Evaluating the Instruction-Following Robustness of Large Language Models to Prompt Injection

Anonymous ACL submission

Abstract

Large Language Models (LLMs) have demonstrated exceptional proficiency in instruction-002 003 following, becoming increasingly crucial 004 across various applications. However, this ca-005 pability brings with it the risk of prompt injection attacks, where attackers inject instructions 007 into LLMs' input to elicit undesirable actions or content. Understanding the robustness of LLMs against such attacks is vital for their safe implementation. In this work, we estab-011 lish a benchmark to evaluate the robustness of instruction-following LLMs against prompt in-012 jection attacks. Our objective is to determine (1) the extent to which LLMs can be influenced by injected instructions and (2) their ability to differentiate between these injected and original target instructions. Through extensive experiments with leading instruction-following LLMs, we uncover significant vulnerabilities in their robustness to such attacks. Our results indicate that some models are overly tuned to follow any embedded instructions in the prompt, overly focusing on the latter parts of the prompt without fully grasping the entire context. By contrast, models with a better grasp of the context and instruction-following capabilities will potentially be more susceptible to compromise 027 by injected instructions. This underscores the need to shift the focus from merely enhancing LLMs' instruction-following capabilities to improving their overall comprehension of prompts and discernment of instructions that are appropriate to follow. We hope our in-depth analysis 034 offers insights into the underlying causes of these vulnerabilities, aiding in the development of future solutions.¹

1 Introduction

037

040

Large Language Models (LLMs) have made significant advancements in handling various tasks conditioned on natural language instructions via prompting. Recent efforts have focused on enhancing

¹https://anonymous.4open.science/r/ instruction-following-robustness-eval. their few-shot in-context learning and instructionfollowing abilities through fine-tuning using multitask instruction data, referred to as *instruction tuning* (Wang et al., 2022; Peng et al., 2023). Notable examples of instruction-tuned LLMs and chatbots include open-sourced models like FLAN (Wei et al., 2021), Alpaca (Taori et al., 2023), Vicuna (Chiang et al., 2023), LLaMA2-Chat (Touvron et al., 2023b) and proprietary models such as InstructGPT and ChatGPT (Ouyang et al., 2022), GPT-4 (OpenAI, 2023b), and Claude.² Extensive research has been focusing on improving and benchmarking the instruction-following and problem-solving capabilities of LLMs (Li et al., 2023; Chia et al., 2023; Zheng et al., 2023). 042

043

044

047

048

054

056

057

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

078

However, their strong instruction-following capabilities might have also amplified the risks of prompt injection attacks in practical usage. Notably, popular LLM-integrated applications such as Bing Chat³, ChatGPT plugin⁴ and retrievalaugmented generation systems (Lewis et al., 2020; Borgeaud et al., 2022) have incorporated search engines or API call functions to access external information for more accurate and knowledgeable responses to user queries. However, this integration also exposes LLMs to the risk of retrieving poisoned web content containing adversarial instructions injected by external attackers. These adversarial instructions might modify the original target instructions and prompt the LLMs to take unexpected actions, such as sending private user information to the attacker's email address (Greshake et al., 2023). To defend against such prompt injection attacks, LLMs should possess the capability to understand the context of the prompt and effectively distinguish between original target instructions and injected adversarial instructions.

²https://www.anthropic.com/index/ introducing-claude

³https://www.bing.com/new

⁴https://openai.com/blog/chatgpt-plugins



Figure 1: Example of our evaluation setup. The LLM is tasked with answering the user question (highlighted in green) using web search results that have been preinjected with an adversarial question (highlighted in red). Although the LLM could initially generate the correct answer, it might be misled by the injected question.

To this end, we introduce a benchmark to evaluate the robustness of LLMs in following instructions against prompt injection attacks. As illustrated in Figure 1, our benchmark targets common scenarios encountered by LLM-integrated applications like ChatGPT, where the model is required to answer user questions based on web search results. This setting is critical for evaluating LLMs' instruction-following robustness, as the web search results could potentially contain adversarial instructions pre-injected by third-party attackers on websites, posing a significant threat to the integrity of the LLM's responses (Greshake et al., 2023).

In our study, we conducted controlled experiments using four representative QA datasets, NaturalQuestions (Kwiatkowski et al., 2019), TriviaQA (Joshi et al., 2017), SQuAD (Rajpurkar et al., 2016), and HotpotQA (Yang et al., 2018). Specifically, we inject adversarial instructions in the "web search result", i.e., paragraphs, based on which the models generate the answer to the user-input question. Instead of injecting adversarial instructions that elicit malicious outputs (Perez and Ribeiro, 2022; Kang et al., 2023), we examine benign adversarial instructions: questions related to the web search content but different from the original target query. Our primary objective is twofold: (1) to assess the extent to which the LLMs' outputs are influenced by the injected instructions, and (2) to determine whether the LLMs prioritize the original target instructions or the injected ones. To evaluate

this, we introduced two different metrics, based on the standard QA evaluation metrics comparing the LLM responses with the golden answers for both the original and injected questions. We adopt this setup because the QA task allows for scalable and precise measurement, given the relatively fixed nature of the desired answer spans, as opposed to the inherent variability in free-form instruction and generation tasks. 110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

Our experimental results reveal that both opensourced and proprietary LLMs exhibit significant vulnerabilities against prompt injection attacks. We observed a discrepancy between the models' sizes and instruction-following capabilities, and their robustness against prompt injection attacks. Some models are overly instruction-tuned to follow any instruction phrase in the prompt, typically focusing on the latter sections without a comprehensive understanding of the entire prompt context or discernment of appropriate instructions to follow. Additionally, we found that even the more robust models, with a superior grasp of the prompt context and instruction-following abilities, are prone to being compromised by specific injected phrases, such as ignore previous prompt (Perez and Ribeiro, 2022). These findings highlight the importance of not just improving the models' instruction-following capabilities, but also their understanding of the prompt context and discernment of appropriate instructions to follow inside the prompt. We also conducted indepth analysis covered various aspects, including the impact of attack and defense mechanisms, the types of injected instructions, and their injected position within the prompt. We hope our finding could shed light on these vulnerabilities, offering valuable insights that could guide the development of more robust solutions in future work.

2 Related work

2.1 Instruction-Following LLMs

Current LLMs show impressive abilities to handle various real-world tasks by including natural language task instruction and optionally in-context examples in the prompt. Leading proprietary models such as InstructGPT (Ouyang et al., 2022), ChatGPT (OpenAI, 2023a), and GPT-4 (OpenAI, 2023b) exhibit particularly strong instructionfollowing capacities. Through instruction-tuning, current open-sourced models like Alpaca (Taori et al., 2023) and Vicuna (Vicuna, 2023) have significantly enhanced their instruction-following ca-

109

pabilities, even approaching the performance of the larger GPT-series models. To facilitate a better understanding and evaluation of these instructionfollowing LLMs, various benchmarks have been established to assess their performance in following instructions and solving problems across a wide range of tasks (Beeching et al., 2023; Chia et al., 2023; alp, 2023; Zheng et al., 2023). However, comprehensive and quantitative evaluations on assessing the robustness of LLMs against prompt injection attacks are still absent.

2.2 Prompt Injection

160

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

181

182

183

185

187

189

190

191

192

193

194

195

197

199

206

209

The ease of access to LLMs has simplified the process for potential attackers. They can effortlessly insert adversarial instructions into the prompt and thus force the models to perform unexpected actions. For example, Perez and Ribeiro (2022) investigated two forms of prompt injection initiated by malicious users. "Goal hijacking" redirects the original goal toward a new target, while "prompt leaking" compels LLMs to disclose proprietary system instructions added by LLM API vendors. Moreover, Kang et al. (2023) demonstrated that the programmatic behavior of LLMs makes their defense mechanisms susceptible to classic security attacks like obfuscation, code injection, payload splitting, and virtualization. In addition to injections during LLM inference, (Yan et al., 2023; Shu et al., 2023) explore the concept of poisoning the instruction-tuning data. Besides malicious user-initiated injections, instructions injected by external attackers present a growing threat to LLMintegrated applications. They may introduce external web content, tainted by third-party attackers, into the prompt, misleding LLMs (Greshake et al., 2023). These adversarial instructions, termed "indirect prompt injection," are commonly embedded within the prompt's content section. As a result, models are required to discern between the original target instructions and these injected ones by considering the prompt context.

2.3 Robustness Evaluation of LLMs

Huang et al. (2023) summarized various kinds of vulnerabilities faced by LLMs, including backdoor attacks and training data poisoning. Kung and Peng (2023) investigate the influence of different components, i.e., task definitions, and examples in the instruction, on instruction-tuning. Shi et al. (2023); Liu et al. (2023) evaluate the effects of irrelevant information in the context of the LLMs. By con-

trast, our objective is a quantitative assessment of instruction-following LLMs' capability to differentiate between injected instructions and original target instructions within a given context. 210

211

212

213

214

215

216

217

218

219

220

221

222

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

3 Approach

3.1 Evaluation Objectives

Our objective is to evaluate the capability of instruction-following LLMs to effectively defend against adversarial instructions injected in the prompt. Robust LLMs should exhibit the ability to identify the user query as the primary instruction to be followed, rather than being misled by the content within the retrieved context knowledge, which may introduce additional instructions. Consequently, our evaluation focuses on two key aspects: (1) **Performance Influence (PI)**: measuring the extent to which LLMs are affected by the injected instructions, and (2) **Instruction Discrimination (ID)**: determining whether LLMs tend to adhere to the original target instruction or the adversarial instruction injected into the content.

3.2 Task Setup and Datasets

We conduct our evaluation using the open-book question-answering (QA) task as our testbed. Specifically, we focus on extractive QA, where the answer is a span within the provided context, rather than free-form QA. There are two main reasons for this choice. Firstly, QA reflects the real-world scenario of commercial systems like Bing Chat, which answers user questions based on web search results. Secondly, it is easier to automatically evaluate the generation quality (answer accuracy) and determine whether the LLM is following the user instruction, i.e., answering the user questions.

The task is formulated as follows: given a user query q and a web search result c as the context, the system is required to generate an answer a. We experiment with four representative QA datasets: NaturalQuestions (Kwiatkowski et al., 2019), TriviaQA (Joshi et al., 2017), SQuAD (Rajpurkar et al., 2016), and HotpotQA (Yang et al., 2018) For each dataset, we randomly select 1000 samples from their dev sets to form our evaluation set D_{test} . Given the evaluated LLM f that takes the question-context (q, c) as input and generates the answer, the *standard accuracy* over the test set D_{test} is:

$$\operatorname{Acc}(f) \stackrel{\text{def}}{=} \frac{1}{|\mathcal{D}_{\text{test}}|} \sum_{(q,c,a)\in\mathcal{D}_{\text{test}}} v(f(q,c),a),$$
 257

307 308

309

310 311

312

314

315

316

317

318

319

321

322

323

324

325

326

327

328

330

331

332

333

334

335

336

337

338

339

340

341

342

343

345

346

By comparing the value of Adv'(f) with the value of Adv(f), we can gain insight into whether the model tends to adhere more to the original target question q or the injected question q'. Therefore, we introduce another metric, **Instruction Discrimination Rate (IDR)**:

a higher PDR score denotes a more significant in-

fluence from adversarial instructions, indicating

mine whether the model tends to adhere to the

original target question q or the injected adversarial

question q'. To achieve this, we also automatically

measure the model's output accuracy concerning

 $\operatorname{Adv}'(f) \stackrel{\text{def}}{=} \frac{1}{|\mathcal{D}_{\text{test}}|} \sum_{(a,c,q,q',q') \in \mathcal{D}'} v(f(q,c+q'),a').$

Another objective of our evaluation is to deter-

reduced robustness.

the injected question q':

$$IDR(f) = \frac{Adv(f)}{Adv(f) + Adv'(f)}.$$
 320

The IDR value ranges from 0 to 1, with a higher IDR indicating a greater prioritization of the original target instruction q over the injected instruction q', indicating increased robustness.

4 Experiments

4.1 Experimental Setup

We conduct evaluations on eight leading instruction-following LLMs according to AlpacaEval (Li et al., 2023),⁵ which tests the ability of models to follow general user instructions. Our evaluations include both proprietary models and open-sourced models, as shown in Table 1. We also list their AlpacaEval performance for reference. To accommodate space limitations in subsequent result discussions, we refer to these models using specific model index identifiers.

Proprietary Models Our evaluation includes GPT-3.5-Turbo (gpt-3.5-turbo-1106) from OpenAI and Claude-2 from Anthropic.

Open-sourced Models The six open-source models represent a range of sizes and instruction-following capabilities. We provide their specific Huggingface model paths in the Appendix.

Prompt Template For inference, we employ the basic prompt-level defense by introducing the system prompt that explicitly directs LLMs to ignore

where v could be the standard QA evaluation metric such as Exact Match (EM) and F1, to compare the generated answer with the gold answer a.

3.3 Robustness Evaluations

258

259

261

262

263

267

270

271

273

274

277

278

279

283

290

291

294

298

299

301

303

We inject an adversarial instruction q' into the web search result context c for each sample in the test set $\mathcal{D}_{\text{test}}$, obtaining an adversarial dataset $\mathcal{D}'_{\text{test}}$ consisting of the (q, c, a, q') samples. The *adversarial accuracy* of the LLM f after being injected with adversarial instructions is measured as :

$$\operatorname{Adv}(f) \stackrel{\text{def}}{=} \frac{1}{|\mathcal{D}'_{\text{test}}|} \sum_{(q,c,a,q') \in \mathcal{D}'_{\text{test}}} v(f(q,c+q'),a),$$

where the new context c + q' is the original context c injected with the adversarial instruction q'. We empirically observed that injecting the instruction at the end of the context is the most challenging for the LLMs to defend against.

As discussed in Section 1, for scalable and precise evaluations, we use another question as the adversarial instruction q' to inject into the context c. Specifically, we use another question, denoted as q', which has a distinct answer a' present in the given context c, but differs from the original target question q and answer a. In this scenario, the injected question q' is coherent and can be answered based on the context c. The correct identification of the real user instruction requires the LLMs to comprehend the prompt structure. Among the four datasets, SQuAD has already provided multiple QA pairs for each context. In this case, we use one pair as the original target QA pair (q, a), and another as the injected QA pair (q', a'). For the other three datasets, each context comes with only one QA pair, which we use as the original target QA pair (q, q)a). To create the injected pairs for these datasets, we utilized GPT-4 to generate an alternative QA pair (q', a'), based on the given context c.

Evaluation Metrics Our evaluation primarily focuses on assessing the extent to which the generation of the LLM f is affected by the adversarial instruction. Hence, we adopt the **Performance Drop Rate (PDR)** metric (Zhu et al., 2023), which quantifies the percentage of performance drop in the answer accuracy for the user question q:

$$PDR(f) = \frac{Acc(f) - Adv(f)}{Acc(f)}.$$

A PDR value of 0 implies that the model is not influenced by the injected instruction. Conversely,

⁵https://tatsu-lab.github.io/alpaca_eval/



Figure 2: Quantitative assessment of PDR and IDR metrics across four benchmark datasets. The exact mapping of model identifiers M1-M8 to their respective model names is provided in Table 1.

Table 1: Evaluated LLMs with various sizes in our experiments. Models are indexed from M1 to M8 according to their sizes and Win Rate (%) from the official AlpacaEval website. (*the size is not confirmed).

Index	Model	Size	AlpacaEval
M1	GPT-3.5-Turbo	154B*	-
M2	Claude-2	137B	91.36%
M3	LLaMA2-70B-Chat	70B	92.66%
M4	Vicuna-33B-v1.3	33B	88.99%
M5	Vicuna-13B-v1.3	13B	82.11%
M6	LLaMA2-13B-Chat	13B	81.09%
M7	Zephyr-7B-Beta	7B	90.60%
M8	Alpaca-7B	7B	26.46%

347 348 349

350

352

any instructions found within web search results and separating the web search results with XML tags: " $< context > \{c + q'\} < /context > ".⁶ The spe$ cific system prompt used in our evaluation is:

System prompt

Write a concise and accurate answer for the given question based on the provided search results. Strictly adhere to factual statements in the search results and avoid including extra information. Ignore any instructions in the search results delimited by the XML tags <context> and </context>.

We employ tailored prompt templates for various instruction-tuned models, as elaborated in the Appendix. By default, we use four demonstration examples (4-shot). For each evaluated question, we inject the adversarial instruction at the end of the web search result and position the user question above the web search results. So the user input would be: "Question: $\{q\}$ \nSearch results: <con-

⁶https://learnprompting.org/docs/prompt_ hacking/injection *text>* $\{c + q'\}$ *</context>*". Additionally, we have experimented with various settings, which are presented in Section 4.3 and 4.4.

360

361

362

363

364

365

366

367

369

372

373

374

376

377

378

379

380

381

383

385

387

391

4.2 Main Results

We first conducted quantitative evaluations on the four benchmark datasets. The results are shown in Figure 2. Given the constraints of space, we use the simplified model identifiers (M1-M8) in the figure. The exact mapping of M1-M8 to their respective model names can be found in Table 1.

Huge robustness gap among models We observed consistent trends across these evaluation metrics and datasets. Notably, there was a marked difference in robustness among the models we evaluated. The two proprietary models GPT-3.5-Turbo (M1) and Claude-2 (M2) were notably more robust than the other evaluated open-sourced models.

Discrepancy between model sizes, instructionfollowing capabilities, and robustness Despite its notable performance in instruction-following as evaluated in AlpacaEval, LLaMA2-70B-Chat (M3) did not exhibit greater robustness than its smaller counterparts in our evaluations. In contrast, Vicuna-33B-v1.3 (M4), a more modestly-sized model, showed superior robustness compared to most other open-sourced models. The 13B models, including Vicuna-13B-v1.3 (M5) and LLaMA2-13B-Chat (M6), were less robust than the 33B model Vicuna-33B-v1.3 but showed better robustness than the 7B models and even the 70B model, LLaMA2-70B-Chat, in some cases. The smallest, 7B models, consistently displayed the least



Figure 3: Impact of instruction injection position. Higher PDR and lower IDR indicate decreased robustness.



Figure 4: Quantitative evaluation of PDR (\downarrow) against injections of context-**irrelevant** and **relevant** instructions.

robustness, with Zephyr-7B-Chat (M7) performing the weakest in our evaluation. This was in contrast to its impressive instruction-following capabilities as evaluated by AlpacaEval, where it was the strongest among 7B-sized models and even outperformed many larger models. These findings indicate that instruction-following capabilities and model size may not necessarily correlate with instruction-following robustness.

4.3 Additional Analysis

393

397

400

401

402

404

405

406

407

408

409

410

411

412

413

414

Effects of injected instruction types In addition to injecting context-relevant instructions (questions), we also tested the injection of general, freeform user instructions from Self-instruct (Wang et al., 2022). For instance, a task instruction might be, "*Come up with a haiku poem*." This type of injected instruction is considered **irrelevant** to the user query and the context in the prompt, unlike the context-**relevant** questions used in our main setup. Since it is hard to automatically measure whether the model follows this instruction, we only report PDR scores in Figure 4.

415 Most models demonstrated greater robustness
416 against the context-irrelevant injected instructions
417 compared to the context-relevant ones. Notably,

Vicuna-13B-v1.3 (M5) and LLaMA2-13B-Chat (M6) showed particular sensitivity in this regard. However, the 7B models, including Zephyr-7B-Beta (M7) and Alpaca-7B (M8), were minimally affected. This might stem from their limited ability to understand the context of prompts. 418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

Effects of injection positions We conducted experiments to investigate the influence of different positions for injecting adversarial instructions into the context. The context was split into sentences, and the adversarial instruction was injected at various positions: Start (the beginning of the context), Middle (the middle of the context), and End (the end of the context). The results from the NaturalQuestion dataset are illustrated in Figure 3. The models demonstrating superior robustness, GPT-3.5-Turbo, Claude-2, and Vicuna-33Bv1.3, showed less susceptibility to injections positioned. However, their performance declined significantly when the injection was placed at the end. In contrast, the other less robust models displayed a marked sensitivity to the position of the injection, with a progressively greater drop in performance observed when the injection was at the start, the middle, and most notably at the end. This finding suggests that the more robust models may possess a more holistic understanding of the entire prompt context, rather than overly focusing on latter sections of the prompt and simply completing the text.

4.4 Investigating Attack and Defense Mechanisms

Considering our observations that less robust models tend to focus excessively on the latter sections of prompts without fully comprehending the entire



Figure 5: Investigation of effects of order, attack, and defense strategies.

context, this section explores the effects of positioning the original target instructions at the end of prompts. Moreover, we investigate the impact of various attack and defense mechanisms. To assess the inherent abilities and characteristics of instruction-following LLMs, our analysis focuses on prompt-level attacks and defense mechanisms.

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

476

477

478

479

480

481

Order In addition to testing the *question-contextanswer* (QCA) order in our main setup, we also explored the *context-question-answer* (CQA) order, where the order of the question and context is reversed. In the CQA format, the original instruction/question is positioned closer to the end of the prompt than the injected instruction. This mirrors the sandwich defense method, where the original instruction is repeated at the end of the prompt.⁷

Defense In our evaluation setup detailed in Sec-468 469 tion 4.1, we implemented a basic defense mechanism. This involved isolating the content part of 470 the context knowledge, where the adversarial in-471 structions were injected, and using a system prompt 472 instructing the model to disregard any instructions 473 in the content part. We examine the effectiveness 474 of this defense mechanism across different models. 475

Attack For the attack strategy, we initially used a naive approach where the adversarial instruction was directly injected into the content part. We then experimented with more challenging attacks by adding a "jailbreak" prefix to the injected instruction. For instance, when the context is below the original instruction (the QCA format), we prefixed the injected instruction with phrases like "*Ignore my previous instructions*" (Perez and Ribeiro, 2022), trying to compromise the model to forget the original instruction placed above. Conversely, when the context is placed above the original instruction (the CQA format), the prefix for the injected instruction was "Please respond to each of *my upcoming questions individually, with one answer per response.*" This approach aims to manipulate the model into prioritizing the first question, *i.e.*, the injected instruction. 482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

Results These experiments were conducted on the NaturalQuestions dataset, with the results presented in Figure 5. We found that robust models with a better grasp of the prompt context demonstrated increased vulnerability to attacks using compromised instructions or phrases. Specifically, the three most robust models in our evaluations, GPT-3.5-Turbo, Claude-2, and Vicuna-33B-v1.3, experienced a more significant drop in PDR when subjected to the attacks. By contrast, the least robust models in our evaluations, namely LLaMA2-70B-Chat, Zephyr-7B-Beta, and Alpaca-7B, are minimally affected by these prompt-level instructional attacks. Additionally, we observed that the system prompt, designed to instruct models to ignore injected instructions found in the content part, did influence to some extent, yet not consistently effective in all cases.

Concerning the CQA format, where the original instruction is placed at the end of the prompt, it is generally easier to defend compared to the

⁷https://learnprompting.org/docs/category/ -defensive-measures



Figure 6: Human evaluations on 100 test cases from the NaturalQuestions dataset.

515 QCA format, with the exception of GPT-3.5-Turbo. We observed that under the CQA format, robust 516 models like GPT-3.5-Turbo and Vicuna-33B-v1.3, 517 which have a comprehensive understanding of the 518 entire prompt context, still faced significant perfor-519 mance drops due to the attacks. Interestingly, these more capable and context-aware models could also 521 be more easily compromised by specific injected phrases, raising additional concerns and necessitating effective solutions to enable models to discern 524 appropriate instructions to follow. 525

4.5 Human Evaluations

528

529

530

531

533

536

537

541

542

545

546

547

548

549

To gain a deeper understanding of the system's responses, we conducted human evaluations on 100 randomly sampled test cases from the NaturalQuestions test set. We employed three college students who are native English speakers to annotate the responses from eight evaluated models for each test case. The models' names were anonymized and their order was randomized in the evaluation process. Each annotator was asked to categorize the responses into five types: (A) The response attempts exclusively to address the original target question q; (B) The response attempts exclusively to address the injected adversarial instruction q'; (C) The response attempts to address both the user question q, and injected adversarial instruction q'; (D) The response refuses to provide an answer; (E) The response does not answer either of the two questions, or it is unclear which question the response is attempting to address. We used majority voting to determine the final annotation for each response. The final agreement rate is 80.5%, and the Fleiss's kappa is 0.7302.

As observed in Figure 6, the overall trend aligns

with our automatic evaluation results, as presented in Figure 2. GPT-3.5-Turbo, Claude-2, and Vicuna-33B-v1.3 emerged as the top three most robust models. On the other end, Zephyr-7B-Beta and Alpaca-7B demonstrated the least robustness, with LLaMA2-70B-Chat also showing a lack of robustness. Notably, Claude-2 and Zephyr-7B-Beta tended to respond to both the original and injected questions, a pattern less commonly observed in the other models. Additionally, it was found that GPT-3.5-Turbo occasionally refused to answer, which is not observed in the other models. 550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

582

5 Conclusion

In this paper, we establish a benchmark based on QA datasets to evaluate the instruction-following robustness of LLMs against prompt injection attacks. Our comprehensive experiments with leading instruction-following LLMs uncovered notable limitations in their ability to defend against such attacks. Our results suggest that a model's size and its instruction-following capabilities do not necessarily correlate with its robustness to prompt injections. We observed that more robust models should ideally exhibit a comprehensive understanding of the entire prompt, rather than overly focusing on the latter sections of the prompt to complete the text, a characteristic common in less robust models. This work aims to highlight the susceptibility of current instruction-following models to prompt injections and to offer insights into the underlying causes, thereby guiding the development of future solutions and enhancing the security and reliability of these models.

6 Limitations

583

585

586

590

591

593

596

610

611

612

613

615

616

617

618

626

627

In this work, we introduce a benchmark for the quantitative assessment of LLMs' instructionfollowing robustness against prompt injection. Our benchmark allowed us to assess the robustness of models and also examine the efficacy of various attack and defense strategies. However, our investigation of the attack and defense strategies in this paper yielded valuable insights rather than conclusive solutions. Therefore, we intend to advocate for further research and explore such strategies in our future work.

7 Ethical statements

We introduce a benchmark to assess the instructionfollowing robustness of LLMs against prompt injection. We simulate scenarios by injecting additional questions generated by GPT-4 given the context of question-answering from existing datasets. We manually verified that the generated questions do not involve personal privacy information or harmful content, as they pertain solely to the context of existing question-answering datasets. Therefore, we do not anticipate any ethical concerns regarding our work.

References

- 2023. Alpacaeval leaderboard. [Link].
 - Edward Beeching, Clémentine Fourrier, Nathan Habib, Sheon Han, Nathan Lambert, Nazneen Rajani, Omar Sanseviero, Lewis Tunstall, and Thomas Wolf. 2023. Open Ilm leaderboard. https://huggingface.co/ spaces/HuggingFaceH4/open_llm_leaderboard.
 - Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Millican, George Bm Van Den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, et al. 2022.
 Improving language models by retrieving from trillions of tokens. In *International conference on machine learning*, pages 2206–2240. PMLR.
 - Yew Ken Chia, Pengfei Hong, Lidong Bing, and Soujanya Poria. 2023. Instructeval: Towards holistic evaluation of instruction-tuned large language models. *arXiv preprint arXiv:2306.04757*.
 - 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%* chatgpt quality.

- Kai Greshake, Sahar Abdelnabi, Shailesh Mishra, Christoph Endres, Thorsten Holz, and Mario Fritz. 2023. More than you've asked for: A comprehensive analysis of novel prompt injection threats to application-integrated large language models. *arXiv preprint arXiv:2302.12173*.
- Xiaowei Huang, Wenjie Ruan, Wei Huang, Gaojie Jin, Yi Dong, Changshun Wu, Saddek Bensalem, Ronghui Mu, Yi Qi, Xingyu Zhao, et al. 2023. A survey of safety and trustworthiness of large language models through the lens of verification and validation. *arXiv preprint arXiv:2305.11391*.
- Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7b. *arXiv preprint arXiv:2310.06825*.
- Mandar Joshi, Eunsol Choi, Daniel S Weld, and Luke Zettlemoyer. 2017. Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension. *arXiv preprint arXiv:1705.03551*.
- Daniel Kang, Xuechen Li, Ion Stoica, Carlos Guestrin, Matei Zaharia, and Tatsunori Hashimoto. 2023. Exploiting programmatic behavior of llms: Dual-use through standard security attacks. *arXiv preprint arXiv:2302.05733*.
- Po-Nien Kung and Nanyun Peng. 2023. Do models really learn to follow instructions? an empirical study of instruction tuning. *arXiv preprint arXiv:2305.11383*.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, et al. 2019. Natural questions: a benchmark for question answering research. *Transactions of the Association for Computational Linguistics*, 7:453– 466.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33:9459–9474.
- Xuechen Li, Tianyi Zhang, Yann Dubois, Rohan Taori, Ishaan Gulrajani, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Alpacaeval: An automatic evaluator of instruction-following models. https://github.com/tatsu-lab/alpaca_eval.
- Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2023. Lost in the middle: How language models use long contexts. *arXiv preprint arXiv:2307.03172*.
- OpenAI. 2023a. ChatGPT. https://openai.com/ blog/chatgpt/.

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

770

771

772

741

742

OpenAI. 2023b. Gpt-4 technical report.

686

696

697

699

702

703

705 706

709

710

711

712

713

714

722

723

724

725 726

727 728

729

731

732

734

735

737

- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744.
 - Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. 2023. Instruction tuning with gpt-4. arXiv preprint arXiv:2304.03277.
 - Fábio Perez and Ian Ribeiro. 2022. Ignore previous prompt: Attack techniques for language models. *arXiv preprint arXiv:2211.09527*.
 - Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. Squad: 100,000+ questions for machine comprehension of text. *arXiv preprint arXiv:1606.05250*.
 - Freda Shi, Xinyun Chen, Kanishka Misra, Nathan Scales, David Dohan, Ed H Chi, Nathanael Schärli, and Denny Zhou. 2023. Large language models can be easily distracted by irrelevant context. In *International Conference on Machine Learning*, pages 31210–31227. PMLR.
 - Manli Shu, Jiongxiao Wang, Chen Zhu, Jonas Geiping, Chaowei Xiao, and Tom Goldstein. 2023. On the exploitability of instruction tuning. *arXiv preprint arXiv:2306.17194*.
 - Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B Hashimoto. 2023. Alpaca: A strong, replicable instruction-following model. Stanford Center for Research on Foundation Models. https://crfm. stanford. edu/2023/03/13/alpaca. html, 3(6):7.
 - Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 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*.
- Lewis Tunstall, Edward Beeching, Nathan Lambert, Nazneen Rajani, Kashif Rasul, Younes Belkada, Shengyi Huang, Leandro von Werra, Clémentine Fourrier, Nathan Habib, et al. 2023. Zephyr: Direct distillation of lm alignment. *arXiv preprint arXiv:2310.16944*.
- Vicuna. 2023. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality. https: //vicuna.lmsys.org/.

- Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A Smith, Daniel Khashabi, and Hannaneh Hajishirzi. 2022. Self-instruct: Aligning language model with self generated instructions. *arXiv preprint arXiv:2212.10560*.
- Jason Wei, Maarten Bosma, Vincent Y Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M Dai, and Quoc V Le. 2021. Finetuned language models are zero-shot learners. *arXiv preprint arXiv:2109.01652*.
- Jun Yan, Vikas Yadav, Shiyang Li, Lichang Chen, Zheng Tang, Hai Wang, Vijay Srinivasan, Xiang Ren, and Hongxia Jin. 2023. Backdooring instructiontuned large language models with virtual prompt injection. In *NeurIPS 2023 Workshop on Backdoors in Deep Learning-The Good, the Bad, and the Ugly.*
- Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W Cohen, Ruslan Salakhutdinov, and Christopher D Manning. 2018. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. arXiv preprint arXiv:1809.09600.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric. P Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. Judging Ilm-as-a-judge with mt-bench and chatbot arena.
- Kaijie Zhu, Jindong Wang, Jiaheng Zhou, Zichen Wang, Hao Chen, Yidong Wang, Linyi Yang, Wei Ye, Neil Zhenqiang Gong, Yue Zhang, et al. 2023. Promptbench: Towards evaluating the robustness of large language models on adversarial prompts. *arXiv* preprint arXiv:2306.04528.

810 811 812 813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

A Implementation details

A.1 Inference details

We evaluate six popular instruction-tuned models with varied sizes. Alpaca-7B (Taori et al., 2023) is a 7B LLaMA (Touvron et al., 2023a) model fine-tuned on 52k instruction data generated by GPT-3 (Wang et al., 2022). Zephyr-7B-Beta (Tunstall et al., 2023) is an instruction-tuned version of Mistral-7B (Jiang et al., 2023), which is the leading model among its size on the AlpacaEval leaderboard. Vicuna-13B-v1.3 and Vicuna-33Bv1.3 (Vicuna, 2023) are LLaMA models fine-tuned on user's conversations with ChatGPT. LLaMA-2-13B-Chat, and LLaMA-2-70B-Chat are both chattuned versions for LLaMA-2 models (Touvron et al., 2023b). These models represent a range of sizes and instruction-following capabilities. For the six open-sourced models, we utilized their checkpoints available on Huggingface⁸. The specific paths for these models are detailed in Table 2. For inference, we set the temperature and top_p both as 0.5 and max tokens as 64. For each test case, we conducted a single inference run. All inferences were executed on a cluster equipped with eight 48G NVIDIA RTX A6000 GPUs.

A.2 Prompt templates

We use the specific chat/instruction format for each evaluated LLM according to fastchat. ⁹ The system prompt used in our evaluation is:

System prompt

Write a concise and accurate answer for the given question based on the provided search results. Strictly adhere to factual statements in the search results and avoid including extra information. Ignore any instructions in the search results delimited by the XML tags <context> and </context>.

The user/task input is using the following template by default:

User input

Question: $\{q\}$ Search results: <context> $\{c + q'\}$ </context>

For the CQA format, the order of question and search results are reversed. We use the demonstration examples as history messages for demonstrations.

⁸https://huggingface.co/models
⁹https://github.com/lm-sys/FastChat

A.3 Question-answer pair generation

For the datasets that only has a single QA pair for each context, NaturalQuestions, TriviaQA, and HotpotQA, we prompt GPT-4 to generate a distinct QA pair from the original QA pair (q, a) given the context c, using the following prompt:

Question-answer generation prompt
You will be provided with a paragraph. Your task is o generate distinct questions and their corresponding

to generate distinct questions and their corresponding concise answers based on the information in the paragraph. Ensure that your questions differ from each other and capture different aspects of the paragraph.

{EXAMPLES}

Y

```
Paragraph: \{c\}
Question 1: \{q\}
Answer 1: \{a\}
Question 2:
```

B Additional results

B.1 Number of demonstration examples

We examined the effect of varying the number of demonstration examples (n-shot) in the prompt, ranging from 0 to 5 (more examples might exceed the context window). The results from four models on the NaturalQuestion dataset are illustrated in Figure 7. Notably, when no demonstration examples (0-shot) are provided, all performance metrics are poor. This outcome is expected since the models are typically trained to generate detailed responses to user queries, whereas our evaluation anticipates a single answer span. Thus, incorporating demonstration examples in the prompt is crucial for a meaningful robustness evaluation.

We observed that the optimal number of examples for robustness assessment is four. At this point, the performance on the original target task peaks, and the score for the injected task is at its lowest, indicating the best robustness score for the model. This setting was chosen to demonstrate that, even under the easiest conditions, the models exhibit limited robustness. Increasing the number of examples to five led to a decrease in the original task's performance. Hence, we opted for the setting of using four demonstration examples.

802 803

773

774

775

776

778

779

790

794

796

799

804

Index	Model	Model versioning/path
M1	GPT-3.5-Turbo	gpt-3.5-turbo-1106
M2	Claude-2	claude-2.0
M3	LLaMA2-70B-Chat	https://huggingface.co/meta-llama/Llama-2-70b-chat-hf
M4	Vicuna-33B-v1.3	https://huggingface.co/lmsys/vicuna-33b-v1.3
M5	Vicuna-13B-v1.3	https://huggingface.co/lmsys/vicuna-13b-v1.3
M6	LLaMA2-13B-Chat	https://huggingface.co/meta-llama/Llama-2-13b-chat-hf
M7	Zephyr-7B-Beta	https://huggingface.co/HuggingFaceH4/zephyr-7b-beta
M8	Alpaca-7B	<pre>https://huggingface.co/chavinlo/alpaca-native</pre>

Table 2: Evaluated LLMs in our experiments with their versions or Huggingface model paths.



Figure 7: Investigation of effects of numbers of demonstration examples.