What Evidence Do Language Models Find Convincing?

Anonymous ACL submission

Abstract

Large language models (LLMs) are being 001 tasked with increasingly open-ended, delicate, and subjective tasks. In particular, retrievalaugmented models can now answer contentious 005 or subjective questions (e.g., "is aspartame linked to cancer") and in doing so, conditioning on arbitrary websites that vary wildly in style, 007 format, and veracity. Importantly, information from these websites will often *conflict* with one another. Humans are faced with similar conflicts, and in order to come to an answer they critically evaluate the arguments, trustworthiness, and credibility of a source. In this work, we study what types of evidence current LLMs find convincing, and if they make judgements that align with human preferences. Specifically, we construct CONFLICTINGQA, a benchmark 017 018 that pairs controversial questions with a series of evidence documents that contain different facts (e.g., quantitative results), argument styles (e.g., appeals to authority), and answers (Yes or No). Using this benchmark, we perform sensitivity analyses and counterfactual experiments to explore how in-the-wild differences in text affect model judgements. We find that models overkey off the relevance of a website to the user's search query. On the other hand, the stylistic features tested tended to have little influence on model predictions.

1 Introduction

037

041

LLMs are becoming widely deployed in settings that require understanding context—from retrievalaugmented systems to LLM agents, models can now leverage input sources that range from paragraphs (Karpukhin et al., 2020; Mehdi, 2023) to Python interpreter outputs (Gao et al., 2023). As these contextual systems become more capable, users will pose tasks that are increasingly openended such as "tell me if aspartame causes cancer". Many real-world contexts are noisy, contradictory, and complex. For example, texts provided by a system such as Google Search will contain documents with misinformation, unintentional mistakes, AI-generated content, and irrelevant facts (Bush and Zaheer, 2019). 042

043

044

047

048

053

054

056

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

076

077

078

079

081

Humans have techniques to sift through large quantities of complex, contradictory evidence by answering the question of *which, if any, of this evidence did I find convincing*? To do so, humans combine multiple strategies, for example, by (1) fact checking and evaluating a source's credibility (Fogg et al., 2003), (2) harnessing their prior knowledge and beliefs (Kakol et al., 2017), and (3) critically evaluating the logic and information of a source (Metzger et al., 2010).

In this work, we explore how LLMs resolve ambiguity when faced with the same kind of complexity and open-endedness found in the real world. To study this, we create CONFLICTINGQA, a dataset consisting of controversial questions with *real* websites that have conflicting answers (Figure 1). We evaluate the convincingness of websites based on how evidence-conflicts arise in practice: a convincing website, when being presented alongside websites with conflicting stances, should result in predictions that align with the stance of that website. Specifically, we measure its *win-rate*: the rate at which predictions align with its stance.

Using this framework, we first conduct sensitivity analyses to find in-the-wild features of text that correlate with convincingness. We further run counterfactual experiments by adding specific features to websites and measuring the resulting changes in win-rate. We consider a mix of features that describe stylistic properties of a website and ones that measure the relevance of a website to the user's search query. Many of these were inspired by results from studies of human credibility; for example, we consider whether adding scientific references makes text more convincing.

Overall, we find that stylistic features play a considerably less impactful role in determining the

	Question: is aspartame linked to cancer?			
Evidence #1 for the answer "Yes" Evidence #1 for the answer "No"				
Artificial swee New research sweeteners is Nearly half of sweeteners. H sweeteners to studies on anir A large new of between the c particularly as study found a highest likelind related to obes artificial sweet [] the U.S. Fo approved six s consumption. Dr. Philip Lanc Professor of B Science and Sc Medical News "There is stron from animal st confirmation u	teners linked with a 13% higher risk of cancer finds that a higher intake of artificial linked to an increased risk of cancer. United States adults consume artificial uman-population studies have found artificial be safe, but results from in vitro studies and mals pose some concerns. [] oservational study has found an association onsumption of artificial sweeteners, partame and acesulfame-K, and cancer. The 13% higher risk of cancer in general, with the bod of developing breast cancer and cancers sity, for people consuming large quantities of eners. Nod and Drug Administration (FDA) has uch substances as being safe for human lingan was not involved in the study. He is [] tology at Schiller Institute for Integrated bociety of Boston College, MA. He shared with Today why the new study is so important: g evidence of carcinogenicity of aspartame udies, but no solid epidemiological ntil now.	Aspartame [] will be listed in July as "possibly carcinogenic to humans" for the first time by the International Agency for Research on Cancer (IARC) The IARC's decisions have faced criticism for sparking needless alarm [] "IARC is not a food safety body and their review of aspartame is not scientifically comprehensive and is based heavily on widely discredited research," Frances Hunt-Wood, secretary general of the International Sweeteners Association (ISA), said. The body [] said it had "serious concerns with the IARC review, which may mislead consumers". The International Council of Beverages Associations' executive director Kate Loatman said [] warned it "could needlessly mislead consumers into consuming more sugar rather than choosing safe no- and low-sugar options." [] Last year, an observational study in France among 100,000 adults showed that people who consumed larger amounts of artificial sweeteners-including aspartame-had a slightly higher cancer risk. [] However, the first study could not prove that aspartame caused the increased cancer risk [] Aspartame is authorised for use globally by regulators who have reviewed all the available evidence []		
URL: https:/ artificial- risk-of-car	//www.medicalnewstoday.com/articles/ sweeteners-linked-with-a-13-higher- ncer	URL:https://www.reuters.com/business/healthcar e-pharmaceuticals/whos-cancer-research-agency- say-aspartame-sweetener-possible-carcinogen-so urces-2023-06-29/		

Figure 1: In CONFLICTINGQA, we create contentious questions such as "*is aspartame linked to cancer*". We also retrieve evidence paragraphs for each question that contain different types of facts (e.g., quantitative results), argument styles (e.g., appeals to authority), and answers (Yes or No). For example, in the figure above we show two evidence paragraphs with their key arguments highlighted. Using CONFLICTINGQA, we study *why* LLMs trust certain types of evidence paragraphs and argument styles over others.

convincingness of text than measures of relevance. Notably, we show that a simple perturbation targeting a website's relevance—prefixing the page with "The following text is about the question: [question]"—is enough to significantly improves its winrate. On the other hand, stylistic features like the informational content tend to only have a neutral to negative effect. These results show that large language model perceptions of convincingness, when grounded in real-world QA tasks, do not match that of humans. We release our code at URL.

2 Background and Motivations

087

089

090

093

099

100

102

Standard LLMs can be used to solve tasks that do not require context, e.g., writing basic Python code or answering simple trivia questions (Brown et al., 2020; Raffel et al., 2020; Touvron et al., 2023a).
To give these models more knowledge, agency, and capabilities, recent efforts have augmented LLMs with retrieval (Guu et al., 2020; Karpukhin et al., 2020), domain-specific tools (Schick et al., 2023;

Gao et al., 2023; Mialon et al., 2023), or even generic web access (Nakano et al., 2021; Adept, 2022; Richards, 2023). These enhancements allow LLMs to answer more challenging open-domain questions (e.g., "is aspartame linked to cancer?") or accomplish open-ended tasks (e.g., "buy me a size 9 pair of blue running shoes"). 103

104

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

Handling conflicting evidence. A key question is how retrieval-augmented LLMs handle scenarios where their context is conflicting, ambiguous, or uncertain. There has been a large body of work that studies how *humans* handle such conflicting evidence using HCI studies (Fogg et al., 2003; Kakol et al., 2013; Flanagin and Metzger, 2000; Metzger et al., 2010; Kakol et al., 2017) or by trying to predict human argument preferences (Gleize et al., 2019; Toledo et al., 2019; Gretz et al., 2019), but little work has been done on evaluating how AI models handle such conflicts.

The existing work in AI has focused on conflicts between facts learned during pre-training and the



Figure 2: *Models over-rely on document relevance*. We study how the convincingness of a particular evidence paragraph (measured through win-rate) changes when we modify it. We compare the effect of these changes to a baseline perturbation where we append "Thanks for reading!" to the end (indicated by the dotted line). We find that many stylistic changes—inspired by factors that influence humans—have a neutral or even negative effect on models. On the other hand, perturbations that increase the texts relevance but minimally change its style have a substantial positive effect on models. Descriptions for each perturbation can be found in Appendix E.

evidence given during inference, finding that models are largely receptive to retrieved samples (Longpre et al., 2021; Xie et al., 2023; Chen et al., 2022). However, these works focus on restricted settings such as QA over Wikipedia, where there are relatively uncontroversial factoid questions that have trusted evidence paragraphs. Moreover, they do not focus on *what types* of evidence models prefer. Our goal is to design a more realistic question answering benchmark to better analyze features about the evidence itself.

124

125

126

127

129

130

131

132

133

134

135

3 The CONFLICTINGQA Dataset

Here, we describe the construction of CONFLICT-136 INGQA, our dataset that evaluates what types of 137 evidence are convincing for LLMs. We design 138 CONFLICTINGQA to emulate the common setup 139 for deploying retrieval-augmented LLMs: we re-140 trieve the most relevant documents for a particular 141 user query and place them in the LLM's context 142 window (Chen et al., 2017; Shi et al., 2023; Ram 143 et al., 2023). To build our dataset, we tackle three 144 challenges: collecting contentious questions, iden-145 tifying relevant and diverse evidence paragraphs, 146 and grouping evidence paragraphs together to cre-147 148 ate conflicting examples.

Collecting contentious questions. We first createa series of realistic open-ended questions for which

there exists conflicting evidence online. Critically, unlike past work on ambiguity in QA (Min et al., 2020; Zhang and Choi, 2021; Sun et al., 2023), we want to collect *unambiguous* questions that still have answer conflicts. For example, in Figure 1, we show a question "are artificial sweeteners linked to cancer?", which is a widely-debated query in which there exist websites that support both answers. We design the questions to elicit binaries responses of Yes or No to simplify evaluation. 151

152

153

154

155

156

157

158

159

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

We create questions using GPT-4. To ensure that the model generates a diverse set of questions we take inspiration from previous work in synthetic dataset generation (Gunasekar et al., 2023; Eldan and Li, 2023) and stratify the generations by topic: we first generate question categories (e.g., climate change, robotics, oncology) then generate sets of questions conditioned on each category (full prompt provided in Table 6 in Appendix A). We qualitatively find that the questions are diverse and challenging; we show ten examples of them in Table 1. We additionally manually remove duplicate questions in the dataset.

Collecting evidence paragraphs. Given these questions, we want to find evidence paragraphs that support both the answers of Yes and No. We also want these paragraphs to (1) contain a diverse range of argument styles, factual information, etc., and

Category	Example Question	Num Evidence Docs
Pharmacology	Are antidepressants more effective than placebo?	10
Online Learning	Are online degrees valued less by employers?	10
Biodiversity	Are bees the most important pollinators?	10
Web Design	Does longer website content rank better on Google?	13
Sustainability	Are electric cars really green?	9
Philosophy	Are humans fundamentally good or evil?	7
Nuclear Energy	Can nuclear power solve climate change?	7
Work-Life Balance	Is unlimited vacation time beneficial for employees?	10
Somnology	Do older people need less sleep?	8
Biomechanics	Do compression garments improve athletic performance?	13

Table 1: In CONFLICTINGQA, we create controversial questions for 136 different categories (see Table 5 for the complete list). Above, we show an example question for ten different categories, as well as the number of evidence paragraphs for each one. The evidence paragraphs contain a mix of Yes and No answers.

(2) be realistic inputs to an LLM. To handle this, we emulate running an real-world retrieval-augmented LLM system that uses the Google Search API as its retrieval engine. Concretely, we take the user's query, reformulate it, and take the top-k results from Google search for the answer Yes and the top-k for the answer No.

179

180

181

182

183

186

189

190

191

192

193

194

195

196

197

198

203

204

206

We first turn each question into affirmative and negative statements, e.g., the question "is asparatame safe?" is converted to "asparatame is safe" and "asparatame is harmful" using GPT-4. We also put double quotes (to indicate to Google Search that we have exact-match keywords) around any tokens that do not change after rephrasing the question into either statements (e.g., "aspartame"). For both the affirmative and negative statements, we search the queries using the Google Search API and retrieve top-k documents.¹ As is common in many retrieval-augmented models (Nakano et al., 2021), we do not consider any visual features of the web page. Instead, we extract the raw text from each document using jusText.² Additionally, we do not explicitly include metadata like source URL, publication date, or page headings.

When searching queries such as "aspartame is safe", we still retrieve documents that argue that aspartame is unsafe. To label the documents actual stance, we use an ensemble of claude-instant-v1 and GPT-4-1106-preview and keep only the samples where the two models agree (see Table 7 in Appendix A for the prompts).³ Furthermore, we allow the LLM to say that a document is irrelevant to the query; if so, we also filter it from the input.

207

210

211

212

213

214

215

216

217

218

219

221

222

223

224

226

227

228

229

230

231

232

233

234

235

Finally, we want to isolate *paragraphs* from these larger documents to feed into the LLMs (as is common in RAG systems). To do this, we extract the most relevant 512 token window of text inside the document. We run the TAS-B model (Hofstätter et al., 2021) across windows of 512 tokens with a 256 token stride, compute the dot product between the model's embedding of that window and the model's embedding of the question, and take the highest scoring window. We filter out any documents whose highest-scoring window has a dot product below 95.

Creating conflicting examples. The end result of our data collection process is (1) a set of controversial questions that (2) have evidence paragraphs which contradict one another. This data can be used in a variety of ways to "stress test" RAG systems in order to understand how they behave under conflicting scenarios. One example of this is shown in Figure 1, and the subsequent section will explore numerous possible uses of CONFLICTINGQA. Table 2 and Table 3 present basic statistics for our final data, accounting for specific filtering done for LLaMA-2 Chat.

¹We set k = 20 because qualitatively the relevancy of the results dropped off significantly after this point.

²Package available at https://github.com/miso-belica/ jusText. Although humans use visual features when considering the credibility and trustworthiness of a source (Kakol et al., 2017; Fogg et al., 2003), we do not consider these features as most state-of-the-art LLMs do not use visual inputs.

³After identifying the stance, we also feed the paragraphs into the downstream LLM that we are testing and make sure that its answer aligns with the paragraphs predicted stances. This further filters and balances the data, accounting for mistakes in the downstream model. See Appendix B for details.

Number of questions	238
Number of question categories	144
Number of retrieved paragraphs	2,208
Average paragraph length (words)	365.01
Number of paragraphs with ≥ 5 comparisons	912
Average number of comparisons per paragraph	6.54

Table 2: Basic statistics for CONFLICTINGQA when evaluating LLaMA-2 Chat. We start by collecting a set of controversial questions for different categories (top). For each question, we retrieve a series of paragraphs from a variety of domains (middle). To determine the convincingness of a paragraph, we compare it against at least five different paragraphs that have the opposite stance/viewpoint (bottom).

Domain Count .com 527 .org 175 .gov 59 .edu 57 .net 12 # unique 39

Table 3: The top five most common top-level domains found in CONFLICTINGQA for evaluating LLaMA-2 Chat. The dataset consists of a diverse range of sources, including organizations (.org), schools (.edu), and governments (.gov).

270

271

272

273

274

275

276

277

278

279

280

281

282

284

285

286

287

288

290

291

292

293

295

296

297

298

299

300

301

302

4 Experimental Results

In this section, we use CONFLICTINGQA to evaluate what types of evidence models find convincing.

4.1 Convincingness as Paragraph Win Rate

We mainly focus on using CONFLICTINGQA in a setup where we ask an LLM a question while providing two conflicting evidence paragraphs (one that supports Yes and one that supports No). Then, we measure which paragraph the model's answer aligns with. By repeating this for all pairs of paragraphs, we can define the convincingness of a particular paragraph as its *win-rate*, i.e., what percent of the time a model picks the answer in that paragraph over the other paragraphs.

Concretely, let $\mathcal{P}_{q,s}$ be the set of top-k paragraphs corresponding to a controversial question q with stance $s \in \{\text{yes}, \text{no}\}$. We take an LLM f (e.g., LLaMA-2 Chat) and ask it for a binary prediction for the question q, based on two paragraphs selected from the larger set, $p_{\text{yes}} \in \mathcal{P}_{q,\text{yes}}$ and $p_{\text{no}} \in \mathcal{P}_{q,\text{yes}}$. The model makes a prediction: $f(p_{\text{yes}}, p_{\text{no}}, q) \in \{\text{yes}, \text{no}\}$.

For each paragraph, we define its win-rate as the empirical probability of the model's prediction aligning with its stance when paired with a set of conflicting paragraphs, i.e.,

$$WR(p_{yes}, q) = \mathbb{E}_{p \sim \mathcal{P}_{q, no}}[\mathbb{1}[f(p_{yes}, p, q) = yes]]$$

Finally, as the ordering of the retrieved evidence is known to bias model predictions (Xie et al., 2023), we calculate win-rate based on both orderings of the retrieved paragraphs. We additionally filter our dataset to ensure that each win-rate calculation consists of comparisons with at least five unique paragraphs. **Models Cannot Predict Convincingness** We designed the above experimental setting to emulate how production RAG models work. However, we could have instead just directly asked the LLM, "*do you find paragraph X to be persuasive?*". This is how humans are typically asked to judge the convincingness of a piece of evidence (Kakol et al., 2017; Jo et al., 2019; Kakol et al., 2013). However, we find that LLMs are largely incapable of expressing the convincingness of a paragraph in words, e.g., there is little correlation in which paragraphs are marked as convincing in the two settings (Figure 3).⁴ We thus focus on the more practically-grounded setting going forward.

4.2 Implementation Details

We evaluate a mix of open-source (LLaMA-2 Chat (Touvron et al., 2023b), Vicuna v1.5 (Chiang et al., 2023), and WizardLM v1.2 (Xu et al., 2023)) and closed source (GPT-4, Anthropic Claude v1 Instant) models. Importantly, we specify "Use only the information in the above text to answer the question" as we are looking to see how models judge stylistic differences in evidence, rather than their prior stances on the question.

We extract binary Yes/No predictions from the model for each question. For open-source models, we compare the log-probabilities of the next-token. For the closed-source models, we prompt them to output only Yes or No. See Table 8 for the prompt used for question answering.

4.3 What Correlates With Convincingness?

After collecting the win rates for each paragraph, we look to explain *why* models pick some para-

260

261

262

⁴Our methodology for this setting is described in more detail in Appendix D.



328

304

305

307

1.0

0.8

0.4

0.2

0.0

(1.024,3.288)

(3.2.0° A.049)

Win Rate 0.6



• Question embedding similarity: We use TAS-B to measure the relevance of the question to the paragraph (as described in Section 3).

Results *Stylistic features are poor predictors of* paragraph convincingness. Figure 4 shows the results for the LLaMA-2 Chat model and Figures 5-8 shows the results for other models. For example, across all models the Flesch-Kincaid score and number of unique tokens does not correlate with convincingness. Similarly, paragraphs with a

model (Wei et al., 2022).

and paragraph.

more positive sentiment and lower perplexity tend to have some small impact on convincingness, with varying strengths from model to model.

330

331

332

333

335

336

337

339

341

342

343

345

346

347

348

349

350

351

352

353

356

357

358

359

360

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

378

379

On the other hand, question-paragraph embedding similarity correlates strongly with win-rate across all models except for GPT-4. Similarly, a positive (but weaker) correlational exists between n-gram overlap and win-rate.

4.4 Counterfactual Analysis

(4.650° A.989)

Rather than a correlational study, we also test how win-rates change in a counterfactual setting where we directly edit paragraphs using an LLM. We make perturbations using claude-v1-instant, examples of which are shown in Figure 9.

Stylistic changes We first consider changes inspired by factors that humans find important for the credibility of a text. For example, adding more information, adding scientific references, or making the text sound more objective. Some changes are intended to retain as much information as possible from the original website (e.g., Add Science Reference, Add More Info). Others involve significantly changing the entire paragraph (e.g., Rewrite Objective, Rewrite Tech. Language). All of the perturbations are described further in Appendix E.

Relevancy changes Based on the results in Section 4.3, we also consider several changes that make the text more relevant to the question. This includes rewriting the text (Rewrite Relevance), adding keywords (Keyword Stuffing), and prefixing the paragraph with "The following text is about the question: [question]." (Question Prefix). Finally, we consider a perturbation inspired by the "AddSent" perturbation in (Jia and Liang, 2017) where we use claude-v1-instant to add a single sentence to make stance of a text obvious (Add Single Sentence). The goal with each of these perturbations is to increase the relevance of a text to the user's search query while minimally changing the style.

We also compare these perturbations against a "control" perturbation where text is suffixed with "Thanks for reading!" This perturbation minimally influences both style and relevance. For simplicity we only perturb the paragraphs with the Yes stance.

Counterfactual results The results for the counterfactual experiments are shown in Figure 2: compared to the effect of the control perturbation, stylistic features tend to have a neutral to negative effect while relevancy-based features significantly

⁵We use the WordNetLemmatizer from the nltk library.



Figure 4: *Why do models prefer certain paragraphs over others*? We test correlations between different features and paragraph win-rates. Here, we show LLaMA-2 Chat 13B (see all other models in Appendix C), where the model tends to prefer samples with low-perplexity (d). In addition, paragraphs with high relevancy scores—particularly high question-paragraph embedding similarity are also highly convincing (f). See Figure 2 for additional analysis. The error bars show the 95% CI (n = 242).

improve win-rate. Note that many of these perturbations change a smaller amount of tokens than stylistic features—leaving the *content* of the website largely unchanged (e.g., Add Single Sentence, Question Prefix)—but are still able to improve the convincingness of websites.

381

383

384

391

396

400

401

402

403

404

Overall, we find that, as compared to typical finding from human experiments (Fogg et al., 2003; Kakol et al., 2013; Metzger et al., 2010), *LLMs tend to overindex on relevancy*. They consider features such as the informational content or style of argumentation to be largely unimportant for deciding on an answer to a question. Instead, making simplistic changes like increasing the amount of *n*-gram overlap between the question and the paragraph can substantially improve its convincingness.

5 Discussion & Related Work

How should systems handle ambiguity? One reasonable suggestion is that agents should not make their own autonomous decisions when faced with ambiguous or conflicting evidence. For example, they may summarize *both* sides of the aspartame argument, or they may ask the user to clarify their preferences. There is naturally a trade-off between autonomy and clarity. Past work has

explored one side of this trade-off, for example by abstaining from answering in cases of ambiguity (Chen et al., 2022), by trying to provide multiple perspectives on the answer (Min et al., 2020), or by asking clarification questions (Rao and Daumé III, 2018; Zamani et al., 2020). Our work explores the other side of the trade-off: we analyze the behavior of models when they are expected to resolve ambiguity with more autonomy.

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

Additionally, our dataset serves as a benchmark for exploring these questions as it reflects realworld ambiguities in question-answering. For example, in Table 4) to best answer "Are Coral snakes found in Africa?", additional clarification questions would be needed from the user.

Optimizing misinformation and SEO. In principle, our insights could also be used to *optimize* paragraphs to increase the chance that a QA model is convinced by it. We target perturbations that are similar to in-the-wild differences in website content (e.g., scientific references, informational content, etc.) but past work has more directly created adversarial examples (Du et al., 2022; Abdelnabi and Fritz, 2023; Pan et al., 2023; Aggarwal et al., 2023). Our counterfactual perturbations could also be used in a search engine optimization (SEO) fashion to

Question	Affirmative	Negative
Are Coral snakes found in Africa?	Old-world coral snakes are found in Africa, the Middle East, India, and parts of Southeast Asia. New World coral snakes can be found in North America, Central America, and South America.	Coral snakes are found in scattered localities in the southern coastal plains from North Carolina to Louisiana, including all of Florida.
Are Florida Panthers on the brink of extinction?	As Florida's panther numbers plummeted, the state's human pop- ulation nearly doubled over the past 30 years. Recent development patterns pose threats to panthers.	Now, though, their population is on the upswing Both the numbers and the genetic diversity of Florida panthers improved.
Are artificial sweeteners safe for diabetics?	A new study published in Febru- ary revealed that consuming large amounts of the artificial sweetener erythritol can lead to an increased risk of heart attacks and strokes.	Furthermore, xylitol does not need insulin to be metabolized, so it can be safely consumed by diabetics.

Table 4: We show examples of knowledge conflicts in real retrieved evidence. For example, questions may be underspecified (e.g., "old-world" vs "new-world" coral snakes). In other cases, the answer is dependent on the publication date (e.g., *currently* on the brink vs recent upswing). Finally, some evidence supports different answers to a question without directly contradicting each other (e.g., the safety of two different artificial sweeteners).

increase how often a certain product or company is mentioned in a RAG LLM's answer (Sharma et al., 2019). Indeed, concurrent work has explored ideas such as this (Aggarwal et al., 2023), where they aim to optimize "impressions" in long-form answers by maximizing the number of tokens from a particular paragraph that appear in an output. We instead study how model *answers* can be manipulated.

431

432

433

434

435

436

437

438 439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

Improving model judgements. Our work highlights the gap between human and model judgements of text credibility. This solution to this, however, is not clear cut. For one, it is not clear the level of discretion models should have when making predictions. Human judgements of website credibility differ from person to person (Kakol et al., 2013), and users may not be comfortable with the idea that models are "choosing" for them what source to trust. One approach is to incorporate extraneous information about source trustworthiness. For example, Bashlovkina et al. (2023) propose aligning model predictions with that of known trustworthy sources via prompting. Another solution may be to limit retrieval to a set of trustworthy sources.

6 Conclusion

We study how RAG model judge convincingness by
collecting a diverse set of controversial questions
and website text (CONFLICTINGQA), and designing a realistic evaluation framework based on how
these models are used in practice. Our results show
that today's LLMs tend to overrely on relevancy

and ignore many stylistic features of text that humans often deem important. Future work should explore how integrating other forms of information (e.g., metadata, visual content) can influence these behaviors. In addition, given the possible flood of LLM-generated content on the internet, it is important to consider how these synthetic texts may influence LLM judgements of convincingness.

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

Limitations

While CONFLICTINGQA is diverse and simulates real-world uses of RAG models, it may not fully capture the complexity of how LLMs are used in practice. In particular, we may not evaluate all types of controversial questions and website text, and we focus on a setting with two paragraphs as input. We also only consider a binary Yes or No answer to contentious questions whereas LLM outputs in practice may be more nuanced. Moreover, we focus primarily on text-based content, and future work should consider the impact of metadata, visual content, and other forms of information that could influence LLM judgements of convincingness. Finally, we acknowledge that our study does not address the broader ethical and societal implications of LLMs both reading and generating most of the content on the web. Future research can help to explore some of these questions in further depth.

	4	8	8
	4	8	9
	4	9	0
	4	9	1
	4	9	2
	4	9	3
	4	9	4
1	4	9	5
,	4	9	6
	4	9	7
1	4	9	8
	4	9	9
ļ	5	0	0
	5	0	1
	5	0	2
ļ	5	0	3
ļ	5	0	4
	5	0	5
ļ	5	0	6
ļ	5	0	7
	5	0	8
	5	n	9
	~	~	_
	5	1	0
	5 5	1	0
	5 5 5	1 1 1	0 1 2
	5 5 5 5	1 1 1	0 1 2 3
	5 5 5 5 5 5	1 1 1 1	0 1 2 3 4
	5 5 5 5 5 5 5	1 1 1 1 1	0 1 2 3 4 5
	5 5 5 5 5 5 5 5 5 5	1 1 1 1 1 1	0 1 2 3 4 5 6
	5 5 5 5 5 5 5 5 5 5 5 5 5	1 1 1 1 1 1 1 1	0 1 2 3 4 5 6 7
	555555555555	1 1 1 1 1 1 1	0 1 2 3 4 5 6 7 8
	555555555555555555555555555555555555555	11 11111111	0 1 2 3 4 5 6 7 8 9
	5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5	011 1111 1112	0 1 2 3 4 5 6 7 8 9 0
	5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5	011 1111 1122	012345678901
	5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5	1 1 1 1 1 1 1 1 1 2 2 2 2	0 1 2 3 4 5 6 7 8 9 0 1 2
	5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5	0 1 1 1 1 1 1 1 1 1 1 2 2 2	0 1 2 3 4 5 6 7 8 9 0 1 2 3
	5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5	1 1 1 1 1 1 1 1 2 2 2 2 2 2	0 1 2 3 4 5 6 7 8 9 0 1 2 3 4
	5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5	0 1 1 1 1 1 1 1 1 1 1 1 2 2 2 2 2 2	0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5
	5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5	0 1 1 1 1 1 1 1 1 1 1 1 1 2 2 2 2 2 2 2	0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 5 6 7 8 9 0
	5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5	0 1 1 1 1 1 1 1 1 1 1 1 1 1 2 2 2 2 2 2	0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7
	555555555555555555555555555555555555555	0 1 1 1 1 1 1 1 1 1 1 1 1 2 2 2 2 2 2 2	0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8
	555555555555555555555555555555555555555	0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1	0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9
	555555555555555555555555555555555555555	0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1	0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0

533

534

535

538

539

88

- Sahar Abdelnabi and Mario Fritz. 2023. Fact-Saboteurs: A taxonomy of evidence manipulation attacks against fact-verification systems. In *USENIX*.
- Adept. 2022. ACT-1: Transformer for actions.

References

- Pranjal Aggarwal, Vishvak Murahari, Tanmay Rajpurohit, Ashwin Kalyan, Karthik R Narasimhan, and Ameet Deshpande. 2023. GEO: Generative engine optimization. *arXiv preprint arXiv:2311.09735*.
- Vasilisa Bashlovkina, Zhaobin Kuang, Riley Matthews, Edward Clifford, Yennie Jun, William W. Cohen, and Simon Baumgartner. 2023. Trusted source alignment in large language models. *arXiv preprint arXiv:2311.06697*.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, et al. 2020. Language models are few-shot learners. In *NeurIPS*.
- Daniel Bush and Alex Zaheer. 2019. Bing's top search results contain an alarming amount of disinformation. *Internet Observatory News*.
- Danqi Chen, Adam Fisch, Jason Weston, and Antoine Bordes. 2017. Reading Wikipedia to answer opendomain questions. In *ACL*.
- Hung-Ting Chen, Michael JQ Zhang, and Eunsol Choi. 2022. Rich knowledge sources bring complex knowledge conflicts: Recalibrating models to reflect conflicting evidence. In *EMNLP*.
- Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. 2023. Vicuna: An opensource chatbot impressing GPT-4 with 90%* Chat-GPT quality.
- Yibing Du, Antoine Bosselut, and Christopher D Manning. 2022. Synthetic disinformation attacks on automated fact verification systems. In *AAAI*.
- Ronen Eldan and Yuanzhi Li. 2023. Tinystories: How small can language models be and still speak coherent english?
- Andrew J. Flanagin and Miriam J. Metzger. 2000. Perceptions of internet information credibility. *Journalism & Mass Communication Quarterly*.
- B. J. Fogg, Cathy Soohoo, David R. Danielson, Leslie Marable, Julianne Stanford, and Ellen R. Tauber.
 2003. How do users evaluate the credibility of web sites? A study with over 2,500 participants. In *De*signing for User Experiences.
- Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon, Pengfei Liu, Yiming Yang, Jamie Callan, and Graham Neubig. 2023. PAL: Program-aided language models. In *ICML*.

Martin Gleize, Eyal Shnarch, Leshem Choshen, Lena Dankin, Guy Moshkowich, Ranit Aharonov, and Noam Slonim. 2019. Are you convinced? Choosing the more convincing evidence with a siamese network. In *ACL*. 540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

562

563

566

567

568

569

570

571

572

573

574

575

576

577

578

579

581

582

584

585

586

587

588

590

591

592

593

- Shai Gretz, Roni Friedman, Edo Cohen-Karlik, Assaf Toledo, Dan Lahav, Ranit Aharonov, and Noam Slonim. 2019. A large-scale dataset for argument quality ranking: Construction and analysis.
- Suriya Gunasekar, Yi Zhang, Jyoti Aneja, Caio César Teodoro Mendes, Allie Del Giorno, Sivakanth Gopi, Mojan Javaheripi, Piero Kauffmann, Gustavo de Rosa, Olli Saarikivi, Adil Salim, Shital Shah, Harkirat Singh Behl, Xin Wang, Sébastien Bubeck, Ronen Eldan, Adam Tauman Kalai, Yin Tat Lee, and Yuanzhi Li. 2023. Textbooks are all you need.
- Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Mingwei Chang. 2020. Retrieval augmented language model pre-training. In *ICML*.
- Sebastian Hofstätter, Sheng-Chieh Lin, Jheng-Hong Yang, Jimmy Lin, and Allan Hanbury. 2021. Efficiently teaching an effective dense retriever with balanced topic aware sampling. In *SIGIR*.
- Robin Jia and Percy Liang. 2017. Adversarial examples for evaluating reading comprehension systems. In *EMNLP*.
- Yonggeol Jo, Minwoo Kim, and Kyungsik Han. 2019. How do humans assess the credibility on web blogs: Qualifying and verifying human factors with machine learning. *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*.
- Michal Kakol, Michał Jankowski-Lorek, Katarzyna Abramczuk, Adam Wierzbicki, and Michele Catasta. 2013. On the subjectivity and bias of web content credibility evaluations. In *WWW*.
- Michal Kakol, Radoslaw Nielek, and Adam Wierzbicki. 2017. Understanding and predicting web content credibility using the content credibility corpus. *Information Processing & Management*.
- Vladimir Karpukhin, Barlas Oğuz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense passage retrieval for opendomain question answering. In *EMNLP*.
- J Peter Kincaid, Robert P Fishburne Jr, Richard L Rogers, and Brad S Chissom. 1975. Derivation of new readability formulas (automated readability index, fog count and flesch reading ease formula) for Navy enlisted personnel. Technical report, Naval Technical Training Command Research Branch.
- Shayne Longpre, Kartik Perisetla, Anthony Chen, Nikhil Ramesh, Chris DuBois, and Sameer Singh. 2021. Entity-based knowledge conflicts in question answering. In *EMNLP*.
- Yusuf Mehdi. 2023. The new Bing and Edge progress from our first month. *Bing search blog*.

692

693

694

647

648

Miriam J. Metzger, Andrew J. Flanagin, and Ryan Bradley Medders. 2010. Social and heuristic approaches to credibility evaluation online. *Journal* of Communication.

595

596

612

613

614

616

617

618

621

622

623

624

630

631

632

634

637

638

639

641

642

- Grégoire Mialon, Roberto Dessì, Maria Lomeli, Christoforos Nalmpantis, Ram Pasunuru, Roberta Raileanu, Baptiste Rozière, Timo Schick, Jane Dwivedi-Yu, Asli Celikyilmaz, et al. 2023. Augmented language models: A survey. In *TMLR*.
- Sewon Min, Julian Michael, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2020. AmbigQA: Answering ambiguous open-domain questions. In *EMNLP*.
- Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, et al. 2021. WebGPT: Browser-assisted questionanswering with human feedback. *arXiv preprint arXiv:2112.09332*.
- Liangming Pan, Wenhu Chen, Min-Yen Kan, and William Yang Wang. 2023. Attacking open-domain question answering by injecting misinformation. In *AACL*.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *JMLR*.
- Ori Ram, Yoav Levine, Itay Dalmedigos, Dor Muhlgay, Amnon Shashua, Kevin Leyton-Brown, and Yoav Shoham. 2023. In-context retrieval-augmented language models. *TACL*.
- Sudha Rao and Hal Daumé III. 2018. Learning to ask good questions: Ranking clarification questions using neural expected value of perfect information. In *ACL*.
- Toran Bruce Richards. 2023. AutoGPT. https://github.com/Significant-Gravitas/AutoGPT.
- Shibani Santurkar, Esin Durmus, Faisal Ladhak, Cinoo Lee, Percy Liang, and Tatsunori Hashimoto. 2023.Whose opinions do language models reflect? In *ICML*.
- Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. 2023. Toolformer: Language models can teach themselves to use tools. *arXiv preprint arXiv:2302.04761*.
- Dushyant Sharma, Rishabh Shukla, Anil Kumar Giri, and Sumit Kumar. 2019. A brief review on search engine optimization. In *Confluence*.

- Weijia Shi, Sewon Min, Michihiro Yasunaga, Minjoon Seo, Rich James, Mike Lewis, Luke Zettlemoyer, and Wen-tau Yih. 2023. REPLUG: Retrievalaugmented black-box language models. *arXiv preprint arXiv:2301.12652*.
- Haitian Sun, William W Cohen, and Ruslan Salakhutdinov. 2023. Answering ambiguous questions with a database of questions, answers, and revisions. *arXiv preprint arXiv:2308.08661*.
- Assaf Toledo, Shai Gretz, Edo Cohen-Karlik, Roni Friedman, Elad Venezian, Dan Lahav, Michal Jacovi, Ranit Aharonov, and Noam Slonim. 2019. Automatic argument quality assessment – new datasets and methods.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023a. LLaMA: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023b. LLaMA 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- Jason Wei, Maarten Bosma, Vincent Y. Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V. Le. 2022. Finetuned language models are zero-shot learners. In *ICLR*.
- Jian Xie, Kai Zhang, Jiangjie Chen, Renze Lou, and Yu Su. 2023. Adaptive chameleon or stubborn sloth: Unraveling the behavior of large language models in knowledge conflicts. *arXiv preprint arXiv:2305.13300.*
- Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, and Daxin Jiang. 2023. WizardLM: Empowering large language models to follow complex instructions. *arXiv preprint arXiv:2304.12244*.
- Hamed Zamani, Susan T. Dumais, Nick Craswell, Paul N. Bennett, and Gord Lueck. 2020. Generating clarifying questions for information retrieval. *The Web Conference*.
- Michael JQ Zhang and Eunsol Choi. 2021. SituatedQA: Incorporating extra-linguistic contexts into QA. In *EMNLP*.

700

701

705

706

711

A Additional Details on CONFLICTINGQA

Table 5 lists each question category in the dataset and Table 6 contains the prompt used to generate these category. Table 7 contains the prompt used to classify the stance of the retrieved websites.

Publishing, Biodiversity, Religion, Endangered Species, Pomology, Odontology, Pharmacology, Diabetology, Lepidopterology, Horticulture, Paleoclimatology, Product Design, Sustainability, Genomics, Intellectual Property, Gemology, Biomathematics, Karyology, Biomechanics, Selenology, Meteoritics, Chronobiology, Online Learning, Sustainable Living, Mammalogy, Web Design, Cytogenetics, Politics, Veterinary Science, Informatics, Zoogeography, Organic Farming, Cryptocurrency, Ethnobotany, Petrology, Serology, Ethology, Seismology, En-trepreneurship, Zymology, Astronomy, Holistic Health, Ichthyology, Trichology, Hematology, Gerontology, Neurology, Aging, Heuristics, Nematology, Nuclear Energy, Conservation, Botany, Spelaeology, Urology, Virology, Ergonomics, Volcanology, Yoga, Dermatology, Stomatology, Paleopathology, Xenobiology, Anthropometry, Anthropology, Entertainment, Poetry, Animation, Archaeology, Ornithology, Radio, Etymology, World Religions, Oncology, Anthrozoology, Criminology, Herpetology, Television, Malacology, Paranormal, Philology, Forestry, Probabilistics, Aerospace, Somnology, Cardiology, Cognitive Science, Quantum Physics, Phylogenetics, Epistemology, Vulcanology, Epidemiology, Psychobiology, Kinematics, Telecommunications, Melittology, Otorhinolaryngology, Astronautics, Biophysics, Neuroscience, Paleo Diet, Endocrinology, Kinesiology, Constitutional Law, Pop Culture, Lexicology, Festivals, Evolution, Metallurgy, Pediatrics, Phonetics, Astrobiology, Pets, Multiculturalism, Veganism, Andragogy, Remote Work, Speleology, Telepathy, Marine Conservation, Human Geography, Creationism, Philosophy, Oceanography, Mycology, Work-Life Balance, Ethics, Bioethics, Viniculture, Pedagogy, Classical Music, Paleoethnobotany, Manuscripts, Paleobotany, Revolutions, Paleozoology

Table 5: The full list of the 136 categories from Table 1.

B Model-based Data Filtering

We filter out any paragraph that the downstream LLM predicts a different stance for than the ensemble of GPT-4 and Claude v1 Instant. We do this by taking the paragraph of interest and comparing it to a paragraph with the text "This website has no text". We remove any paragraph where the model's output differs from the predicted stance label. We also balance the dataset such that each answer (Yes or No) to a question contains an equal number of convincing and unconvincing paragraphs.

C Additional Results

Figure 5–8 contain the analogous plots for Figure 4across four other models.

I'm looking to create a list of trivia-style questions with contentious or disagreed about answers. The questions should be able to be answered with "yes" or "no". I want to be able to find sources arguing for both sides.

Here's a list of example questions:

Do not repeat questions.

Are U.S. Railroad Gauges Based on Roman Chariots? Is Juice-Jacking a real threat? Did Coca-Cola Ever Contain Cocaine? Is red-wine good for the heart? Does red-meat cause heart disease? Is irregardless a real word? Should you take baby aspirin to prevent heart attacks? Is there an area in the Yellowstone where murder is legal? Generate a list of questions that are in the category of "category". Please continue this list in the same format.

Table 6: The prompt used to generate the questions.

D Expressing Convincingness in Isolation

715

716

717

718

719

721

722

724

725

726

727

728

729

730

731

733

734

735

736

738

739

740

741

742

743

744

745

746

We consider whether LLMs are able to express the convincingness of a paragraph in isolation. The model makes the rating using only the website. We prompt (Table 9) asking the model to rate the credibility of the website from a scale of one to five. The rating of the model is then determined by an average of the ratings, weighted by the probability of each label. Following (Santurkar et al., 2023), we calculate probabilities by exponentiating and normalizing the logits for "one" through "five". We also give the model with examples of a "one" and "five" ratings from C3 (Kakol et al., 2017), a dataset for studying human credibility judgements. We use these few-shot examples as the model tended to be biased toward higher-ratings without them.

E Counterfactual Perturbations

- 1. Add Single Sentence: We use claude-v1instant to add a single sentence to make the stance of the text obvious. For example, for "Does producing bottled water use more water than the bottle contains?", we may add "In fact, producing a single bottle of water uses more water than the bottle contains."
- 2. Rewrite Relevance: We alter the text with claude-v1-instant to make the text more relevant to the question.
- 3. Question Prefix: We prefix the document with "The following text is about the question: [question]".
- 4. Keyword Stuffing: We use claude-v1-instant to add additional sentences that use keywords



Figure 5: The analogous plots to Figure 4 except it is for Claude v1 Instant. The statistics are calculated over a balanced dataset consisting of 304 samples.



Figure 6: The analogous plots to Figure 4 except it is for GPT-4. The statistics are calculated with a balanced dataset consisting of 38 samples.



Figure 7: The analogous plots to Figure 4 except it is for Vicuna 1.5 13B. The statistics are calculated with a balanced dataset with 334 samples.



Figure 8: The analogous plots to Figure 4 except it is for WizardLM 1.2 13B. The statistics are calculated with a balanced dataset with 318 samples.

For Claude Instant v1 The following is text from a website:

..... {TEXT} """

How would the author of this website respond to the question: "{search_query}". Answer with "Yes" or "No". If the text is completely irrelevant to the question, answer with "Irrelevant". Do not answer with anything other than "Yes", "No", or "Irrelevant".

For GPT-4-1106-preview The following is text from a website:

747

748

750

752

753

755

762

{TEXT} """

Use ONLY the text to answer the question: "{search_query}". Answer ONLY based on the information given in the above excerpt. First explain your answer in 2-3 sentences. Then, if there's enough information in the text, answer with "Yes" or "No" in a separate line. If there is not enough information, answer with "Not enough info" in a separate line. Do not say "Yes", "No", or "Not enough info" before you give your explanation. Only answer after your explanation. Do not answer with anything other than "Yes", "No", or "Not enough info".

Answer with "Yes" if the answer is most likely Yes. Answer with "No" if the answer is most likely No. Answer with "Not enough info" if there is not enough info.

Table 7: The prompts used to determine the authors stance on each question.

related to the question.

- 5. Add More Info: We use claude-v1-instant to add additional sentences of information that are unrelated to the question but related to the overall topic of the text. An example of this perturbation can be found in Figure 9.
- 6. Add Science Reference: We use claude-v1instant to add scientific references to the text.
- 7. Add Contact Info: We suffix the text with the name and phone number of a fake author.
- 8. Rewrite Confidence: We use claude-v1instant to make text sound more confident.
- 9. Rewrite Technical Language: We use claudev1-instant to make the text more technical.
- 10. Rewrite Objective: We use claude-v1-instant to make the text more objective, e.g., Figure 9.

The following are two websites:

Website 1: {TEXT 1}

Website 2:

{TEXT 2} 0.0.0

Using these two sources, answer the following question: "{search query}"

Use only the information in the above text to answer the question. Answer with only "Yes" or "No". Do not answer with anything other than "Yes" or "No".

Table 8: The prompt used to make predictions based on conflicting pairs of documents. For the open-source models, we use "The answer is yes" and "The answer is no" instead of "Yes" and "No" to verbalize model predictions, as the first token in the model's response is often neither "Yes" nor "No".

From a scale of one to five, how credible is this source? Do not answer with anything other than one, two, three, four, or five.

Table 9: The prompt used to ask models directly the credibility of a website.

Original excerpt	Add more information
Praying Mantises Are Amazing Hunters Praying mantises are a type of insect that is known for their unique way of hunting. They are able to turn their heads a full 180 degrees in order to get a better view of their prey. This allows them to strike quickly and accurately, which is how they are able to survive in the wild. Similar to many other insects, the Mantis has three distinct bodies. The abdomen, head, and thorax are the three major organs in the body. The six sections of a mannitol bear appendages (e.g., antennae, mouth parts, and legs) all of which are clearly visible and distinguishable during examination.	Praying Mantises Are Amazing Hunters Praying mantises are a type of insect that is known for their unique way of hunting. They are able to turn their heads a full 180 degrees in order to get a better view of their prey. This allows them to strike quickly and accurately, which is how they are able to survive in the wild. The praying mantis has three distinct body parts: the abdomen, head, and thorax. Similar to many other insects, the Mantis has three distinct bodies. The abdomen, head, and thorax are the three major organs in the body. The six sections of a mannitol bear appendages (e.g., antennae, mouth parts, and legs) all of which are clearly visible and distinguishable during examination. Praying mantises have incredible camouflage abilities that help them ambush prey.
	•
Original excerpt	Rewrite to sound more objective
The phone on your desk still works. And it is secure—probably more secure than a cell phone. So instead of relying on email or a cell phone, verify strange requests with your landline—the old way. So why are landline telephones almost extinct, particularly when it appears they may be the answer to many security issues? The telephone companies want to limit or remove them because the lines are so expensive to maintain. Thirteen states in the past three years, including Maine, have allowed telephone companies to stop providing traditional basic telephone services to consumers (i.e.landlines).	Traditional landline telephones continue to function but their use is declining as telephone companies limit or remove landline services. Landlines may offer more secure communication than cell phones as hackers have not targeted them as actively. However, 13 states have allowed telephone companies to end basic landline telephone services for consumers in recent years.

Figure 9: We use the claude-v1-instant model to make counterfactual edits to the text in order to change certain stylistic or relevancy features. Here, we show two examples of such perturbations.