
Illusion of Control	96	65.3
Status Quo Bias	95	67.0
Confirmation Bias	95	87.2
False Consensus Effect	95	63.9
Selection Bias	94	64.7
Pro-Innovation Bias	93	65.6
Outcome Bias	92	65.6
Just-World Hypothesis	91	65.0
Sunk Cost Fallacy	91	65.9
Dunning–Kruger Effect	90	94.2
Gambler’s Fallacy	89	64.0
Halo Effect	87	77.1
Loss Aversion	86	79.0
Bandwagon Effect	85	74.9
Negativity Bias	84	64.3
Authority Bias	84	82.0
Actor–Observer Bias	82	67.0
Self-Serving Bias	81	87.0
Base Rate Fallacy	81	63.8
Overconfidence	80	92.0
Illusory Correlation	80	65.4
Anchoring	79	77.7
Availability Heuristic	78	69.7
Group Attribution Error	71	65.5
Survivorship Bias	69	65.0

Table 3: Frequency and average annotator confidence for all cognitive bias categories in COGNIBIAS.

- 218 • If biased, assign exactly one **dominant bias type** from the taxonomy (see Appendix A).
- 219 • If the response was neutral, label it as **No Bias**.
- 220 • Assign a **confidence score** (0–100) indicating certainty in the label.
- 221 • Provide optional notes explaining the rationale in ambiguous cases.

222 C.2 Bias Labeling Criteria

- 223 • **Dominant bias rule:** When multiple biases could apply, annotators selected the one most
- 224 central to the response.
- 225 • **Surface vs. latent cues:** Explicit linguistic markers (e.g., “since it was first mentioned”)
- 226 were prioritized, but pragmatic and discourse-level reasoning (e.g., optimism bias without
- 227 explicit markers) was also considered valid.
- 228 • **No Bias category:** Responses that were rational, balanced, or contextually neutral were
- 229 marked as *No Bias*.

230 C.3 Resolving Disagreements

- 231 • Each dialogue was annotated independently by two experts.
- 232 • Disagreements were discussed in weekly calibration meetings.
- 233 • In cases of persistent disagreement, the annotation was adjudicated by a third expert.
- 234 • Confidence scores were averaged across annotators for the final dataset.

235 These guidelines ensured that annotations were consistent, reproducible, and grounded in established
236 cognitive psychology definitions while acknowledging the inherent subjectivity in bias identification.

237 **D Dataset Format**

238 The COGNIBIAS dataset is released under an open license for research purposes. It is provided in
239 JSON format with one file per participant. Each file contains a list of dialogues, where each dialogue
240 is annotated with bias type and confidence.

241 **D.1 Fields**

242 Each dialogue instance includes the following fields:

- 243 • `participant_id` – Unique identifier for the participant (e.g., P01).
- 244 • `question_id` – Identifier for the dialogue question.
- 245 • `question` – Decision-making prompt shown to the participant.
- 246 • `ai_suggestion` – Response generated by a language model.
- 247 • `human_response` – Human-like response reflecting possible bias.
- 248 • `bias_type` – Assigned label from the taxonomy (Appendix A).
- 249 • `confidence` – Annotator confidence score (0–100).

250 **D.2 Example JSON Snippet**

Example JSON Snippet

```
{
  "participant_id": "P01",
  "dialogue": [
    {
      "question_id": "Q15",
      "question": "Should I invest in cryptocurrency right now?",
      "ai_suggestion": "It could be profitable, but it is highly
        volatile.",
      "human_response": "Yes, I ve heard so many success stories
        recently. I dont want to miss out.",
      "bias_type": "Availability Heuristic",
      "confidence": 78
    }
  ]
}
```

251

252 **E Distribution of Bias Types**

253 CogniBias encapsulates over 30 unique cognitive biases, amounting to 3,000 annotated dialogues
254 collected from 50 participants. Participants showed an average confidence of 70.6 on a scale of 0–100
255 for each annotation. The most common biases included: Pessimism Bias, Optimism Bias, Spotlight
256 Effect, Framing Effect, and No Bias (all present in >100 dialogues). Important, but less frequent
257 biases included: Halo Effect, Planning Fallacy, Authority Bias, and Illusory Correlation. Figure 3
258 displays the frequency of the top 15 most common bias types.

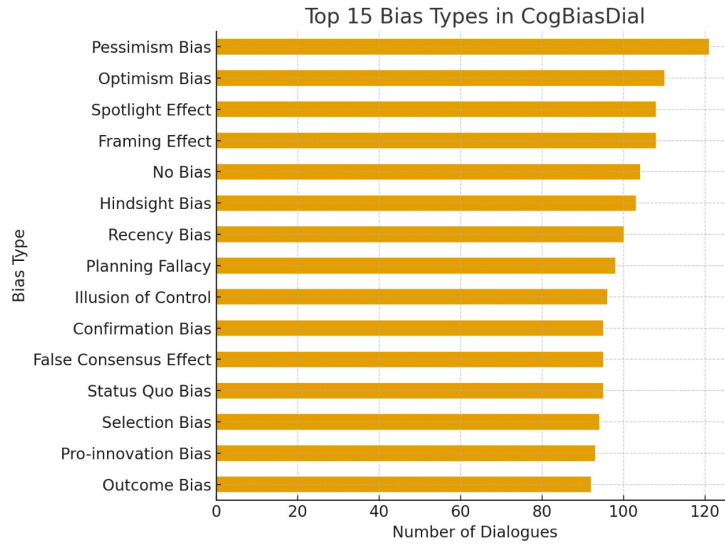


Figure 3: Distribution of the most common bias types in CogniBias. Long-tail categories ensure diversity across 30+ psychological bias types.

259 **F Dataset Access**

260 The dataset and accompanying documentation are available at: [https://github.com/Mri1306/](https://github.com/Mri1306/Cognitive-Bias-Dataset)
 261 [Cognitive-Bias-Dataset](https://github.com/Mri1306/Cognitive-Bias-Dataset)

262 **NeurIPS Paper Checklist**

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

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

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

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

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

287 **IMPORTANT, please:**

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

291 **1. Claims**

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

294 Answer: [Yes]

295 Justification: CogniBias is stated openly in the abstract and introduction as a benchmark
296 dataset for cognitive biases in dialogue, which matches the contributions (Sections 1 and 3)

297 Guidelines:

- 298 • The answer NA means that the abstract and introduction do not include the claims
299 made in the paper.
- 300 • The abstract and/or introduction should clearly state the claims made, including the
301 contributions made in the paper and important assumptions and limitations. A No or
302 NA answer to this question will not be perceived well by the reviewers.
- 303 • The claims made should match theoretical and experimental results, and reflect how
304 much the results can be expected to generalize to other settings.
- 305 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
306 are not attained by the paper.

307 **2. Limitations**

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

309 Answer: [Yes]

310 Justification: A separate "Discussion Limitations" section (Section 6) provides explicit
311 notes on annotation ambiguity, coverage gaps, domain imbalance, as well as the reliance on
312 LLM-generated dialogues.

313 Guidelines:

- 314 • The answer NA means that the paper has no limitation while the answer No means that
315 the paper has limitations, but those are not discussed in the paper.
- 316 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 317 • The paper should point out any strong assumptions and how robust the results are to
318 violations of these assumptions (e.g., independence assumptions, noiseless settings,
319 model well-specification, asymptotic approximations only holding locally). The authors
320 should reflect on how these assumptions might be violated in practice and what the
321 implications would be.
- 322 • The authors should reflect on the scope of the claims made, e.g., if the approach was
323 only tested on a few datasets or with a few runs. In general, empirical results often
324 depend on implicit assumptions, which should be articulated.
- 325 • The authors should reflect on the factors that influence the performance of the approach.
326 For example, a facial recognition algorithm may perform poorly when image resolution
327 is low or images are taken in low lighting. Or a speech-to-text system might not be
328 used reliably to provide closed captions for online lectures because it fails to handle
329 technical jargon.
- 330 • The authors should discuss the computational efficiency of the proposed algorithms
331 and how they scale with dataset size.
- 332 • If applicable, the authors should discuss possible limitations of their approach to
333 address problems of privacy and fairness.
- 334 • While the authors might fear that complete honesty about limitations might be used by
335 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
336 limitations that aren't acknowledged in the paper. The authors should use their best
337 judgment and recognize that individual actions in favor of transparency play an impor-
338 tant role in developing norms that preserve the integrity of the community. Reviewers
339 will be specifically instructed to not penalize honesty concerning limitations.

340 3. Theory assumptions and proofs

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

343 Answer: [NA]

344 Justification: This paper does not include theoretical results, theorems, or proofs, since it is
345 a dataset paper.

346 Guidelines:

- 347 • The answer NA means that the paper does not include theoretical results.
- 348 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
349 referenced.
- 350 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 351 • The proofs can either appear in the main paper or the supplemental material, but if
352 they appear in the supplemental material, the authors are encouraged to provide a short
353 proof sketch to provide intuition.
- 354 • Inversely, any informal proof provided in the core of the paper should be complemented
355 by formal proofs provided in appendix or supplemental material.
- 356 • Theorems and Lemmas that the proof relies upon should be properly referenced.

357 4. Experimental result reproducibility

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

361 Answer: [Yes]

362 Justification: Full details on dataset construction pipeline, annotation process, and analysis
363 are provided in Section 3 (Dataset Design) and 4 (Dataset Analysis). Format information of
364 the dataset is presented in Appendix E.

365 Guidelines:

- 366 • The answer NA means that the paper does not include experiments.
- 367 • If the paper includes experiments, a No answer to this question will not be perceived
368 well by the reviewers: Making the paper reproducible is important, regardless of
369 whether the code and data are provided or not.
- 370 • If the contribution is a dataset and/or model, the authors should describe the steps taken
371 to make their results reproducible or verifiable.
- 372 • Depending on the contribution, reproducibility can be accomplished in various ways.
373 For example, if the contribution is a novel architecture, describing the architecture fully
374 might suffice, or if the contribution is a specific model and empirical evaluation, it may
375 be necessary to either make it possible for others to replicate the model with the same
376 dataset, or provide access to the model. In general, releasing code and data is often
377 one good way to accomplish this, but reproducibility can also be provided via detailed
378 instructions for how to replicate the results, access to a hosted model (e.g., in the case
379 of a large language model), releasing of a model checkpoint, or other means that are
380 appropriate to the research performed.
- 381 • While NeurIPS does not require releasing code, the conference does require all submis-
382 sions to provide some reasonable avenue for reproducibility, which may depend on the
383 nature of the contribution. For example
 - 384 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
385 to reproduce that algorithm.
 - 386 (b) If the contribution is primarily a new model architecture, the paper should describe
387 the architecture clearly and fully.
 - 388 (c) If the contribution is a new model (e.g., a large language model), then there should
389 either be a way to access this model for reproducing the results or a way to reproduce
390 the model (e.g., with an open-source dataset or instructions for how to construct
391 the dataset).
 - 392 (d) We recognize that reproducibility may be tricky in some cases, in which case
393 authors are welcome to describe the particular way they provide for reproducibility.
394 In the case of closed-source models, it may be that access to the model is limited in
395 some way (e.g., to registered users), but it should be possible for other researchers
396 to have some path to reproducing or verifying the results.

397 5. Open access to data and code

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

401 Answer: [Yes]

402 Justification: The dataset and annotation scripts will be released publicly with access
403 instructions (Appendix E).

404 Guidelines:

- 405 • The answer NA means that paper does not include experiments requiring code.
- 406 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/
407 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 408 • While we encourage the release of code and data, we understand that this might not be
409 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not
410 including code, unless this is central to the contribution (e.g., for a new open-source
411 benchmark).
- 412 • The instructions should contain the exact command and environment needed to run to
413 reproduce the results. See the NeurIPS code and data submission guidelines ([https://
414 nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 415 • The authors should provide instructions on data access and preparation, including how
416 to access the raw data, preprocessed data, intermediate data, and generated data, etc.

- 417 • The authors should provide scripts to reproduce all experimental results for the new
418 proposed method and baselines. If only a subset of experiments are reproducible, they
419 should state which ones are omitted from the script and why.
- 420 • At submission time, to preserve anonymity, the authors should release anonymized
421 versions (if applicable).
- 422 • Providing as much information as possible in supplemental material (appended to the
423 paper) is recommended, but including URLs to data and code is permitted.

424 6. Experimental setting/details

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

428 Answer: [NA]

429 Justification: No training or model evaluation is reported, since the focus is dataset creation
430 and analysis, not benchmarks.

431 Guidelines:

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

437 7. Experiment statistical significance

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

440 Answer: [NA]

441 Justification: The paper does not contain experiments that require statistical significance
442 testing, only descriptive dataset statistics.

443 Guidelines:

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

464 8. Experiments compute resources

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

468 Answer: [NA]

469 Justification: The work involves dataset construction and annotation, not computationally
470 heavy experiments requiring resource reporting.

471 Guidelines:

- 472 • The answer NA means that the paper does not include experiments.
- 473 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
474 or cloud provider, including relevant memory and storage.
- 475 • The paper should provide the amount of compute required for each of the individual
476 experimental runs as well as estimate the total compute.
- 477 • The paper should disclose whether the full research project required more compute
478 than the experiments reported in the paper (e.g., preliminary or failed experiments that
479 didn't make it into the paper).

480 9. Code of ethics

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

483 Answer: [Yes]

484 Justification: All dialogues were LLM-generated or expert-annotated. No personal data or
485 privacy-sensitive content was collected. Ethical considerations are discussed in Section 7
486 (Broader Impact & Ethics).

487 Guidelines:

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

493 10. Broader impacts

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

496 Answer: [Yes]

497 Justification: Section 7 (Broader Impact & Ethics) addresses both upsides (bias detection
498 and fairness) and risks (misuse for persuasion).

499 Guidelines:

- 500 • The answer NA means that there is no societal impact of the work performed.
- 501 • If the authors answer NA or No, they should explain why their work has no societal
502 impact or why the paper does not address societal impact.
- 503 • Examples of negative societal impacts include potential malicious or unintended uses
504 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
505 (e.g., deployment of technologies that could make decisions that unfairly impact specific
506 groups), privacy considerations, and security considerations.
- 507 • The conference expects that many papers will be foundational research and not tied
508 to particular applications, let alone deployments. However, if there is a direct path to
509 any negative applications, the authors should point it out. For example, it is legitimate
510 to point out that an improvement in the quality of generative models could be used to
511 generate deepfakes for disinformation. On the other hand, it is not needed to point out
512 that a generic algorithm for optimizing neural networks could enable people to train
513 models that generate Deepfakes faster.
- 514 • The authors should consider possible harms that could arise when the technology is
515 being used as intended and functioning correctly, harms that could arise when the
516 technology is being used as intended but gives incorrect results, and harms following
517 from (intentional or unintentional) misuse of the technology.
- 518 • If there are negative societal impacts, the authors could also discuss possible mitigation
519 strategies (e.g., gated release of models, providing defenses in addition to attacks,
520 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
521 feedback over time, improving the efficiency and accessibility of ML).

522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573

11. Safeguards

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

Answer: [Yes]

Justification: The dataset is framed for evaluation only, not persuasion or manipulation. No personal or sensitive data is included (Section 7).

Guidelines:

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

12. Licenses for existing assets

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

Answer: [Yes]

Justification: Related benchmark datasets (e.g., WinoBias, StereoSet, CrowS-Pairs) are cited in Section 2 (Related Work). No external copyrighted data was reused.

Guidelines:

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

13. New assets

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

Answer: [Yes]

Justification: CogniBias is a new dataset, documented in Section 3 (Dataset Design), Section 4 (Analysis), and Appendix E (Dataset Access & Format)

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.

- 574 • The paper should discuss whether and how consent was obtained from people whose
575 asset is used.
576 • At submission time, remember to anonymize your assets (if applicable). You can either
577 create an anonymized URL or include an anonymized zip file.

578 14. **Crowdsourcing and research with human subjects**

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

582 Answer: [NA]

583 Justification: No crowdsourcing was used. Annotations were performed by domain experts.

584 Guidelines:

- 585 • The answer NA means that the paper does not involve crowdsourcing nor research with
586 human subjects.
- 587 • Including this information in the supplemental material is fine, but if the main contribu-
588 tion of the paper involves human subjects, then as much detail as possible should be
589 included in the main paper.
- 590 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
591 or other labor should be paid at least the minimum wage in the country of the data
592 collector.

593 15. **Institutional review board (IRB) approvals or equivalent for research with human 594 subjects**

595 Question: Does the paper describe potential risks incurred by study participants, whether
596 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
597 approvals (or an equivalent approval/review based on the requirements of your country or
598 institution) were obtained?

599 Answer: [NA]

600 Justification: No human-subject data collection was involved; all dialogues were generated
601 with LLMs.

602 Guidelines:

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

613 16. **Declaration of LLM usage**

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

618 Answer: [Yes]

619 Justification: Section 3 (Dataset Design) explicitly describes that dialogues were generated
620 using large language models and then annotated by experts.

621 Guidelines:

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