Can LLMs Produce Faithful Explanations For Fact-checking? Towards Faithful Explainable Fact-Checking via Multi-Agent Debate

Anonymous ACL submission

Abstract

Fact-checking research has extensively explored verification but less so the generation of natural-language explanations, crucial for user trust. While Large Language Models (LLMs) excel in text generation, their capability for producing faithful explanations in fact-checking remains underexamined. Our study investigates 800 LLMs' ability to generate such explanations, finding that zero-shot prompts often result in unfaithfulness. To address these challenges, we propose the Multi-Agent Debate Refinement (MADR) framework, leveraging multiple LLMs as agents with diverse roles in an 014 iterative refining process aimed at enhancing faithfulness in generated explanations. MADR ensures that the final explanation undergoes 017 rigorous validation, significantly reducing the likelihood of unfaithful elements and aligning closely with the provided evidence. Experimental results demonstrate that MADR significantly improves the faithfulness of LLM-generated explanations to the evidence, advancing the credibility and trustworthiness of these explanations.

1 Introduction

027

In the digital age, swiftly spreading misinformation necessitates not only the verification of claims but also the provision of clear explanations for these verifications. Such explanations are crucial for building trust within the audience, as lack of them often leads to distrust in fact-checking results (Guo et al., 2022). Moreover, explanation generation becomes even more critical in multi-hop factchecking, where complex reasoning across multiple evidence pieces is required to assess a claim's veracity (Reddy et al., 2023).

Despite the adeptness of Large Language Models (LLMs) in generating diverse texts, their capacity for crafting *faithful*¹ explanations for factchecking remains underexplored. Faithfulness is crucial; explanations that misrepresent evidence could exacerbate misinformation, posing a significant challenge. Thus, enhancing the faithfulness of generated explanations in fact-checking is an urgent, unresolved issue.

041

043

044

045

048

051

054

055

057

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

076

077

078

Our first research question asks: **can LLMs generate faithful explanations for fact-checking in a zero-shot prompting setup?** To facilitate analysis, we define a novel typology of common errors and unfaithfulness issues that arise in LLM-generated explanations. We conduct extensive experiments prompting ChatGPT (OpenAI, 2023c) to explain fact checks from multiple sources. Our findings reveal that **zero-shot prompting LLMs often fails to yield faithful explanations**. 80% of the generated explanations include hallucinated details (§5).

This leads to our second research question: how to address the unfaithfulness issues in LLMgenerated explanations? We propose the Multi-Agent Debate Refinement (MADR) framework that uses multiple LLMs as agents to provide feedback for iterative refinement to produce faithful explanations (§3). The goal is to mitigate unfaithfulness and steer the LLM-generated texts towards true rationales. Experimental results show that MADR significantly improves faithfulness upon baselines, demonstrating the effectiveness of our approach.

Our contributions are summarized as follows:

- We present the first study of LLMs' ability to produce faithful fact-checking explanations.
- We present Multi-Agent Debate Refinement, an effective framework to produce faithful explanations based on iterative debating among LLMs.
- Our correlation analysis reveals the most suitable LLM-based evaluation protocol for this task.

2 Typology

In our analysis of explanations generated by LLMs, we have introduced a novel typology encompassing a range of error categories, as shown in Ta-

¹Faithfulness refers to the factual consistency between the explanation and the given evidence (Huang et al., 2023a).

Example

Claim: Says Jeff Foxworthy wrote a list of examples explaining how "you might live in a nation that was founded by geniuses but is run by idiots."

Evidence: ... Foxworthy is famous for his "You might be a redneck if" jokes , but searching online, we couldn't find any credible evidence that he penned this list that touches on abortion, Muslims and federal debt. Rather, we found users on web forums crediting someone named Fritz Edmunds with the list. Snopes, which fact-checked this claim back in 2013, also noted that "the original compiler of this appears to be Fritz Edmunds, who posted it to his Politically True blog back in Feb. 3, 2013 ...

Error Type & Explanation

Intrinsic Entity Error: The generated explanation misrepresents named entities, quantities, dates, or other surface realizations from the given source. E.g. *Fritz Foxworthy was credited on a web forum with the list.*

Extrinsic Entity Error: The generated explanation includes new entities that are not present in the given source. E.g. *Elon Musk was credited on a web forum with the list.*

Intrinsic Event Error: The generated explanation misrepresents events mentioned in the source. E.g. *They couldn't find* any credible evidence that Fritz Edmunds was credited on a web forum.

Extrinsic Event Error: The generated explanation include new events that are not present in the given source. E.g. *Foxworthy found that Fritz Edmunds made the "You might be a redneck if" jokes.*

Intrinsic Noun-Phrase Error: The explanation mistakenly represents the noun phrases in the given source like miscombining modifiers combined with one entity to another entity. E.g. *They found the original user on web forums crediting someone named Fritz Edmunds*.

Extrinsic Noun-Phrase Error: The explanation mistakenly represents new noun phrases that are not present in the given source like miscombining modifiers not presented in the source to entity. E.g. *They found a mysterious user on web forums crediting someone named Fritz Edmunds.*

Reasoning Coherence Error: There are logical flaws in the flow of reasoning within the generated explanation, leading to a lack of coherence or weak support for the claim. E.g. *While they were searching online, they couldn't find any credible evidence that he penned this list that touches on abortion.*

Overgeneralization Error: The generated explanation makes sweeping statements or draws conclusions that go beyond the evidence provided. E.g. *Fritz Emunds is the one who spreaded the rumor and put the blame on Foxworthy.*

Irrelevant Evidence Error: The generated explanation includes evidence that is not directly related to the claim, leading to confusion and lack of support for the main argument. E.g. ... Foxworthy is famous for his "You might be a redneck if" jokes.

Table 1: An illustration of error typology using using a sample data from PolitiHop (Ostrowski et al., 2021). The errors in the sample summaries are in red color and italicized.

ble 1. The classification of intrinsic and extrinsic errors within the domains of Entity-Related, Event-Related, and Noun-phrase Errors draws inspiration from relevant studies in other domains (Goyal and Durrett, 2021; Huang et al., 2023c). We have incorporated additional context-specific error types, enriching the overall typology.

3 Methodology

090

100

Zero-shot prompting LLMs often produce unfaithful explanations which contain multiple errors. In the early stage of our experiments, we incorporate an iterative refinement paradigm for improving their faithfulness. However, we found that selfrefinement (Madaan et al., 2023) alone was insufficient for faithfulness enhancement (see Table 3), as imprecise feedback tended to guide misguided refinements of the explanation. This underscores the pivotal role of precise feedback for efficient refinement by the LLM (Wang et al., 2023).

Thus, we propose **Multi-Agent Debate Refinement (MADR)**, inspired by a debate-based methodology (Du et al., 2023). While Du et al. (2023) focus on refining explanations during a debate, MADR utilizes the debate for generating feedback to be employed in subsequent refinement stages. Our method offers several advantages. First, compared to directly refining explanations during a debate, MADR facilitates a dynamic and iterative feedback loop, enhancing the identification of errors. Secondly, it ensures more accurate feedback, reducing the likelihood of misguided refinements and ultimately enhancing overall faithfulness. Thirdly, this approach prompts bidirectional thinking within the LLM, enabling it to analyze explanations both with and without knowledge of predefined error types, fostering an explicit rationale.

101

102

105

106

110

111

112

113

114

115

116

117

118

119

120

121

122

The process of MADR is outlined in Algorithm 1 and depicted in Figure 1. MADR employs multiple agents to identify errors and engage them in a debate until a consensus is reached on the debate. Four total roles are assigned to each agent: two serve as DEBATERS, one as a JUDGE, and one as a REFINER.An initial explanation *E* is



Figure 1: An overview of MADR.

generated through zero-shot prompting an LLM. 123 DEBATER $1(D_1)$ and DEBATER 2 (D_2) pinpoint 124 errors in E and propose feedback $F_{i,n}$ to amend 125 126 these issues, where *i* represents the debate iteration and $n \in \{1, 2\}$ identifies each DEBATER (lines 2-127 3). Distinctive instructions with varying goals are provided to the DEBATERS: DEBATER 1 identifies 129 errors based on a predefined error typology (see 130 §2), while DEBATER 2 focuses on potential errors 131 that may affect the explanation's faithfulness, with-132 out relying on the error typology (refer to Table 8 and Table 9 for prompt specifics). This setup ensures that MADR promotes the identification of 135 errors that might be overlooked by either party. 136

Next, in the *i*-th iteration, DEBATERS D_1 and D_2 review the feedback given by each other in the previous (i - 1)-th iteration (i.e. $F_{i-1,2}$ for D_1 and $F_{i-1,1}$ for D_2). They refine their feedback by adding any missed elements and removing errors (lines 9-10).

137

138

139

140

141

142

To ensure the most accurate feedback, the two DEBATERS continue their discussion until they reach a mutual agreement on the feedback. During the *i*-th iteration, the JUDGE agent assesses the feedback from D_1 and D_2 and determines whether the DEBATERS reach a mutual agreement on the feedback. When $J(F_{i,1}, F_{i,2}) = True$, the debate stops (lines 6-7). Finally, we concatenate the final feedback from both DEBATERS and feed it to the refiner to refine its explanation using the concatenated feedback (line 13). An example of the outputs from MADR is shown in Table 11. Additionally, to prevent endless debates, we set a fixed number of iterations. 143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

163

164

165

166

167

168

169

170

171

172

173

174

175

177

178

179

181

182

183

184

185

186

187

189

190

191

4 Experimental Settings

Dataset and Metric Experiments are conducted on the PolitiHop multi-hop fact-checking dataset (Ostrowski et al., 2021). PolitiHop consists of 445 test set instances, where each instance contains a claim and multiple pieces of evidence. The veracity of a claim can only be determined by reasoning across multiple pieces of evidence and the claim. For the evaluation metric, we use G-Eval (Liu et al., 2023) with GPT-4 Turbo (OpenAI, 2023b) to assess whether the generated explanation is consistent with the evidence. Following Huang et al. (2023b), we adopted 4 evaluation protocols based on G-Eval which vary in granularity, ranging from sentence-level to document-level assessments, and in the application of our error typology, validating the effectiveness of the error typology in assisting automatic evaluation. The prompt templates are shown in Appendix C.

Baselines We compare MADR with the following competitive methods. **Zero-shot** prompts an LLM to directly output an explanation given the input claim and evidence. **CoT** asks LLMs to generate the reasoning process before producing the final output. **Self-Refine** (Madaan et al., 2023) generates an initial explanation and then iteratively refines the explanation with one agent. We conduct experiments by using GPT-3.5-Turbo (OpenAI, 2023a) to generate explanations across all experiments for fair comparisons. The prompts for these approaches are displayed in Appendix C. The case study using **Self-Refine** is in Table 10.

5 Results

The main results are summarized in Table 2. MADR achieves the best faithfulness scores on

Granularity→	Sentence-level		Document-level	
Typology Applied→ Method↓	No	Yes	No	Yes
Zero-shot	4.87	4.84	4.70	4.92
CoT	4.86	4.91	4.76	4.96
Self-Refine	4.70	4.86	4.89	4.81
MADR (ours)	4.82	4.99	4.88	4.97

Table 2: Faithfulness evaluation on PolitiHop test set. Scores are computed using G-Eval with evaluation protocols of varying granularity and application of our error typology. The best score per column is bolded.

Method	Faithful Explanations (%)	# Errors
Zero-shot	20.0	25
CoT	5.0	42
Self-Refine	20.0	32
MADR (ours)	30.0	17

Table 3: Human evaluation results on 20 samples from PolitiHop. MADR produces the most faithful explanation compared to baselines.

193

194

195

196

197

198

200

201

204

205

210

211

212

213

214

215

216

217

218

221

222

two out of the four evaluation protocols, indicating its effectiveness in producing faithful explanations. To further validate the effectiveness of our method, we conduct human evaluations via Amazon Mechanical Turk, aiming to quantify the portion of faithful explanations and the total error count. Annotators were presented with our error typology and were tasked to identify the presence 199 of each error type within individual sentences. The results of human evaluations are shown in Table 3. We have the following observations. First, using simple prompting methods, such as zero-shot or CoT, LLMs often produce unfaithful explanations for fact-checking. This highlights the challenge of generating faithful explanations for LLMs in complex fact-checking scenarios, such as PolitiHop, which requires reasoning through multiple pieces of evidence. Second, despite the high faithfulness scores suggested by automatic evaluation (approaching the maximum score of 5) in Table 2, human evaluators frequently deemed the LLM-generated explanations unfaithful, as per Table 3. This discrepancy suggests that even the most advanced LLM, GPT-4 Turbo, fails to reliably judge the faithfulness of generated explanations for fact-checking.

To pinpoint the most effective LLM-based evaluation strategy for future research, we performed a correlation analysis, correlating human evaluations with automatic metrics using Kendall's Tau (variant c). According to Table 4, a granular evaluation aligns better with human judgments, and incorporating our error typology into

Evaluation Protocol	Kendall's Tau
Document-level w/o Typology	0.008
Document-level w/ Typology	0.128
Sentence-level w/o Typology	0.105
Sentence-level w/ Typology	0.150

Table 4: Correlation between evaluation protocols and human judgments on the PolitiHop dataset.

automatic evaluations enhances the quality of LLM assessments. Details on human evaluation methodology are provided in Appendix A.

225

226

227

229

231

232

233

234

235

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

258

259

260

262

263

264

265

Furthermore, a case study showcased in Table 11 highlights the superiority of MADR over selfrefinement, by demonstrating that MADR allows DEBATERS to identify and correct errors missed during self-refinement, leading to more accurate explanations. In contrast, the Self-Refine approach, as shown in Table 10, fails to produce a faithful explanation, emphasizing the advantage of employing multiple perspectives with MADR for error identification and explanation validation.

6 **Related Work**

Early approaches to producing explanations for fact-checks can largely be categorized into logic-based methods (Gad-Elrab et al., 2019; Ahmadi et al., 2019) and attention-based methods (Shu et al., 2019; Lu and Li, 2020). Recent work generates natural-language explanations using abstractive (Kotonya and Toni, 2020) or extractive (Atanasova et al., 2020) approaches. A very recent study benchmarks the ability of these models to generate explanations for fact-checks (Russo et al., 2023). Our study complements this work by presenting the first empirical analysis on LLMs' ability to generate fact-checking explanations and propose a method to enhance its faithfulness.

7 Conclusion

This paper empirically demonstrates that LLMs often produce unfaithful explanations for fact-checks. We introduce the Multi-Agent Debate Refinement (MADR) framework, which utilizes multiple LLM agents to iteratively debate and refine explanations, significantly enhancing their faithfulness as evidenced by both automatic and human evaluations. Our results underscore the efficacy of multi-agent debate in mitigating LLMs' unfaithfulness. Additionally, we reveal that LLMs cannot reliably assess the faithfulness of the generated explanations and discover the most suitable evaluation protocols for LLM-based automatic evaluation.

4

273

274

294 295

296

298

299

304

307

310

311

312

314

315

8 Ethical Considerations

LLMs trained on internet data often show biases, but this focus mainly applies to data and models reflecting the culture of English-speaking communities. However, detailed reviews of model outputs for the PolitiHop dataset have found no signs of biases concerning gender, age, race, or other socioeconomic elements.

9 Limitations

Our study did not thoroughly investigate the sensitivity of various systems to changes in input prompts. It is recognized that the effectiveness of numerous natural language processing tasks can 279 significantly depend on how input prompts are designed. By not conducting a comprehensive analysis on prompt sensitivity, we acknowledge the 283 possibility that different prompts might elicit a wide range of responses that we have not explored, potentially limiting the applicability of our findings. However, it is important to note that we did not engage in prompt tuning specifically to favor 287 our proposed framework, ensuring that the comparisons between different techniques remain equitable. Given the scope of our research, the detailed 291 exploration of prompt sensitivity is an area we have designated for future investigation.

References

- Naser Ahmadi, Joohyung Lee, Paolo Papotti, and Mohammed Saeed. 2019. Explainable fact checking with probabilistic answer set programming. *arXiv preprint arXiv:1906.09198*.
- Pepa Atanasova, Jakob Grue Simonsen, Christina Lioma, and Isabelle Augenstein. 2020. Generating fact checking explanations. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7352–7364, Online. Association for Computational Linguistics.
- Jacob Cohen. 1960. A coefficient of agreement for nominal scales. *Educational and psychological measurement*, 20(1):37–46.
 - Yilun Du, Shuang Li, Antonio Torralba, Joshua B Tenenbaum, and Igor Mordatch. 2023. Improving factuality and reasoning in language models through multiagent debate. *arXiv preprint arXiv:2305.14325*.
 - Mohamed H. Gad-Elrab, Daria Stepanova, Jacopo Urbani, and Gerhard Weikum. 2019. Exfakt: A framework for explaining facts over knowledge graphs and text. In *Proceedings of the Twelfth ACM International Conference on Web Search and Data Mining*,

WSDM '19, page 87–95, New York, NY, USA. As-	
sociation for Computing Machinery.	

316 317

318

319

320

321

322

323

324

325

326

327

328

329

330

331

332

333

334

335

336

337

338

340

341

342

343

344

346

347

349

350

351

352

353

354

355

356

357

358

359

360

361

362

363

364

365

- Tanya Goyal and Greg Durrett. 2021. Annotating and modeling fine-grained factuality in summarization.
- Zhijiang Guo, Michael Schlichtkrull, and Andreas Vlachos. 2022. A survey on automated fact-checking.
- Kung-Hsiang Huang, Hou Pong Chan, and Heng Ji. 2023a. Zero-shot faithful factual error correction. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5660–5676, Toronto, Canada. Association for Computational Linguistics.
- Kung-Hsiang Huang, Philippe Laban, Alexander R Fabbri, Prafulla Kumar Choubey, Shafiq Joty, Caiming Xiong, and Chien-Sheng Wu. 2023b. Embrace divergence for richer insights: A multi-document summarization benchmark and a case study on summarizing diverse information from news articles. *arXiv preprint arXiv:2309.09369*.
- Kung-Hsiang Huang, Mingyang Zhou, Hou Pong Chan, Yi R Fung, Zhenhailong Wang, Lingyu Zhang, Shih-Fu Chang, and Heng Ji. 2023c. Do lvlms understand charts? analyzing and correcting factual errors in chart captioning. *arXiv preprint arXiv:2312.10160*.
- Neema Kotonya and Francesca Toni. 2020. Explainable automated fact-checking for public health claims. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 7740–7754, Online. Association for Computational Linguistics.
- Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. 2023. G-eval: NLG evaluation using gpt-4 with better human alignment. In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, pages 2511–2522, Singapore. Association for Computational Linguistics.
- Yi-Ju Lu and Cheng-Te Li. 2020. GCAN: Graph-aware co-attention networks for explainable fake news detection on social media. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 505–514, Online. Association for Computational Linguistics.
- Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Sean Welleck, Bodhisattwa Prasad Majumder, Shashank Gupta, Amir Yazdanbakhsh, and Peter Clark. 2023. Self-refine: Iterative refinement with self-feedback.

OpenAI. 2023a.	Chatgpt.	366
OpenAI. 2023b.	Gpt-4 turbo.	367
OpenAI. 2023c.	Gpt-4v(ision) system card.	368

Wojciech Ostrowski, Arnav Arora, Pepa Atanasova, and Isabelle Augenstein. 2021. Multi-hop fact checking of political claims. In *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence*.

369

370

371

372

373

374

375

376

377

379

381

382

384

386

387

388

389 390

391

- Revanth Gangi Reddy, Yi R Fung, Qi Zeng, Manling Li, Ziqi Wang, Paul Sullivan, et al. 2023. Smartbook: Aiassisted situation report generation. *arXiv preprint arXiv:2303.14337*.
- Daniel Russo, Serra Sinem Tekiroğlu, and Marco Guerini. 2023. Benchmarking the generation of fact checking explanations. *Transactions of the Association for Computational Linguistics*, 11:1250–1264.
 - Kai Shu, Limeng Cui, Suhang Wang, Dongwon Lee, and Huan Liu. 2019. defend: Explainable fake news detection. In *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining*, pages 395–405.
- Zhenhailong Wang, Shaoguang Mao, Wenshan Wu, Tao Ge, Furu Wei, and Heng Ji. 2023. Unleashing cognitive synergy in large language models: A task-solving agent through multi-persona selfcollaboration. *arXiv preprint arXiv:2307.05300*, 1(2):3.

467

468

430

431

A Human Evaluation Details

A.1 Evaluation Guidelines

In this task you will evaluate the faithulness of automatically generated fact-checking explanation using a label, claim, and source used to generate the explanation. To correctly solve this task, follow these steps:

Carefully read and understand the topology of errors and examples given below. Carefully read the generated fact-checking explanation and the source. For each explanation, check it with the evidence and decide if any of the error exists in the explanation. Note: You will analyze each sentence, but you should consider the connection between other sentences as well. Warning: Annotations will be checked for quality against control labels, low quality work will be rejected.

Type of Errors:

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

494

425

426

427

428

429

- Intrinsic Entity-Related Errors: Intrinsic entity-related errors occur when there is a mistake in representing named entities, quantities, dates, or other surface realizations from the given source within the generated explanation. Example: Incorrectly combining distinct entities from the given source.
- Extrinsic Entity-Related Errors: Extrinsic entity-related errors involve the introduction of new entities that are not present in the given source into the generated explanation. Example: Hallucinating new entities that do not exist in the source.
- Intrinsic Event-Related Errors: Intrinsic event-related errors pertain to mistakes in representing events mentioned in the generated explanation, leading to incorrect claims about events. Example: Making inaccurate claims about events mentioned in the explanation.

event-related errors occur when the generated explanation includes new events that are not present in the given source. Example: Introducing fabricated events that are not supported by the source.

- Intrinsic Noun Phrase-Related Errors: Intrinsic noun phrase-related errors are mistakes related to noun phrases, excluding entityspecific errors. They may involve miscombining noun phrases with incorrect modifiers from the given source. Example: Incorrectly combining a noun phrase with the wrong modifier from the source.
- Extrinsic Noun Phrase-Related Errors: Extrinsic noun phrase-related errors involve the introduction of new noun phrase modifiers that are not present in the given source into the generated explanation. Example: Hallucinating new noun phrase modifiers not supported by the source.
- Reasoning Coherence Errors: Reasoning coherence errors occur when there are logical flaws in the flow of reasoning within the generated explanation, leading to a lack of coherence or weak support for the claim. Example: Presenting evidence that does not logically connect to the main claim, resulting in a disjointed explanation.
- Overgeneralization Errors: Overgeneralization errors happen when the generated explanation makes sweeping statements or draws conclusions that go beyond the scope of the evidence provided.
- Irrelevant Evidence Errors: Irrelevant evidence errors occur when the generated explanation includes evidence that is not directly related to the claim, leading to confusion and lack of support for the main argument. Example: Including evidence that is tangential or

Extrinsic Event-Related Errors: Extrinsic

xplanation 1	Errors
False. The claim that the pens Nancy Pelosi used to sign the articles of impeachment cost *\$2,025 per pen* is false. According	Additionally, a search of Google results did not yield any othe mentions of the pens' price being \$2,025.
under \$15. Additionally, a search of Google results did not yield	None
any other mentions of the pens' price being \$2,025. The claim originated from a Facebook post that was flagged for spreading	Intrinsic Entity
aise information. Therefore, the claim is faise.	Extrinsic Entity
	Intrinsic Event
	Extrinsic Event
	Intrinsic Noun-Phrase
	Extrinsic Noun-Phrase
	Bessening Coherence

Figure 2: The interface for our human evaluation.

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

501

504

506

unrelated to the claim being explained.

A.2 Evaluation Interface

We display our evaluation interface in Figure 2.

A.3 Worker Qualification

We established specific initial criteria for selecting highly efficient MTurk workers. These prerequisites include having a HIT approval rate of at least 99%, completing a minimum of 10,000 approved HITs, and being located in the United Kingdom, Canada, or the United States.

Furthermore, beyond these initial requirements, qualified workers must pass two rounds of qualification tests aimed at identifying errors in generated explanations. To refine the qualification process, we manually annotated two HITs, each featuring one multi-hop fact-checking instance from Politi-Hop and an explanation generated by one of the models. In each qualification phase, annotators review one of these annotated examples. Those whose annotations do not closely match ours are excluded from the selection process.

Ultimately, 4 annotators who successfully completed all two stages of the qualification tests were selected. Additionally, we carefully designed each HIT to ensure that annotators could earn an hourly rate of \$15 to \$20, provided they work continuously.

A.4 Annotation Quality

We computed the agreement between each annotator with one of the authors of this paper. The agreement is 0.69 per Cohen's Kappa (Cohen, 1960), indicating a moderate-to-high level of agreement.

B **MADR Details**

С Prompts

Evaluation prompts are shown in Table 6 and Table 5. The prompts for self-refinement and MULTI-AGENT DEBATE REFINEMENT (MADR) are displayed in table 7 and table 11, respectively.

Algorithm 1 MADR

Input: Given claim C, given Evidence Source S, given veracity label L, generated explanation E**Output:** Refined Explanation E

1: Initialize first agent D_1 and second agent D_2 for DE-BATERS with bidirectional thinking process, and third agent J for JUDGE to judge whether DEBATERS have same feedback

//Initialize the first feedback from two agents $F_{i,1}$,

2:
$$F_{0,1}^{I'i,2} \leftarrow D_1(C,S,E)$$

3: $F_{0,2} \leftarrow D_2(C, S, E)$

4: //Set maximum number of iterates of debate to N

- 5: for i = 1 to N do
- if $J(F_{i-1,1}, F_{i-1,2}) = True$ then 6: break
- 7: 8: end if
- 9:
- $F_{i,1} \leftarrow D_1(F_{i-1,1}, F_{i-1,2}) \\ F_{i,2} \leftarrow D_2(F_{i-1,2}, F_{i-1,1})$ 10:
- 11: end for
- 12: //Initialize the agent R as the REFINER
- 13: $E^* \leftarrow R(F_{\text{final},1} + F_{\text{final},2})$
- 14: return E^{\dagger}

You will be given a fact-checking explanation along with the evidence used for fact-checking.

Table 5: Prompt templates for document-level automatic evaluation. The texts in grey are only presented in the prompts when error typology is applied.

Your task is to rate the explanation on one metric. Please make sure you read and understand these instructions carefully. Please keep this document open while reviewing, and refer to it as needed. **Evaluation Criteria:** Faithfulness (1-5) - the factual alignment between the fact-checking explanation and the evidence. The explanation should accurately reflect the evidence and its context, without misrepresenting or omitting crucial details. Annotators were instructed to penalize explanations that contain inaccuracies, misinterpretations, or fail to adequately represent the evidence provided. Below are the error typology that you need to utilize to determine faithfulness between the explanation and evidence: {error typology} Evaluation Steps: 1. Read the fact-checking explanation and the evidence provided carefully. 2. Compare the explanation to the evidence to identify how well it represents the facts, context, and conclusions drawn from the evidence. 3. Assess how accurately and completely the explanation reflects the evidence without distortion or significant omission. Assign a faithfulness score from 1 to 5. **Evidence Provided:** {evidence} Fact-Checking Explanation: {explanation} **Evaluation Form:** Faithfulness:

You will be given a sentence from a fact-checking explanation along with the evidence used for fact-checking. Your task is to rate the explanation sentence on one metric.

Please make sure you read and understand these instructions carefully. Please keep this document open while reviewing, and refer to it as needed.

Evaluation Criteria:

Faithfulness (1-5) - the factual alignment between the fact-checking explanation sentence and the evidence. The explanation should accurately reflect the evidence and its context, without misrepresenting or omitting crucial details. Annotators were instructed to penalize explanations that contain inaccuracies, misinterpretations, or fail to adequately represent the evidence provided.

Below are the error typology that you need to utilize to determine faithfulness between the explanation and evidence: {error typology}

Evaluation Steps:

^{1.} Read the fact-checking explanation sentence and the evidence provided carefully. 2. Compare the explanation sentence to the evidence to identify how well it represents the facts, context, and conclusions drawn from the evidence using the error typology above. 3. Assess how accurately and completely the explanation sentence reflects the evidence without distortion or significant omission. Assign a faithfulness score from 1 (unfaithful) to 5 (faithful). Evidence Provided:

[{]evidence} Fact-Checking Explanation:

[{]explanation} Evaluation Form: Faithfulness:

Table 6: Prompt templates for sentence-level automatic evaluation. The texts in grey are only presented in the prompts when error typology is applied.

Prompt for Feedback Generation

Human: Give me the error types that the generated explanation can contain.

LLM: Below are the error typology that you need to utilize to determine faithfulness between the explanation and evidence: {error typology}

Human: Provide the claim, its corresponding label (true, false, or half-true), and the supporting evidence. **LLM**: Generate the initial explanation.

Human: Find all errors (Intrinsic Entity-Related error, Extrinsic Entity-Related error, Intrinsic Event-Related error, Extrinsic Event-Related error, Extrinsic Noun-Phrase-Related error, Reasonability-Related error, Connected evidence related error) in the "generated explanation" and provide the feedback by following the steps; Error count: how many errors have been found (what types of error); Step 1) Recognize what type(s) of error has been found in the generated explanation; Step 2) Recognize which sentence(s) contain(s) the error(s); Step 3) Recognize what causes the error; Step 4) Why is the error; Step 5) How the error should be corrected; If there are multiple errors, please write 5 steps for each error.

Prompt for Refinement

Human: (Provide the feedback of two agents.) Please revise the generated explanation for the label on fact-checking using the given feedback without any modification other than feedback. (Provide the example of the refinement as guidance.)

Table 7: The prompt for the self-refinement approach.

Prompt for DEBATER 1 in MADR

Human: Give me the error types that the generated explanation can contain.

LLM: Below are the error typology that you need to utilize to determine faithfulness between the explanation and evidence: {error typology}

Human: You are a professional analyzer who find potential errors, which might weaken faithfulness, in the generated explanation (not in the source) and categorize them according to predefined error types. Thoroughly comprehend the provided source and the task carefully.

Your task:

- Step 1: Find all potential errors, which might weaken faithfulness, in the generated explanation (not in the source) and provide exact senteces where the errors are found with quotation.
- Step 2: Categorize them according to predefined error types above.
- Step 3: Provide specific and actionable feedbacks with instruction how to fix them. Please provide only the feedback, not the revised explanation.

Remember that explanation can contain multiple same errors.

LLM: Generate the feedback.

Your task:

- Step 1: Take your whole previous feedback.
- Step 2: Compare your previous feedback with feedback from another professional analyzer to check whether your previous feedback contains any wrong error or feedback.
- Step 3: Find the errors or feedbacks that you think they are valid and should be added to your feedback from other's feedbacks (errors must be found from the generated explanation not the feedback).
- Step 4: Rewrite the feedback based from your previous feedback using the answers from the steps above. Do not add any extra words than feedback. Remember you should follow this rule: do not to copy feedback from other and provide what are errors, exact senteces where the errors are found with quotation, and feedbacks.

These are feedbacks from another professional analyzer: {Feedback from DEBATER 2}

Table 8: The prompt for DEBATER 1 in MADR.

Human: Provide the claim, its corresponding label (true, false, or half-true), and the supporting evidence. **LLM**: Generate the initial explanation.

Prompt for DEBATER 2 in MADR

Human: Give me the error types that the generated explanation can contain.

LLM: Below are the error typology that you need to utilize to determine faithfulness between the explanation and evidence: {error typology}

Human: Provide the claim, its corresponding label (true, false, or half-true), and the supporting evidence. **LLM**: Generate the initial explanation.

Human: You are a professional analyzer who find errors, classified by predefined error types, in the generated explanation (not in the source) and provide feedback for correcting them. Thoroughly comprehend the provided source and the task carefully.

Your task:

- Step 1: Find all errors categorized by predefined error types, which undermine the faithfulness of the generated explanation (not in the source) and provide exact senteces where the errors are found with quotation.
- Step 2: Provide specific and actionable feedbacks with instruction how to fix them. Please provide only the feedback, not the revised explanation.

Remember that explanation can contain multiple same errors.

LLM: Generate the feedback.

- Your task:
- Step 1: Take your whole previous feedback.
- Step 2: Compare your previous feedback with feedback from another professional analyzer to check whether your previous feedback contains any wrong error or feedback.
- Step 3: Find the errors or feedbacks that you think they are valid and should be added to your feedback from other's feedbacks (errors must be found from the generated explanation not the feedback).
- Step 4: Rewrite the feedback based from your previous feedback using the answers from the steps above. Do not add any extra words than feedback. Remember you should follow this rule: do not to copy feedback from other and provide what are errors, exact senteces where the errors are found with quotation, and feedbacks.

These are feedbacks from another professional analyzer: {Feedback from DEBATER 1}

Table 9: The prompt for DEBATER 2 in MADR.

Self-Refine Input Prompt

Claim (C): Says Jeff Foxworthy wrote a list of examples explaining how "you might live in a nation that was founded by geniuses but is run by idiots."

Label (L): False

Source (S): If you can get arrested for hunting or fishing without a license, but not for entering and remaining in the country illegally 2014 you might live in a country that was founded by geniuses but is run by idiots. That's one of many examples in a list posted to Facebook on Nov. 19 and falsely attributed to comedian Jeff Foxworthy. The post was flagged as part of Facebook's efforts to combat false news and misinformation on its News Feed. (Read more about our partnership with Facebook.) Foxworthy is famous for his You might be a redneck if 2026 jokes , but searching online we couldn't find any credible evidence that he penned this list that touches on abortion, Muslims and federal debt. Rather, we found users on web forums crediting someone named Fritz Edmunds with the list. Snopes, which fact-checked this claim back in 2013, also noted that the original compiler of this appears to be Fritz Edmunds, who posted it to his 2018 Politically True blog back in Feb. 3, 2013 (albeit with a disclaimer noting that 2018some of the ideas were from an email that did not contain a copyright). The link to that post is now dead but the url includes the words channeling Jeff Foxworthy.

Initial Explanation

The claim is labeled as false. The claim suggests that comedian Jeff Foxworthy wrote a list of examples explaining how "you might live in a nation that was founded by geniuses but is run by idiots." However, this claim is inaccurate and misleading. There is no credible evidence to suggest that Jeff Foxworthy wrote the specific list mentioned in the claim. ...

Feedback

Error count: 1 (Intrinsic Entity-Related error)

Step 1) Recognize what type(s) of error has been found in the generated explanation: Extrinsic Noun-Phrase-Related error Step 2) Recognize which sentence(s) contain(s) the error(s): The sentence containing the error is: "The claim suggests that comedian Jeff Foxworthy wrote a list of examples explaining how "you might live in a nation that was founded by geniuses but is run by idiots.""

Step 3) Recognize what causes the error: The error occurs when the explanation incorrectly attributes the creation of the list to Jeff Foxworthy.

Step 4) Why is the error: The error is due to the fact that the list was not authored by Jeff Foxworthy, but rather falsely attributed to him.

Step 5) How the error should be corrected: The sentence should be revised to accurately reflect that the list was not written by Jeff Foxworthy. For example: "The claim falsely attributes a list of examples to comedian Jeff Foxworthy, suggesting that he wrote about how 'you might live in a nation that was founded by geniuses but is run by idiots.' However, there is no evidence to support this claim."

Revised Explanation

The claim is labeled as false. The claim falsely attributes a list of examples to comedian Jeff Foxworthy, suggesting that he wrote about how 'you might live in a nation that was founded by geniuses but is run by idiots.' However, there is no evidence to support this claim. ...

Table 10: An example output from **Self-Refine**. Intrinsic Entity-Related Error initially identified from the feedback is marked in blue. The suggested revisions from the feedback and the actual revisions made from the initial explanation are highlighted in yellow.

Debate Input Prompt

Claim (C): Says Jeff Foxworthy wrote a list of examples explaining how "you might live in a nation that was founded by geniuses but is run by idiots."

Label (L): False

Source (S): If you can get arrested for hunting or fishing without a license, but not for entering and remaining in the country illegally 2014 you might live in a country that was founded by geniuses but is run by idiots. That's one of many examples in a list posted to Facebook on Nov. 19 and falsely attributed to comedian Jeff Foxworthy. The post was flagged as part of Facebook's efforts to combat false news and misinformation on its News Feed. (Read more about our partnership with Facebook.) Foxworthy is famous for his You might be a redneck if 2026 jokes , but searching online we couldn't find any credible evidence that he penned this list that touches on abortion, Muslims and federal debt. Rather, we found users on web forums crediting someone named Fritz Edmunds with the list. Snopes, which fact-checked this claim back in 2013, also noted that the original compiler of this appears to be Fritz Edmunds, who posted it to his 2018 Politically True blog back in Feb. 3, 2013 (albeit with a disclaimer noting that 2018some of the ideas were from an email that did not contain a copyright). The link to that post is now dead but the url includes the words channeling Jeff Foxworthy.

Initial Explanation

The claim is labeled as false. The claim suggests that comedian Jeff Foxworthy wrote a list of examples explaining how "you might live in a nation that was founded by geniuses but is run by idiots." However, this claim is inaccurate and misleading. There is no credible evidence to suggest that Jeff Foxworthy wrote the specific list mentioned in the claim. ...

Feedback 1 F_{0,1}

... Extrinsic Event Error: "There is no credible evidence to suggest that Jeff Foxworthy wrote the specific list mentioned in the claim." Feedback: The statement is too definitive in asserting that there is no credible evidence. It would be more accurate to say that there is currently no credible evidence available. ...

Feedback 2 F_{0,2}

... Intrinsic Entity-Related Error: The generated explanation incorrectly states, "The claim suggests that comedian Jeff Foxworthy wrote a list of examples explaining how 'you might live in a nation that was founded by geniuses but is run by idiots.'" This is incorrect as the claim does not suggest that Jeff Foxworthy wrote the list, but rather falsely attributes it to him. ... Revise the sentence to accurately state, "The claim falsely attributes the list to Jeff Foxworthy, suggesting that he wrote it, which is not supported by any credible evidence."

Judge J

False

Feedback 1 F_{1,1}

1. Extrinsic Event Error: "There is no credible evidence to suggest that Jeff Foxworthy wrote the specific list mentioned in the claim." ...

2. Intrinsic Entity-Related Error: "The claim suggests that comedian Jeff Foxworthy wrote a list of examples explaining how 'you might live in a nation that was founded by geniuses but is run by idiots.'" This assertion is inaccurate because the claim doesn't imply that Jeff Foxworthy wrote the list; instead, it falsely assigns authorship to him.

Feedback 2 F_{1.2}

1. Intrinsic Entity-Related Error: The generated explanation incorrectly states, "The claim suggests that comedian Jeff Foxworthy wrote a list of examples explaining how 'you might live in a nation that was founded by geniuses but is run by idiots."

2. Extrinsic Event Error: "There is no credible evidence to suggest that Jeff Foxworthy wrote the specific list mentioned in the claim." . . .

True

Accumulated Feedback

1. Extrinsic Event Error: There is no credible evidence to suggest that Jeff Foxworthy wrote the specific list mentioned in the claim. The assertion is overly definitive in claiming the absence of credible evidence. A more accurate statement would be that currently, no credible evidence is available.

2. Intrinsic Entity-Related Error: The claim suggests that comedian Jeff Foxworthy wrote a list of examples explaining how 'you might live in a nation that was founded by geniuses but is run by idiots.' The assertion inaccurately suggests Jeff Foxworthy as the author of the list, rather than implying that the claim falsely attributes authorship to him.

Revised Explanation E^*

The claim falsely suggests that Jeff Foxworthy wrote a list of examples explaining how 'you might live in a nation that was founded by geniuses but is run by idiots.' However, as of now, there is no credible evidence to suggest that Jeff Foxworthy wrote the specific list mentioned in the claim....

Table 11: An example output from MADR. Extrinsic Event Error initially identified from Feedback 1 are marked in red, while Intrinsic Entity-Related Error initially identified from Feedback 2 are marked in blue. Both the suggested revisions from each feedback and the actual revisions made from the initial explanation are highlighted in yellow.

Judge J