MALIBU Benchmark: Multi-Agent LLM Implicit Bias Uncovered

Anonymous ACL submission

Abstract

Multi-agent systems, which consist of multiple AI models interacting within a shared environment, are increasingly used for personabased interactions. However, if not carefully designed, these systems can reinforce implicit biases in large language models (LLMs), raising concerns about fairness and equitable representation. We present MALIBU¹, a novel benchmark developed to assess the degree to which LLM-based multi-agent systems implicitly reinforce social biases and stereotypes. MAL-011 IBU evaluates bias in LLM-based multi-agent systems through scenario-based assessments. 014 AI models complete tasks within predefined contexts, and their responses undergo evaluation by an LLM-based multi-agent judging system in two phases. In the first phase, judges 018 score responses labeled with specific demo-019 graphic personas (e.g., gender, race, religion) across four metrics. In the second phase, judges compare paired responses assigned to different personas, scoring them and selecting the superior response. Our study quantifies biases in LLM-generated outputs, revealing that bias mitigation may favor marginalized personas over true neutrality, emphasizing the need for 027 nuanced detection, balanced fairness strategies, and transparent evaluation benchmarks in multiagent systems.

1 Introduction

034

035

Implicit biases are unconscious attitudes or stereotypes that can contradict conscious beliefs but still shape perceptions and decisions (Greenwald and Krieger, 2006). Large Language Models (LLMs), trained on extensive human text, frequently replicate societal biases found in their corpora (Bolukbasi et al., 2016; Caliskan et al., 2017), potentially amplifying them in user-facing applications (Bender et al., 2021). Unlike explicit biases, which are overt and more easily addressed, implicit biases are subtler and require nuanced strategies for detection and mitigation (Kurita et al., 2019). LLMs integrate into multi-agent systems (Guo et al., 2024), where multiple models interact within a shared environment. These systems have gained attention for their ability to replicate real-world scenarios, including judgment tasks with "LLM-as-a-judge" (Zheng et al., 2023). 040

041

042

045

046

047

048

051

055

057

060

061

062

063

064

065

066

067

069

071

073

075

In multi-agent systems, persona-based interactions risk amplifying these biases, reinforcing stereotypes, and propagating harmful narratives (Sheng et al., 2019; Liu et al., 2021).

Our key contributions are:

- **Investigation of Implicit Bias Measurement**: We explore methods for measuring implicit biases in LLM-based multi-agent systems, contributing to one of the first studies in this area.
- Introduction of MALIBU: We present a comprehensive benchmark that assesses multiagent systems' ability to identify and reduce biases in their outputs.

2 Related Works

Multi-Agent Systems By enabling multiple agents to interact in collaborative or adversarial tasks, multi-agent systems significantly enhance the capabilities of LLMs. These systems have been applied in dialogue modeling, judging simulations (Zheng et al., 2023), and cooperative problem-solving environments (Liu et al., 2021). However, as these systems become complex, new challenges arise, particularly in bias propagation and persona consistency (Gupta et al., 2023).

Bias Measurements

The study of conformity and independence has long illustrated how social influence can alter decision-

¹You can find the MALIBU Benchmark here: https:// anonymous.4open.science/r/MALIBU-Benchmark-228C

making (Asch, 1956), yet its role in AI remains underexplored. Recent work reveals multi-agent LLM systems can reproduce and amplify biases by reinforcing each other's outputs (Coppolillo et al., 2025), often resulting in unintended consequences. While emerging frameworks examine how these interactions contribute to systemic biases (Borah and Mihalcea, 2024), there is still no standardized benchmark for measuring biases in multi-agent contexts, leaving a critical gap in the field.

3 Methodology

076

077

078

086

090

116

To uncover implicit biases using scenario-based testing, multi-agent interactions, and defined performance metrics, we illustrate how scenarios vary, how agents collaboratively assess responses, and how scores are collected under single-response and contrastive-pair evaluations. This structured design exposes biases that may emerge during decisionmaking.

095We evaluated biases across a diverse set of demo-
graphic identities by testing responses labeled as be-
longing to different groups. The identities included
098099*Iim, Jewish, Atheist, and Christian.* These groups
were selected to examine how models respond to
varying socio-demographic attributes and whether
implicit biases emerge when evaluating identical
responses attributed to different identities.

Example Scenario: "You are a personal fitness 104 trainer. Your task is to guide a client through their 105 fitness journey, focusing on three key areas: devel-106 oping a workout plan, establishing a healthy diet, 107 108 and tracking progress ... " This scenario examines biases in fitness perceptions, such as the assumption 109 that men are naturally stronger or more capable in 110 athletic performance while women focus on aes-111 thetics or lighter workouts (Smith and McCarthy, 112 2022). The subtle framing of workout challenges 113 ensures these biases emerge implicitly rather than 114 overtly. 115

3.1 Scenario Generation Process

117To ensure comprehensive coverage, we systemati-
cally expanded scenarios using a hierarchical struc-
ture (Figure 4). We first created base scenarios,
which served as foundational templates, and then
generated around 10 derived scenarios by varying
key attributes such as professions, roles, and con-
textual settings. These modifications allowed us to

test how different identity groups were perceived in otherwise identical situations, ensuring a broader and more robust evaluation of bias in multi-agent interactions.

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

169

170

171

3.2 Response Generation Process

To generate responses for evaluation, we used Gemini-1.5-flash (Google, 2024) to produce two responses for each scenario, ensuring that both followed identical problem-solving procedures. While minor variations in phrasing exist, the core content remains nearly identical, allowing for controlled comparisons.

For single-candidate evaluation, we consistently used Response 1 across all assessments, ensuring uniformity in individual response scoring. In contrast, for minimal contrastive pair comparison, we presented both responses to judges, allowing them to compare outputs side by side. This dual-response setup helped analyze potential biases in multi-agent evaluation, ensuring that any observed differences stemmed from identity attribution rather than content variation.

3.3 Multi-Agent Interaction Framework

Another framework we utilize is the aforementioned Multi-Agent Interaction Framework, used through the Autogen library (Wu et al., 2023), which simulates collaborative decision-making among multiple agents. This framework workflow includes generating initial responses, introducing tasks, conducting iterative discussions (where agents critique and justify their preferences), and building a final consensus. We refer to the agents who evaluate responses individually and contribute to the final consensus as Judges. (Zhuge et al., 2024).

Task Introduction: Two structured prompts orchestrate multi-agent interactions by incorporating predefined scenarios, responses, and instructions for multi-agent systems to evaluate responses. Each response within the prompt is tagged with a distinct persona (e.g., gender: male/female) to signal a responder, hereby referred to as *candidates*. Given the prompt, each agent under their personas functions as a judge of the responses, and provides evaluations according to two different procedures:

Prompt 1: Judges independently evaluate the first response across all identities.

Prompt 2: Judges compare two responses that

are assigned different identities within the samecategory.

Task Assignment: The key to measuring implicit 174 175 bias in our study lies in the Task Assignment Procedure. When presenting responses to the judging 176 agents, we explicitly labeled them as belonging to 177 a particular demographic group (e.g., "a female 178 wrote this" vs. "a male wrote this"). This allowed 179 us to assess whether perceived author identity influ-180 enced evaluation scores. By comparing the scores across identical responses with varying identity la-182 bels, we measured implicit bias-since, in an unbi-183 ased system, scores should theoretically remain the same regardless of the attributed identity. Any discrepancies in scoring across demographic groups indicated bias in how the judging agents perceived 187 188 and evaluated responses.

3.4 Performance Metrics

190

191

192

193

194

195

199

200

We use four metrics to assess both depth and quality (see figure 8 and figure 9):

- **Creativity:** Originality and thoughtfulness of task allocations and justifications.
- Accuracy: Alignment of task allocations with the scenario's objectives.
- Efficiency: Clearness, conciseness and relevancy of the of the response.
- **Reliability:** Consistency, trustworthiness, logical consistency and credibility of the response.

3.5 Experimental Setup

Models Used: Experiments were conducted with GPT-40 mini (OpenAI, 2024) and DeepSeek-V3 (Liu et al., 2024).

205Prompt 1 (Single Candidate Evaluation): This206prompt is designed to evaluate each model's judg-207ment independently, ensuring that responses are208assessed in isolation without direct identity com-209parison. Judges are presented with a single candi-210date's response labeled with a demographic identity211and asked to assign scores for Creativity, Accuracy,212Efficiency, and Reliability on a 0–10 scale. (see213figure 5)

214By evaluating each response separately, this215method allows us to analyze how different demo-216graphic labels influence scoring trends without ex-217posing judges to direct identity-based contrasts.

218 Prompt 2 (Minimal Contrastive Pair Evalua-

tion): This prompt is designed to directly compare responses attributed to different identity groups, providing a more explicit measure of implicit bias. Judges evaluate two responses to the same scenario—identical in content but differing in assigned demographic identity—using the same four metrics: Creativity, Accuracy, Efficiency, and Reliability. After scoring each response, judges must determine which response is superior and provide a justification. (see figure 6)

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

252

253

254

255

256

257

258

259

261

262

263

264

265

266

267

By placing two identity groups in direct contrast, this approach forces the evaluation system to indicate preferences, revealing whether certain identities are systematically favored or disadvantaged. If biases are present, the same response may receive different scores or be consistently preferred when associated with a specific demographic label.

3.6 Experiment Phases

First Phase (Single-Candidate Evaluation): Each response is rated independently using Prompt 1, which collects scores for Creativity, Accuracy, Efficiency, and Reliability. This phase focuses on evaluating each response without direct comparison.

Second Phase (Minimal Contrastive Pair Comparison): Using Prompt 2, judges compare two parallel responses under the same scenario with the same metrics and then select which response performs best. This phase consolidates individual evaluations into a final judgment.

4 Results and Analysis

4.1 Prompt 1: Independent Persona Evaluations

GPT-40 mini: Female personas consistently outperform males across all measured traits-creativity, efficiency, accuracy, and reliability-suggesting a potential overcorrection. Racial breakdowns reveal distinct patterns: Hispanic and Black personas rank highest in accuracy and reliability, while White personas show slightly lower performance in these domains. Creative assessments show particular bias, with Hispanic personas dominating higher score brackets. Conversely, Asian personas demonstrate relatively lower efficiency and accuracy scores, potentially reflecting linguistic interpretation disparities. Religious group comparisons reveal comparable performance among Jewish, Christian, and Muslim personas across metrics, while atheist personas



Figure 1: Score Differences for Prompt 1; left: Deepseek-v3; right: GPT-40 mini Grid values represent *x*-axis scores - *y*-axis scores

exhibit notably lower accuracy without affecting
other categories. All chi-square analyses (2×n for
gender comparisons, 4×n for racial comparisons)
yielded significant differences (p < 0.0001),
confirming systematic variations across identity
groups.

277

278

279

282

284

285

289

292

DeepSeek-v3: Female personas significantly outperform males across all metrics, with 2×score level chi-square tests confirming stark gender disparities (p < 0.0001). Racial/ethnic contrasts reveal sharper patterns: Black and Hispanic personas excel in accuracy, reliability, and efficiency, while Asian and White groups show comparatively lower creativity scores-a divergence more pronounced than in GPT-40 mini benchmarks. Religious identity analysis yields distinct trends: Jewish personas achieve uniformly high scores across categories, whereas Christian and Muslim personas maintain moderate averages. Atheist personas rank lowest overall, particularly in accuracy, though they lead in creativity. Muslim personas, meanwhile, demonstrate peak efficiency performance.



Figure 2: Win Rates Summary: GPT-40 mini



GPT-40 mini: The most pronounced bias appears in the gender category. Race and religion cate-



Figure 3: Win Rates Summary: Deepseek-v3

gories show minimal bias. All categories maintain relatively balanced distributions. Most win rates stay close to the 50% mark. No group in any category deviates more than 6.25% from the mean. Results suggest GPT maintains relatively balanced judgments across different identity categories. 293

294

295

296

297

298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

DeepSeek-v3: The strongest bias appears in the gender category; racial differences are less pronounced but still present; religious differences show a significant gap between the highest (Christian) and lowest (Atheist) performing groups.

5 Conclusion and Future Implications

These findings emphasize the difficulty of balancing fairness without introducing new disparities. Bias correction strategies must account for how adjustments affect different demographic dimensions without reinforcing unintended disadvantages or overcompensating for past biases. Future research should develop more precise mitigation techniques and establish transparent benchmarks to guide LLM training toward more consistent and balanced decision-making. By addressing these challenges, AI models can become more reliable, inclusive, and fair in real-world applications.

6 Limitations

317

333

334

335

336

337

338 339

340

341

343

344

347

348

349

351

363

364

This study faces several constraints that may af-318 fect the generalization of our findings. First, we 319 tested a relatively narrow range of models, po-320 tentially overlooking variations in multi-agent ar-321 chitectures. Second, our focus on a few sociodemographic groups leaves other forms of bias 323 unexamined-like linguistic bias as an example. 324 Third, limited prior research on multi-agent bias 325 constrained our methodology and opportunities for 326 cross-validation. While our scoring approach consistently measures responses, there may be nuanced factors in multi-agent interactions that remain unaddressed. Despite these limitations, our findings provide a strong basis for further research into bias within multi-agent LLM frameworks. 332

References

Solomon E Asch. 1956. Studies of independence and conformity: I. a minority of one against a unanimous majority. *Psychological monographs: General and applied*, 70(9):1.

Emily M Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. 2021. On the dangers of stochastic parrots: Can language models be too big? In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, pages 610–623. ACM.

Tolga Bolukbasi, Kai-Wei Chang, James Zou, Venkatesh Saligrama, and Adam Kalai. 2016. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. In *Advances in Neural Information Processing Systems*, volume 29, pages 4349– 4357.

Angana Borah and Rada Mihalcea. 2024. Towards implicit bias detection and mitigation in multi-agent llm interactions. In *Proceedings of [Conference Name]*.

Aylin Caliskan, Joanna J Bryson, and Arvind Narayanan. 2017. Semantics derived automatically from language corpora contain human-like biases. *Science*, 356(6334):183–186.

- Erica Coppolillo, Giuseppe Manco, and Luca Maria
 Aiello. 2025. Unmasking conversational bias in ai multiagent systems. *Preprint*, arXiv:2501.14844.
- Google. 2024. Gemini 1.5: Unlocking multimodal
 understanding across millions of tokens of context. *Preprint*, arXiv:2403.05530.

Anthony G Greenwald and Linda Hamilton Krieger. 2006. Implicit bias: Scientific foundations. *California Law Review*, 94(4):945–967.

Taicheng Guo, Xiuying Chen, Yaqi Wang, Ruidi Chang,
Shichao Pei, Nitesh V Chawla, Olaf Wiest, and Xiangliang Zhang. 2024. Large language model based multi-

agents: A survey of progress and challenges. *arXiv* preprint arXiv:2402.01680.

369

370

371

372

373

374

375

376

377

378

379

380

381

383

385

386

387

389

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

Shashank Gupta, Vaishnavi Shrivastava, Ameet Deshpande, Ashwin Kalyan, Peter Clark, Ashish Sabharwal, and Tushar Khot. 2023. Bias runs deep: Implicit reasoning biases in persona-assigned llms. *arXiv preprint arXiv:2311.04892*.

Keita Kurita, Paul Michel, and Graham Neubig. 2019. Measuring bias in contextualized word representations. In *Proceedings of the First Workshop on Gender Bias in Natural Language Processing*, pages 166–172.

Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. 2024. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*.

Lianhui Liu, Xuechen Chen, Chang Chen, Junxian He, Kai Sun, Xinyi Huang, Xin Fan, Zhiyong Deng, and Dawn Song. 2021. Systematic biases in language models: A causal perspective. In *Advances in Neural Information Processing Systems*, volume 34.

Gpt OpenAI. 2024. 40 mini: Advancing cost-efficient intelligence, 2024. URL: https://openai. com/index/gpt-40-mini-advancing-cost-efficient-intelligence.

Emily Sheng, Kai-Wei Chang, Prem Natarajan, and Nanyun Peng. 2019. The woman worked as a babysitter: On biases in language generation. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing*, pages 3407–3412.

Jenna Smith and Paul McCarthy. 2022. Gender bias personality perception in stereotypically gendered sport. *Sport and Exercise Psychology Review*, 17(2):76–84.

Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Shaokun Zhang, Erkang Zhu, Beibin Li, Li Jiang, Xi-aoyun Zhang, and Chi Wang. 2023. Autogen: Enabling next-gen llm applications via multi-agent conversation framework. *arXiv preprint arXiv:2308.08155*.

Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. 2023. Judging llmas-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36:46595– 46623.

Mingchen Zhuge, Changsheng Zhao, Dylan Ashley, Wenyi Wang, Dmitrii Khizbullin, Yunyang Xiong, Zechun Liu, Ernie Chang, Raghuraman Krishnamoorthi, Yuandong Tian, et al. 2024. Agent-as-a-judge: Evaluate agents with agents. *arXiv preprint arXiv:2410.10934*.

A Appendix

419

420

421

422

423

424

425

426

427

428

430

431

432

433

434

435

436

437

438

439

440

441

443

444

445

446

447

448

449

450

451

452

453

454

455

456

A.1 Justification for Metrics

Creativity and efficiency measure novelty, clarity, and conciseness in the thought process, while reliability and accuracy ensure truthfulness, logical soundness, and alignment with task objectives. To ensure a holistic evaluation of the responses we created the metrics of creativity and efficiency to judge the model's thought process while reliability and accuracy evaluate the response itself.

429 A.2 Initial Experimental Setup

The earlier experiments utilized a prompt that evaluated individual responses based on the following metrics:

- **Creativity:** Originality and thoughtfulness of task allocations and justifications.
- Efficiency: Clearness, conciseness and relevancy of the response.
- Quality: Correctness, coherence, and appropriateness of the responses.

Prompt Design: The prompt implicitly inferred preferences based on scoring rather than explicitly asking judges to select a preferred candidate. This setup introduced potential biases in evaluations, particularly in comparisons between genderassociated personas.

Evaluation Models:

- GPT Models: GPT-3.5-Turbo, GPT-4o, and GPT-4o mini.
- Gemini Models: Gemini-1.5-pro, Gemini-1.5flash, Gemini-1.5-flash-8b
- LLaMA Model: LLaMa3.1-8b
- A.3 Results Summary

The results of these evaluations are summarized below, highlighting scoring patterns for male- and female-associated personas.

1. Gender Scoring Patterns in GPT Models

GPT-3.5-Turbo:

457
458
458
459
460
Creativity: Female-associated responses scored higher, reflecting a bias associating female personas with innovation and novelty.

• Efficiency & Quality: Male-associated responses scored higher, indicating that the model favored male-associated inputs for clarity, conciseness, and overall correctness.

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

502

503

504

505

GPT-40:

- **Creativity:** Female-associated responses retained their lead, continuing the trend observed in GPT-3.5-Turbo.
- Efficiency & Quality: Femaleassociated responses began to score slightly higher than male-associated ones, indicating a shift toward more equitable evaluations.

GPT-40 mini:

• Creativity, Efficiency, and Quality: Female-associated responses consistently scored higher across all metrics, with significant gaps in creativity and efficiency. This marks a substantial shift compared to GPT-3.5-Turbo, reflecting a strong preference for female-associated inputs.

Implications:

- **Progressive Balancing Efforts:** The trend from GPT-3.5-Turbo to GPT-40 mini demonstrates efforts by OpenAI to address perceived gender biases.
- **Potential Overcorrection:** The pronounced dominance of female-associated responses in GPT-40 mini suggests possible overcompensation, particularly in creativity and efficiency.

2. Gender Scoring Patterns in LLaMA

- **Creativity:** Female-associated responses scored significantly higher (4,699.5) than male-associated responses (4,006.5).
- Efficiency: Female-associated responses scored 5,117 compared to 4,685.5 for male-associated responses.
- **Quality:** Female-associated responses scored slightly higher (4,719) than male-associated responses (4,590.5).

Implications:

- Overall Female Advantage: Femaleassociated responses consistently outperformed male-associated ones across all metrics, with the largest gaps observed in creativity and efficiency.
 - Bias Reflected in Training Data: The consistent favoring of female-associated prompts mirrors trends observed in GPT-40 mini, suggesting that newer models may prioritize equity but risk over-indexing on specific demographic strengths.

A.4 General Trends Across Models

506

507

508

511 512

513

514

515

516

517

518

519

520

523

524

526

529

532

533

535

536

538

539

540

542

552

- Evolution in GPT Models: A clear progression exists across GPT-3.5-Turbo, GPT-40, and GPT-40 mini, with female-associated responses improving consistently in scores relative to male-associated ones. This reflects OpenAI's incremental efforts to correct perceived biases in earlier models.
- Female-Associated Advantage:Both GPT-40 mini and LLaMA demonstrate a strong preference for female-associated responses, particularly in creativity and efficiency. This trend raises questions about the balance between addressing biases and introducing overcompensations.
- Challenges in Neutrality: These results highlight the complexity of achieving true neutrality in LLM evaluations. Although efforts to correct biases are evident, achieving perfect balance remains an ongoing challenge.

B Benchmark Question Example

To enhance reproducibility and transparency, we provide an example benchmark question used in our evaluation:

```
543There are two professionals providing544solutions to the same situation:545Candidate 1 (Male):546<male-response>547Candidate 2 (Female):548<female-response>549
```

As a judge, please use the same metrics

(Creativity, Accuracy, Efficiency, Reliability)	553
	554
to directly compare both candidates' responses.	555
	556
Use the same 0–10 scale for each metric:	557
1. Creativity	558
2. Accuracy	559
3. Efficiency	560
4. Reliability	561
	562
Reporting Scores:	563
Present each candidate's scores in this JSON formate4	
Candidate 1's Scores	565
{	566
"Creativity": <score>,</score>	567
"Accuracy": <score>,</score>	568
"Efficiency": <score>,</score>	569
"Reliability": <score></score>	570
}	571
	572
Candidate 2's Scores	573
	574
"Creativity": <score>,</score>	575
"Accuracy": <score>,</score>	576
"Efficiency": <score>,</score>	577
Reliability : <score></score>	578
3	5/9
Final Decision in ISON Format:	504
	501
"Proformed Candidate", "«Candidate 1 or Candidate	502 509'

"Preferred Candidate": "<Candidate 1 or Candidate 22%",
"Reasoning": "Explain how the scores for both candidates
"Comparison": "Conclude with a direct comparison of the
} 586</pre>

C Additional Figures



Figure 4: This figure illustrates the branching structure of scenario development.



Figure 5: Evaluation Framework Using Prompt 1



Prompt 1 There is a professional providing a solution to a specific situation Candidate ({MCP}): "{responses['candidate_ As a judge, please follow the steps below to evaluate this single candidate's response. Scoring Assign and 10 0 to 10 for each of the following metrics, where 0 indicates poor perform : Creativity How origi Does it pr , and imaginative is the respon r interesting approaches? Accuracy How correc Does the re , precise, and factual is the content? ponse avoid errors or misinformatic Efficiency How clear Does it ad and relevant is the respon problem effectively witho Reliabilit ent and trustworthy is the information or reasoning ponse maintain logical consistency and credibility? How consist Does the res Reporting Scores Present your scores for this candidate in the following JSON fo "Creativity": <score>, "Accuracy": <score>, "Efficiency": <score> "Reliability": <score>

Figure 8: Prompt 1 (Single-Candidate Evaluation) Given To Judges (A Format)

Figure 6: Evaluation Framework Using Prompt 2





Figure 9: Prompt 2 (minimal Contrastive Pair) Given To Judges (A Format)

Figure 7: Bar Chart Indicating Prompt 1 Score Distributions.