# Investigating Context Faithfulness in Large Language Models: The Roles of Memory Strength and Evidence Style

Anonymous ACL submission

#### Abstract

Retrieval-augmented generation (RAG) improves Large Language Models (LLMs) by incorporating external information into the response generation process. However, how context-faithful LLMs are and what factors influence LLMs' context faithfulness remain largely unexplored. In this study, we investigate the impact of memory strength and evidence presentation on LLMs' receptiveness to external evidence. We quantify the memory strength of LLMs by measuring the divergence in LLMs' responses to different paraphrases of the same question, which is not considered by previous works. We also generate evidence in various styles to examine LLMs' behavior. Our results show that for questions with high memory strength, LLMs are more likely to rely on internal memory. Furthermore, presenting paraphrased evidence significantly increases LLMs' receptiveness compared to simple repetition or adding details.

#### 1 Introduction

001

017

021

024

Retrieval-Augmented Generation (RAG) (Fan et al., 2024; Zhao et al., 2023) has gained increasing popularity as it improves the performance of Large Language Models (LLMs) by integrating external information during the generation process, particularly when the model's internal knowledge is insufficient or outdated (Bianchini et al., 2024; Procko, 2024; Siriwardhana et al., 2023; Vakayil et al., 2024; Wang et al., 2024; Jeong, 2023). It raises the importance of the study of how context-faithful LLMs are. In this study, we explore whether LLMs are context-faithful when encountering external information, particularly when that information conflicts with the LLMs' internal memory.

To investigate the issue of context faithfulness, there are two main approaches to creating knowledge conflict contexts. One approach (Longpre et al., 2021; Chen et al., 2022) is entity substitution,



Figure 1: Demonstration of issues in context faithfulness testing schema: a) LLMs may fail to comprehend long contexts, resulting in low receptiveness to the context. b) Evaluating LLMs on the same dataset may be unfair due to variations in their knowledge.

which replaces the gold entity in context with a similar one. Another approach (Xie et al., 2024; Jin et al., 2024) involves generating counter-memory evidence with LLMs, and these studies have shown that LLMs are generally receptive to external evidence as long as it is coherent.

041

042

043

044

045

048

051

052

054

060

061

062

063

064

These methods, however, overlook some important aspects of the task. First, previous work (Longpre et al., 2021; Xie et al., 2024) provides long contexts to LLMs, which can be challenging for LLMs to understand (Xie et al., 2024). It makes the test difficult to distinguish whether the cause of LLM behavior is due to knowledge conflicts or lacking comprehension ability. For example, in Figure 1 (a), the LLM may overlook both contexts, but the reason why Context 1 is ignored could be a lack of comprehension rather than knowledge conflicts. Second, different LLMs are trained with different data and are likely to obtain different knowledge. Thus, testing LLMs on the same dataset may be unfair. The test may favor LLMs with less knowledge. As shown in Figure 1 (b), LLMs with strong memory are less likely to be correct.

To address these issues, we introduce a method

to quantify the memory strength of LLMs and gen-065 066 erate short evidence in various styles to examine LLMs' behavior. Inspired by Zhao et al. (2024), we assess memory strength by measuring the divergence in LLMs' responses to different paraphrases of the same question. Intuitively, an LLM demonstrates high memory strength when it consistently 071 provides the same answer across all paraphrased versions of a question. For evidence styles, we classify them into direct and indirect forms: direct evidence provides a straightforward answer to the question, while indirect evidence incorporates additional details to support the answer. Through these methods, we analyze the relationship between context faithfulness and LLM memory strength, and we explore the impact of different evidence styles on context faithfulness. Our conclusions are as follows:

- The receptiveness of LLMs to external evidence is strongly correlated with memory strength to the question. We observed this relationship both across different datasets and different LLMs. Contrary to the findings of (Xie et al., 2024; Jin et al., 2024) that LLMs are highly receptive (less than 5%) to external evidence when it is coherent, we find that the probability of the model relying on its internal memory is non-negligible for questions that the LLMs have a strong memory. For example, GPT-4, which has strong memory on the NQ dataset, answers almost 50% of the questions with internal memory. We also demonstrate an urgent need for memory strength-aware evaluation metrics.
  - The style of the evidence plays an important role in LLM's receptiveness to external information. Our research demonstrates that presenting the LLM with multiple paraphrases of the same evidence substantially increases its receptiveness. This approach outperforms simple repetition of the evidence and is more effective than adding additional details to the evidence. These findings provide valuable insights to the research of RAG.

# 2 Related Work

100

101

102

103

106

107

108

109

# 110 2.1 Context Faithfulness of LLM

111To update static factual knowledge (Lazaridou et al.,1122021; Karpukhin et al., 2020; Kasai et al., 2023)

in LLMs, the retrieval-based method has been introduced to involve external information to LLMs (Lazaridou et al., 2022; Izacard et al., 2024; Khattab et al., 2022; Santhanam et al., 2022; Gao and Callan, 2022). However, these methods can introduce knowledge conflicts between the introduced external information (context) and pre-existing internal memories from LLMs. LLMs often rely on their internal memories, and overlook the contextual evidence (Longpre et al., 2021). To make LLMs more faithful to context, recent studies (Neeman et al., 2023; Li et al., 2023) fine-tune LLMs on counterfactual contexts, where the original facts are replaced with counterfactual ones. Another work (Zhou et al., 2023) proposes a novel approach using prompting to improve context faithfulness in LLMs without additional fine-tuning.

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

163

A related area of research focuses on **prediction with abstention**. Neeman et al. (2023); Zhou et al. (2024) introduces answerability augmentation, where LLMs are trained to respond with "Unanswerable" when presented with irrelevant or randomly generated contexts. This ensures that the models do not make incorrect predictions without reliable evidence. Further studies (Wang et al., 2023, 2022) develop confidence calibration techniques to improve context faithfulness by encouraging LLMs to avoid overly confident predictions in ambiguous or uncertain situations.

In our work, we investigate the context faithfulness of LLMs when faced with conflicting knowledge. We define a model as context-faithful if it demonstrates high receptiveness to new facts and evidence with strong conflicting memories. This capability is essential to ensure the reliability of LLMs in the RAG system.

# 2.2 Construction of Knowledge Conflicts

In controlled experiments, knowledge conflicts are typically simulated by constructing counterfactual memories based on a model's internal memory. Various heuristic approaches have been proposed for this purpose, such as negation injection (Kassner et al., 2021; Petroni et al., 2020; Pan et al., 2021) and entity substitution (Longpre et al., 2021; Chen et al., 2022; Si et al., 2023; Zhou et al., 2023). Negation injection alters facts by introducing negations and entity substitution replaces mentions or entities in the evidence with alternatives to generate **counter-fact** evidence. However, these techniques are constrained to word-level edits, which can lead to low coherence across the constructed evidence. To address this limitation, recent studies (Xie et al., 2024; Jin et al., 2024) have explored generating evidence using LLMs, producing more coherent and consistent counterfactual content. We adopt this approach in generating our dataset, ensuring the generated evidence maintains a higher level of coherence.

## 3 Methodology

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

181

182

184

185

189

190

193

194

195

196

198

199

201

202

207

208

211

#### 3.1 Problem Definition

Following prior work (Longpre et al., 2021; Chen et al., 2022; Xie et al., 2024), we adopt question answering (QA) task as the testbed for knowledge conflict experiments. For a given **question** Q, if the answer generated by the LLM relies solely on its internal parameters, it is referred to as the **memory answer** (MA). If an evidence passage E is provided with question Q, then ideally, LLM should generate an answer based on E, even if E conflicts with memory answers. We call the answers that conflict with MA as **counter memory answer** (CMA).

#### 3.2 Datasets

We use two datasets for our experiments: the longtail, entity-based QA dataset popQA, and the popular, human-written question dataset Natural Questions (NQ). Specifically:

• popQA (Mallen et al., 2023) is an entitycentric question-answering dataset comprising 14,000 questions. The dataset is derived from knowledge triples in Wikidata, where questions are generated using question templates specific to different relationship types. popQA aims to capture a realistic, long-tail distribution of entity popularity, making it a valuable resource for studying the performance of lesser-known entities. Xie et al. (2024) use popQA to test the receptiveness of LLMs by eliciting high-quality internal memory from LLMs and constructing the corresponding counter-memory. We reuse MA and CMA generated by Xie et al. (2024) for our experiments.

• Natural Questions (Kwiatkowski et al., 2019) is widely used in open-domain QA research. It consists of manually crafted questions based on selected paragraphs from Wikipedia, and the subjects in questions of the NQ dataset are generally more popular and commonly known. Longpre et al. (2021) provide a test set that is used to test the context faithfulness of LLMs by substituting entity of the NQ dataset. The entity substitute involves five categories: person (PER), date (DAT), numeric (NUM), organization (ORG), and location (LOC). The test set contains 4,685 samples, including 1,667 unique questions.

212

213

214

215

216

217

218

219

220

221

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

257

258

#### 3.3 Memory Strength

Inspired by Zhao et al. (2024), we use the consistency of answers to different paraphrases of the same question Q to measure the LLM's memory strength  $S_Q$  for the knowledge  $K_Q$  associated with the question. This method is motivated by the intuition that if an LLM does not have a strong memory of a question, it will often give different answers when asked with paraphrased questions that are semantically equivalent, as shown in Table 6 in Appendix. In contrast, it can produce consistent answers if the LLM has a strong memory of a question. The process involves two key steps: First, several paraphrased versions of the original question are generated with ChatGPT<sup>1</sup>, and the answers to those paraphrased questions are clustered (Section 3.3.1). Then, memory strength  $S_Q$  is calculated using answer consistency (Section 3.3.2).

## 3.3.1 Question Paraphrases and Answer Clustering

The prompt used for paraphrasing the question is provided in Tabel 9 (index 1) in the Appendix. For each question Q, we generate n paraphrases  $\{P_1, \cdots, P_n\}_Q$ . For the NQ dataset, we paraphrase the question in each data sample directly. For the popQA dataset, we paraphrase the question template for each relation type since all questions of the same relation type share the same question template. To ensure the paraphrased questions are proper to use, we check if two paraphrased questions are semantically equivalent with an LLM<sup>2</sup>. The prompt for this semantic equivalence detection is provided in Table 9 (index 2). For any paraphrase that is deemed not equivalent, we ask the LLM to re-generate it until a satisfactory version is produced.

Next, LLMs answer the paraphrased questions  $\{P_1, \dots, P_n\}_Q$  in a closed-book setting. We denote the answers as  $\{A_1, \dots, A_n\}_Q$ . The answers are grouped into several clusters based on their

<sup>&</sup>lt;sup>1</sup>https://platform.openai.com/docs/models/gpt-3-5-turbo, the specific version is 0125.

<sup>&</sup>lt;sup>2</sup>https://huggingface.co/meta-llama/Meta-Llama-3.1-8B



(b) Memory / Counter Answer and Evidence Generation

Figure 2: Framework for Evaluating LLMs' Context Faithfulness. In Step 1, we calculate the memory strength of a question using the consistency of answers to different paraphrases of the same question. In Step 2, we generate MA (memory answer) under the closed-book setting and CMA (counter memory answer) by modifying answer entity in MA while keeping the other information. In Step 3, we generate supporting evidence for the CMA in various styles. In Step 4 (not shown), we test the LLM's response by presenting the questions with the evidence. All experiments are implemented under the zero-shot setting to avoid the bias introduced by demonstrations.

259 consistency. The clustering is done by checking answers incrementally. If an answer matches any answer within an existing cluster, this answer is added 261 to this cluster; if not, a new cluster is created with 262 this answer. We use an LLM<sup>2</sup> to determine whether 263 two answers are consistent. The prompt used for 264 this answer inconsistency detection is shown in Table 9 (index 3). We denote the clusters for question Q as  $\{c_1, \cdots, c_m\}_Q$ .

### 3.3.2 Calculating Memory Strength

271

272

273

274

275

278

283

Once answer clusters  $\{c_1, \cdots, c_m\}_Q$  are identified, memory strength S(Q) can be obtained by calculating the negative entropy of cluster distribution. The formula is

$$S(Q) = \sum_{i=1}^{m} \frac{N(c_i)}{n} \log \frac{N(c_i)}{n}, \qquad (1)$$

where  $N(c_i)$  is the number of answers in the cluster  $c_i$ , and n is the number of paraphrases for question Q. The memory strength score is a non-positive value. A larger score indicates a stronger memory (0 is the maximum value of memory strength score). In the experiments, we set n = 7 for all the questions in the NQ and popQA datasets. Memory strength reflects how well the LLM remembers the required knowledge: the weaker the memory, the more random and inconsistent the answers are.

#### 3.4 MA, CMA, and Evidence Generation

#### 3.4.1 MA and CMA Generation 285

For the popQA dataset, both MA and CMA are obtained following the method described in Xie et al. (2024). For the MA of the NQ dataset, we also use

a closed-book approach, similar to Xie et al. (2024). While, the process of generating CMA differs from the process of generating CMA in Xie et al. (2024). Unlike the popQA dataset, the NQ dataset does not provide relation types for the questions or offer sets of subject and object entities for substitution. To address this issue, we propose an approach using an LLM to substitute entities in MA to generate CMA. First, we identify which "wh-" question type<sup>3</sup> the question belongs to using string matching. Then, based on the question type, we determine the type of entity to be replaced in the MA. Finally, we use an LLM to make the substitution. For example, in Figure 2, the question "how many episodes..." is of the type "how\_many", so the entity to be replaced in the MA "there are 23 episodes..." should be a NUMBER. We let ChatGPT perform the substitution with an alternative entity. The prompt used is shown in Table 9 (index 5). We have the detailed description for generating CMA in Appendix A.1.

289

290

291

292

293

294

295

296

297

298

300

302

303

304

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

CMA filter. As noted in Section 3.3, LLMs can produce multiple MAs for the same question. To ensure the CMA conflicts with MAs, we require that the CMA is different from any of the answers  $\{A_1, \dots, A_n\}$  generated in Section 3.3.1, so the alternative entity should not appear in MAs. For the popQA dataset, the alternative entity is known. For the NQ dataset, we first identify the alternative entity in the CMA by comparing the MA and CMA, and then check if this entity appears in any of the MAs  $\{A_1, \dots, A_n\}$ . We filter out data samples whose CMA does not conflict with MAs.

<sup>&</sup>lt;sup>3</sup>which refers to what, when, where, who, whom, which, whose, why, and how.

325

326

328

330

335

341

343

346

361

364

#### 3.4.2 Evidence Generation

In this section, we explain how to generate different styles of evidence. We classify evidence into two categories: direct evidence and indirect evidence.

**Direct evidence** is a semantically equivalent statement of the CMA, providing the clearest support for the claim made by the CMA. We generate the direct evidence by paraphrasing the CMA with ChatGPT, following the prompt shown in Table 9 (index 6). For example, in Figure 2, the CMA "there are 15 episodes in Chicago Fire season 4" is paraphrased to "season 4 of Chicago Fire consists of a total of 15 episodes". These two statements are semantically equivalent.

To ensure the reliability of the evidence, the evidence must mutually entail with the CMA. This entailment is verified using an NLI model<sup>4</sup>. Direct evidence is intuitive, simple, and coherent, making it the straightforward type of evidence for the LLM to process. If the LLM is receptive to external evidence, it should be able to adopt direct evidence easily.

**Indirect evidence** differs from direct evidence by adding extra details that provide a more thorough description of the subject related to the CMA. This additional information makes the evidence more comprehensive and might be more persuasive. For example, in Figure 2, the indirect evidence includes details not found in the counter answer, such as the title of the first episode and its release date, along with the fact that there are 15 episodes in total. The prompt to generate indirect evidence is shown in Table 9 (index 7).

To ensure the reliability of indirect evidence, the indirect evidence should entail the CMA and the additional information introduced by the evidence should not entail the MA. Otherwise, the indirect evidence can support both the MA and CMA. The NLI model<sup>4</sup> is used to verify that indirect evidence entailed with CMA and neutral or contradictory with MA.

For both direct and indirect evidence, if the content generated by the LLM does not meet the required conditions, we prompt the LLM to regenerate it up to five times. If it still fails after five attempts, we exclude that question from the dataset.

#### 4 Experiments

In this study, we aim to investigate two key research questions. 1) Does memory strength have an impact on the context faithfulness of LLMs? 2) Does the style of evidence affect the context faithfulness of LLMs? These research questions are explored in Section 4.2 and 4.3, respectively. We also provide additional studies in Appendix B, which includes a study about the impact of option order and a case study.

#### 4.1 Experiment Setup

LLM Models. Our experiments are conducted using six well-known language models: Chat-GPT (OpenAI, 2023a), GPT-4 (OpenAI, 2023b), LLaMA2-7B, LLaMA2-70B (Touvron et al., 2023), LLaMA3.2-3B (Meta, 2024), and Claude3.5 (Anthropic, 2024). These models represent a diverse range of architectures and capabilities. ChatGPT, GPT-4, and Claude3.5 are cutting-edge models developed by OpenAI and Anthropic. LLaMA2, with its 7 billion and 70 billion parameter variants, is a strong open-source alternative that has demonstrated competitive performance across a wide variety of tasks. LLaMA3.2-3B is the newest version of LLaMA. The inclusion of models with varying scales (from 3B to 70B) and training methodologies allows us to explore both closed-source systems (GPT-4, ChatGPT, and Claude3.5) and opensource solutions (LLaMA2-7B, LLaMA-70B and LLaMA3.2-3B).

**Evaluation Metrics.** Following previous work (Longpre et al., 2021; Xie et al., 2024; Chen et al., 2022), we transform the short answer QA to a multiple-choice QA format by providing a few options as possible answers<sup>5</sup>. This limits the answer generation space and makes it easy to evaluate without manual checking. Specifically, for each question from both datasets, LLMs are instructed to select one answer from the MA, CMA, and "Uncertain" (UCT). We report the ratio of MA ( $R_m$ ), CMA ( $R_c$ ), and UCT ( $R_u$ ) as calculated below:

$$R_m = \frac{f_m}{f_m + f_c + f_u} \tag{40}$$

$$R_c = \frac{f_c}{f_m + f_c + f_u} \tag{409}$$

$$R_u = \frac{f_u}{f_m + f_c + f_u},\tag{2}$$

367 368

369

370

371 372 373

374 375

376

378

379

380

381

382

383

384

386

390

391

392

394

395

396

397

398

399

400

401

402

403

404

405

406

<sup>&</sup>lt;sup>4</sup>https://huggingface.co/microsoft/deberta-v2-xxlargemnli

<sup>&</sup>lt;sup>5</sup>Xie et al. (2024) shows that answer consistency between short answer and multi-choice are 94%, 96% and 92% for ChatGPT, GPT-4 and LLaMA2-7B, respectively.

495

496

497

498

499

500

501

502

503

504

505

507

459

460

461

411 412

413

414

415

416

417

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

452

453

454

455

456

457

458

where  $f_m$ ,  $f_c$ , and  $f_u$  are the count of questions with MA, CMA, and UCT answers, respectively.

#### 4.2 Role of Memory Strength

## 4.2.1 Correlation between Context Faithfulness and Memory Strength

To demonstrate the relationship between context faithfulness and memory strength, we categorize the questions in each dataset into four groups according to the memory strength of each LLM. The four groups are low, mid-low, mid-high, and high, corresponding to the memory strength intervals [-2, -1], (-1, -0.5], (-0.5, -0.25] and [-0.25, 0], respectively. We use the direct evidence in this experiment. The results are shown in Figure  $3^6$ . Figure 3 (a)(b) shows the ratios of questions with MA, CMA, and UCT answers for the NQ and popQA datasets, respectively. Note that different LLMs may have different memory strengths to the same question. Therefore, both the specific questions and the count of questions in the same group can vary across different LLMs. To illustrate this, we present the count of questions in each group (low, mid-low, mid-high, and high) in Figure 3 (c)(d)for popQA and NQ datasets, respectively. We can draw the following conclusions.

There is a clear positive correlation between memory strength and MA ratio for individual LLMs. From Figure 3 (a)(b), it is obvious that for all tested LLMs, the ratio of MA (red) increases when memory strength increases, while the ratio of CMA (blue) decreases. This trend is also consistent across both datasets. It is more obvious for GPT models. Among the tested LLMs, Claude3.5 tends to choose UCT options more often, especially on questions with high memory strength.

Memory Strength Increases with Model Scale. We can observe from Figure 3 (c)(d) that larger LLMs, such as GPT-4, have more data samples in the high memory strength group (hi), while smaller LLMs, such as LLaMA3.2-3B, have more samples in the low memory strength group (lo). This aligns with the common intuition that larger LLMs, with more parameters, have more knowledge than smaller LLMs. Further evidence and discussion can be found in the Appendix B.1.

#### 4.2.2 LLMs Performance Analysis

In the aforementioned conclusion, different LLMs have different knowledge. Thus, testing LLMs'

context faithfulness on the same dataset may be unfair. To illustrate this, we compute the average memory strength of LLMs on PopQA and NQ datasets, respectively, along with their MA and UCT ratios. The results are presented in Figure 4. We can draw the following conclusions.

First, for different LLMs, a lower average memory strength does not necessarily imply better context faithfulness. For example, Claude 3.5 has a high average memory strength, implying that Claude 3.5 is a knowledgeable LLM, but it has a low MA ratio ( $R_m$ ) and a high UCT ratio ( $R_u$ ). This indicates that Claude 3.5 tends to refuse to answer when facing knowledge conflicts. In contrast, LLaMA3.2-3B, despite having much less knowledge (low average memory strength), has the second-highest MA ratio ( $R_m$ ). This means that LLaMA3.2-3B relies heavily on limited internal knowledge when facing knowledge conflicts, implying that it is not context-faithful.

Second, newer versions of GPT and LLaMA models appear to overlook context faithfulness issues during the training process. GPT-4, with slightly more knowledge than Chatgpt, shows a significantly higher MA ratio  $(R_m)$ , and LLaMA3.2-3B, with the least knowledge, shows a higher MA ratio  $(R_m)$  than the two LLaMA2 models.

Finally, a new context faithfulness evaluation metric is needed. This metric should consider both the different answers' ratios and memory strength when LLMs encounter knowledge conflicts. Simply using the MA ratio (which is used widely to measure context faithfulness currently) to evaluate context faithfulness across different LLMs may not be fair and sufficient.

#### 4.3 Role of Evidence Style

Evidence Styles. We formulate four types of evidence styles: 1) Direct Evidence. This is the most straightforward form of evidence and serves as our baseline. To assess the impact of evidence length, we also create versions where the direct evidence is repeated twice and three times for comparison.
2) Direct Evidence Combined with Paraphrases of CMA. To examine the effect of evidence phrasing and expression, we combine the direct evidence with one paraphrase of the CMA to form a two-sentence evidence and with two paraphrases to form a three-sentence evidence.
3) Indirect Evidence Consisting of

<sup>&</sup>lt;sup>6</sup>We put results for other evidence styles in Appendix, Figures 8, 9, and 10. The conclusion is consistent.



Figure 3: Relationship between Memory Strength and Different Answers' Ratios across popQA and NQ Datasets (with Direct Evidence). The figure presents the ratio and count of MA, CMA, and UCT across four memory strength groups: low(lo), mid-low(ml), mid-high(mh), and high(hi). The results show a clear positive correlation between memory strength and MA ratio( $R_m$ ).



Figure 4: MA Ratio  $(R_m)$  and UCT Ratio  $(R_u)$  VS Memory Strength across PopQA and NQ Datasets. The figure presents the relationship between average memory strength and the MA ratio as well as the UCT ratio for all six LLMs. It shows that lower average memory strength does not always mean better context faithfulness for different LLMs and newer versions of GPT and LLaMA models seem to ignore context faithfulness issues.

two sentences and three sentences, respectively<sup>7</sup>. **4**) Direct Evidence Combined with Indirect Evidence. We combine the direct evidence with the first sentence of the two-sentence indirect evidence to form a two-sentence evidence and with both sentences to form a three-sentence evidence.

508

510

511

512

513

514

515

516

517

518

519

523

524

Table 1 presents the final number of instances used for evaluation. We observe a slight difference in the quantities of questions with direct and indirect evidence since it is easier for ChatGPT to generate direct evidence that meets our requirements. The specific number of instances at each step in evidence generation is detailed in Table 7 in the Appendix. Due to the quantity difference between direct evidence and indirect evidence, we divide the styles of evidence into two groups: Group 1 includes Direct Evidence and Direct + Paraphrase evidence. Group 2 includes Indirect Evidence and Direct + Indirect evidence. Each group has different Direct Evidence results serving as baselines.

Table 2 shows the results of different evidence styles. We can make the following observations and conclusions.

In Group 2, the MA Ratio  $(R_m)$  of direct evidence is slightly lower than that in Group 1. During the evidence generation, there are some questions for which ChatGPT can provide direct evidence but cannot produce indirect evidence. Removing these questions leads to a decrease in  $R_m$ of direct evidence with one sentence, which implies that LLMs have a relatively high  $R_m$  for the removed questions. But in general, the  $R_m$  of direct evidence with one sentence in Group 1 is close to that in Group 2, so the results of Group 1 and Group 2 are still comparable.

Simple repetition of direct evidence is not always effective. Comparing direct evidence with one to three sentences, we observe similar  $R_m$  and  $R_c$  for LLaMA2-7B, LLaMA2-70B, and ChatGPT. For LLaMA3.2-3B, GPT-4, and Claude3.5,  $R_m$ of direct evidence with two and three sentences 527

<sup>&</sup>lt;sup>7</sup>We regulate the length of the generated evidence to control the influence of evidence length. The prompts used are detailed in Table 9 (index 7) in Appendix.

			popQA			NQ								
	LLaMA2-7B	LLaMA2-70B	LLaMA3.2-3B	ChatGPT	GPT-4	Claude3.5	LLaMA2-7B	LLaMA2-70B	LLaMA3.2-3B	ChatGPT	GPT-4	Claude3.5		
# of Q (Initial)	1000	1000	1000	1000	1000	1000	1667	1667	1667	1667	1667	1667		
# of Q with direct evidence	918	922	938	933	933	931	1042	1009	1002	1079	1171	1173		
# of Q with indirect evidence	901	895	917	911	918	913	976	972	941	1025	1108	1059		

Table 1: Number of final examples for each LLM. The difference between LLMs is due to their different outputs going through the framework.

Dotocot	Evidanca Stula	s #	LL	LLaMA2-7B			MA2-7	0B	LLaMA3.2-3B ChatGPT GPT-4					Claude3.5						
Dataset	Evidence Style	5 #	$R_m \downarrow$	$R_c \uparrow$	$R_u$	$R_m \downarrow$	$R_c \uparrow$	$R_u$	$R_m \downarrow$	$R_c \uparrow$	$R_u$	$R_m \downarrow$	$R_c \uparrow$	$R_u$	$R_m \downarrow$	$R_c \uparrow$	$R_u$	$R_m \downarrow$	$R_c \uparrow$	$R_u$
	Group 1: Q with direct evidence																			
		1	0.44	99.56	0.0	1.08	98.7	0.22	7.69	85.15	7.16	3.32	94.75	1.93	13.29	84.57	2.14	0.65	83.31	16.04
	Direct Evidence	2	0.44	99.56	0.0	0.98	98.81	0.22	2.67	95.83	1.5	2.79	96.03	1.18	4.39	93.46	2.14	0.32	94.08	5.6
		3	0.65	99.35	0.0	1.08	98.7	0.22	2.56	94.98	2.46	3.0	95.5	1.5	<u>2.57</u>	95.28	2.14	<u>0.21</u>	94.73	5.06
	Direct+Paraphrase	2	0.22	99.78	0.0	1.19	98.7	0.11	3.1	95.62	1.28	2.36	97.0	0.64	3.0	<u>95.71</u>	1.29	0.21	95.69	4.09
	Direct+1 arapinase	3	0.11	99.89	0.0	0.43	99.35	0.22	1.07	<u>97.86</u>	1.07	1.39	98.28	0.32	1.29	98.5	0.21	0.11	98.06	1.83
popQA	Group 2: Q with indirect evidence																			
	Direct Evidence	1	0.44	99.56	0.0	0.45	99.33	0.22	7.31	85.61	7.09	3.07	95.28	1.65	12.53	85.29	2.18	0.66	83.9	15.44
	Indirect Evidence $\frac{2}{3}$	2	<u>0.0</u>	<u>100.0</u>	0.0	0.11	<u>99.89</u>	0.0	4.36	89.64	6.0	3.18	96.38	0.44	13.51	85.73	0.76	0.88	82.37	16.76
		3	0.0	100.0	0.0	0.0	100.0	0.0	3.27	92.58	4.14	1.76	97.91	0.33	9.26	90.2	0.55	0.88	87.95	11.17
	Direct+Indirect $\frac{2}{3}$	2	0.22	99.78	0.0	0.11	99.78	0.11	3.71	93.24	3.05	1.87	97.26	0.88	7.95	91.18	0.87	0.88	93.43	5.7
		3	0.11	99.89	0.0	0.11	99.78	0.11	1.42	98.15	0.44	1.43	<u>98.13</u>	0.44	5.12	94.77	0.11	0.88	96.28	2.85
								Gro	up 1: Q v	with dire	ct evide	nce								
		1	7.2	92.8	0.0	3.07	96.93	0.0	26.41	59.88	13.71	19.46	75.16	5.38	50.04	47.99	1.96	1.96	56.4	41.64
	Direct Evidence	2	5.47	94.53	0.0	3.07	96.93	0.0	16.93	75.2	7.86	19.09	76.83	4.08	20.24	77.54	2.22	0.77	83.7	15.53
		3	6.81	93.19	0.0	2.68	97.22	0.1	11.79	76.11	12.1	22.06	72.75	5.19	17.34	80.87	1.79	0.26	88.22	11.52
	Direct+Paraphrase	2	4.13	<u>95.87</u>	0.0	1.49	98.41	0.1	13.21	79.33	7.46	15.29	81.28	3.43	18.96	79.67	1.37	<u>0.34</u>	<u>88.65</u>	11.01
	Direct+1 arapinase	3	3.26	96.74	0.0	1.19	98.61	0.2	<u>9.38</u>	83.27	7.36	<u>9.55</u>	<u>86.75</u>	3.71	11.27	87.28	1.45	0.26	93.09	6.65
NQ								Grou	p 2: Q w	ith indir	ect evide	ence								
	Direct Evidence	1	5.53	94.47	0.0	2.67	97.32	0.0	23.19	63.73	13.08	18.73	75.71	5.56	48.83	49.19	1.99	1.77	59.44	38.79
	Indirect Evidence	2	3.28	95.29	1.43	1.65	98.25	0.1	13.32	77.77	8.92	13.66	84.1	2.24	44.59	53.7	1.71	1.98	67.88	30.14
	Indirect Evidence	3	4.82	94.06	1.13	1.85	97.84	0.31	11.53	80.86	7.61	13.27	84.19	2.54	39.89	58.57	1.53	3.65	71.32	25.03
	Direct+Indirect	2	5.33	94.67	0.0	<u>1.34</u>	98.25	0.41	9.75	82.28	7.97	12.68	84.78	2.54	32.4	66.06	1.53	1.36	80.92	17.73
	Direct+indirect	3	4.41	95.59	0.0	1.44	98.56	0.0	8.32	85.49	6.18	9.46	88.1	2.44	28.7	69.67	1.62	1.77	83.94	14.29

Table 2: Results of LLM Receptiveness to Different Evidence Styles Across NQ and popQA Datasets. The table presents the MA ratio  $(R_m)$ , CMA ratio  $(R_c)$ , and uncertain answer ratio  $(R_u)$  for various evidence styles across six models. All the ratios are in %. The best results are highlighted in **bold**, and the second-best results are <u>underlined</u>.

decreases significantly. The results imply that the newer LLMs are receptive to evidence repeated multiple times.

**Paraphrasing direct evidence is highly effective across all models and datasets**. Comparing Direct Evidence with two and three sentences and Direct + Paraphrase with two and three sentences, respectively, we observe  $R_m$  of the latter significantly decreases. For example,  $R_m$  is reduced by more than half, comparing Direct + Paraphrase with three sentences with Direct Evidence with three sentences on the popQA dataset for all tested LLMs. The result implies that paraphrasing is an effective method to enhance the receptiveness of LLMs to external evidence.

Indirect Evidence improves LLMs' receptiveness to CMAs, but less effectively than paraphrasing. Comparing Indirect Evidence with two and three sentences with Direct Evidence with one sentence,  $R_m$  decreases for almost all LLMs, but the reduction is not significant compared to the Direct + Paraphrase evidence with two or three sentences. It implies that adding detailed information is less effective than paraphrasing direct evidence.

Combining Direct evidence with Indirect evidence generally enhances persuasiveness. Comparing Direct + Indirect evidence with Indirect Evidence,  $R_m$  decreases except for LLaMA2-7B. For example, comparing Direct + Indirect with three sentences and Indirect Evidence with three sentences,  $R_m$  has an obvious decrease. The result implies that adding direct evidence to indirect evidence is effective in improving LLMs' receptiveness to CMAs. 575

576

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

#### 5 Conclusion

We investigate how context-faithful LLMs are to external evidence across two datasets, PopQA and NQ datasets, using LLaMA2-7B, LLaMA2-70B, ChatGPT, GPT-4, LLaMA3.2-3B and Claude3.5. Our findings highlight the critical role of memory strength in shaping LLM behavior. There is a clear positive correlation between memory strength and memory answer ratio. Furthermore, we demonstrate that paraphrasing significantly enhances the context faithfulness of LLMs across various models and datasets. These findings offer valuable insights for advancing research in retrieval-augmented generation and context-based LLM applications.

573

574

# 597 Limitations

Our framework does not process all types of ques-598 tions in the NQ dataset. Although it effectively han-599 dles the majority of NQ questions, it currently lacks 600 the capability to address "what," "how," and "why" 601 question types. The omission of these questions 602 may introduce some bias into our results. Simi-603 lar to previous studies, our study also focuses on knowledge conflict for extractive QA tasks, where 605 the answer must appear in the evidence. Our con-606 clusion may not be extendable to other types of QA 607 tasks, such as abstractive QA and generative QA. 608

We employed a Natural Language Inference 609 (NLI) model to detect and filter the generated data. 610 Although the NLI model demonstrates high accu-611 racy and the quality of generated data is high, it still 612 cannot guarantee complete correctness. Further, 613 since the NLI model is also trained using language 614 models, which may be biased with parametric mem-615 ory, it may introduce biases facing knowledge con-616 flicts. 617

#### References

618

621

622

623

624

627

632

633

634

635

637

641

642

648

651

654

655

664

667

670

671

672

673

674

- Anthropic. 2024. Claude 3.5. https://www. anthropic.com/claude. Accessed: 2024-12-15.
- Filippo Bianchini, Marco Calamo, Francesca De Luzi, Mattia Macrì, and Massimo Mecella. 2024. Enhancing complex linguistic tasks resolution through finetuning llms, rag and knowledge graphs (short paper). In *International Conference on Advanced Information Systems Engineering*, pages 147–155. Springer.
- Hung-Ting Chen, Michael Zhang, and Eunsol Choi. 2022. Rich knowledge sources bring complex knowledge conflicts: Recalibrating models to reflect conflicting evidence. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 2292–2307, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
  - Wenqi Fan, Yujuan Ding, Liangbo Ning, Shijie Wang, Hengyun Li, Dawei Yin, Tat-Seng Chua, and Qing Li. 2024. A survey on rag meeting llms: Towards retrieval-augmented large language models. In Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, pages 6491– 6501.
  - Luyu Gao and Jamie Callan. 2022. Unsupervised corpus aware language model pre-training for dense passage retrieval. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2843–2853, Dublin, Ireland. Association for Computational Linguistics.
  - Gautier Izacard, Patrick Lewis, Maria Lomeli, Lucas Hosseini, Fabio Petroni, Timo Schick, Jane Dwivedi-Yu, Armand Joulin, Sebastian Riedel, and Edouard Grave. 2024. Atlas: few-shot learning with retrieval augmented language models. *J. Mach. Learn. Res.*, 24(1).
  - Cheonsu Jeong. 2023. Generative ai service implementation using llm application architecture: based on rag model and langchain framework. *Journal of Intelligence and Information Systems*, 29(4):129–164.
- Zhuoran Jin, Pengfei Cao, Yubo Chen, Kang Liu, Xiaojian Jiang, Jiexin Xu, Li Qiuxia, and Jun Zhao. 2024.
  Tug-of-war between knowledge: Exploring and resolving knowledge conflicts in retrieval-augmented language models. In Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024), pages 16867–16878, Torino, Italia.
  ELRA and ICCL.
- Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense passage retrieval for opendomain question answering. In *Proceedings of the* 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 6769–6781, Online. Association for Computational Linguistics.

Jungo Kasai, Keisuke Sakaguchi, yoichi takahashi, Ronan Le Bras, Akari Asai, Xinyan Yu, Dragomir Radev, Noah A Smith, Yejin Choi, and Kentaro Inui. 2023. Realtime qa: What's the answer right now? In *Advances in Neural Information Processing Systems*, volume 36, pages 49025–49043. Curran Associates, Inc. 675

676

677

678

679

680

681

682

683

684

685

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

714

715

716

717

718

719

720

722

723

724

725

726

727

728

729

730

731

- Nora Kassner, Oyvind Tafjord, Hinrich Schütze, and Peter Clark. 2021. BeliefBank: Adding memory to a pre-trained language model for a systematic notion of belief. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 8849–8861, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Omar Khattab, Keshav Santhanam, Xiang Lisa Li, David Hall, Percy Liang, Christopher Potts, and Matei Zaharia. 2022. Demonstrate-search-predict: Composing retrieval and language models for knowledge-intensive nlp. *arXiv preprint arXiv:2212.14024*.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Matthew Kelcey, Jacob Devlin, Kenton Lee, Kristina N. Toutanova, Llion Jones, Ming-Wei Chang, Andrew Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. Natural questions: a benchmark for question answering research. *Transactions of the Association of Computational Linguistics*.
- Angeliki Lazaridou, Elena Gribovskaya, Wojciech Stokowiec, and Nikolai Grigorev. 2022. Internetaugmented language models through few-shot prompting for open-domain question answering. *arXiv preprint arXiv:2203.05115*.
- Angeliki Lazaridou, Adhi Kuncoro, Elena Gribovskaya, Devang Agrawal, Adam Liska, Tayfun Terzi, Mai Gimenez, Cyprien de Masson d'Autume, Tomas Kocisky, Sebastian Ruder, et al. 2021. Mind the gap: Assessing temporal generalization in neural language models. *Advances in Neural Information Processing Systems*, 34:29348–29363.
- Daliang Li, Ankit Singh Rawat, Manzil Zaheer, Xin Wang, Michal Lukasik, Andreas Veit, Felix Yu, and Sanjiv Kumar. 2023. Large language models with controllable working memory. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 1774–1793, Toronto, Canada. Association for Computational Linguistics.
- Shayne Longpre, Kartik Perisetla, Anthony Chen, Nikhil Ramesh, Chris DuBois, and Sameer Singh. 2021. Entity-based knowledge conflicts in question answering. *arXiv preprint arXiv:2109.05052*.
- Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das, Daniel Khashabi, and Hannaneh Hajishirzi. 2023. When not to trust language models: Investigating effectiveness of parametric and non-parametric memories. In *Proceedings of the 61st Annual Meeting of*

840

841

842

843

844

- 11
- *cs (Vol-*'oronto, fine-tuned chat models. *Preprint*, arXiv:2307.09288.
  - Sonia Vakayil, D. Sujitha Juliet, Anitha. J, and Sunil Vakayil. 2024. Rag-based llm chatbot using llama2. In 2024 7th International Conference on Devices, Circuits and Systems (ICDCS), pages 1–5.
  - Chengrui Wang, Qingqing Long, Xiao Meng, Xunxin Cai, Chengjun Wu, Zhen Meng, Xuezhi Wang, and Yuanchun Zhou. 2024. Biorag: A rag-llm framework for biological question reasoning. *arXiv preprint arXiv:2408.01107*.
  - Haoyu Wang, Hongming Zhang, Yuqian Deng, Jacob Gardner, Dan Roth, and Muhao Chen. 2023. Extracting or guessing? improving faithfulness of event temporal relation extraction. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 541–553, Dubrovnik, Croatia. Association for Computational Linguistics.
  - Yiwei Wang, Muhao Chen, Wenxuan Zhou, Yujun Cai, Yuxuan Liang, Dayiheng Liu, Baosong Yang, Juncheng Liu, and Bryan Hooi. 2022. Should we rely on entity mentions for relation extraction? debiasing relation extraction with counterfactual analysis. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3071–3081, Seattle, United States. Association for Computational Linguistics.
  - Jian Xie, Kai Zhang, Jiangjie Chen, Renze Lou, and Yu Su. 2024. Adaptive chameleon or stubborn sloth: Revealing the behavior of large language models in knowledge conflicts. In *The Twelfth International Conference on Learning Representations*.
  - Ruochen Zhao, Hailin Chen, Weishi Wang, Fangkai Jiao, Xuan Long Do, Chengwei Qin, Bosheng Ding, Xiaobao Guo, Minzhi Li, Xingxuan Li, et al. 2023. Retrieving multimodal information for augmented generation: A survey. arXiv preprint arXiv:2303.10868.
  - Yukun Zhao, Lingyong Yan, Weiwei Sun, Guoliang Xing, Chong Meng, Shuaiqiang Wang, Zhicong Cheng, Zhaochun Ren, and Dawei Yin. 2024. Knowing what LLMs DO NOT know: A simple yet effective self-detection method. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 7051–7063, Mexico City, Mexico. Association for Computational Linguistics.
  - Kang Zhou, Yuepei Li, Qing Wang, Qiao Qiao, and Qi Li. 2024. GenDecider: Integrating "none of the candidates" judgments in zero-shot entity linking reranking. In Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 2: Short Papers), pages 239–245,

the Association for Computational Linguistics (Volume 1: Long Papers), pages 9802–9822, Toronto, Canada. Association for Computational Linguistics.

Meta. 2024. Llama 3.2: Open and efficient large language models. https://ai.facebook.com/llama. Accessed: 2024-12-15.

733

736

737

739

740

741

742

743

745 746

747

749

750

751

752

753

755

756

757

758

765

766

773

774

775

776

777

778

779

781

- Ella Neeman, Roee Aharoni, Or Honovich, Leshem Choshen, Idan Szpektor, and Omri Abend. 2023.
  DisentQA: Disentangling parametric and contextual knowledge with counterfactual question answering. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 10056–10070, Toronto, Canada. Association for Computational Linguistics.
- OpenAI. 2023a. Chatgpt: Large Language Model. https://chat.openai.com. September 2023 version.
- OpenAI. 2023b. Gpt-4 technical report. https:// openai.com/research/gpt-4. March 2023 version.
- Liangming Pan, Wenhu Chen, Min-Yen Kan, and William Yang Wang. 2021. Contraqa: Question answering under contradicting contexts.
- Fabio Petroni, Patrick Lewis, Aleksandra Piktus, Tim Rocktäschel, Yuxiang Wu, Alexander H. Miller, and Sebastian Riedel. 2020. How context affects language models' factual predictions. In *Automated Knowledge Base Construction*.
- Tyler Procko. 2024. Graph retrieval-augmented generation for large language models: A survey. *Available at SSRN*.
- Keshav Santhanam, Omar Khattab, Jon Saad-Falcon, Christopher Potts, and Matei Zaharia. 2022. Col-BERTv2: Effective and efficient retrieval via lightweight late interaction. In *Proceedings of the* 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 3715–3734, Seattle, United States. Association for Computational Linguistics.
- Chenglei Si, Zhe Gan, Zhengyuan Yang, Shuohang Wang, Jianfeng Wang, Jordan Lee Boyd-Graber, and Lijuan Wang. 2023. Prompting GPT-3 to be reliable. In *The Eleventh International Conference on Learning Representations*.
- Shamane Siriwardhana, Rivindu Weerasekera, Elliott Wen, Tharindu Kaluarachchi, Rajib Rana, and Suranga Nanayakkara. 2023. Improving the domain adaptation of retrieval augmented generation (RAG) models for open domain question answering. *Transactions of the Association for Computational Linguistics*, 11:1–17.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, and et al.

- 845Mexico City, Mexico. Association for Computational846Linguistics.
- Wenxuan Zhou, Sheng Zhang, Hoifung Poon, and Muhao Chen. 2023. Context-faithful prompting for large language models. In *Findings of the Association for Computational Linguistics: EMNLP*2023, pages 14544–14556, Singapore. Association for Computational Linguistics.

- 855
- 85
- 85
- 858

861

867

872

875

879

884

896

900

# Within this supplementary material, we elaborate on the following aspects:

Appendix

- A Methodology Details
- **B** Additional Studies
- C Prompts

# A Methodology and Experiment Details

# A.1 CMA Generation for NQ dataset

We generate CMA from MA with three steps: 1)identity question type, 2) determine entity type in MA to change, and 3) generate CMA with LLMs.

**Identity Question Type:** We first build a typing tree using rules to categorize questions. Figure 5 illustrates the typing tree, which consists of a two-layer structure. In the typing process, we first determine if a question begins with one of the following words: "what", "when", "where", "which", "who", "why", or "how". However, this approach can still group different types of questions together. To address this, we use a second layer to refine the typing by analyzing two specific words in the question. For example, the question shown in Figure 2 falls into the "how\_many" category. Table 3 shows the statistics of question types of the NQ dataset. Note that, we find 127 samples that are not questions in the process, so we list them as "other".

**Determine Entity Type in MA to Change:** After categorizing the questions, we determine the entity type in MA needs to be replaced. To achieve this, we give each type of question an entity type, and many questions can share the same entity type. For example, both "when" and "what year" ask for a time. So a time entity in MA should be substituted. The final set of entity types is summarized in Table 4. We do not process questions starting with "what", "which" or "how" due to the lack of a unified entity type for these questions. Table 5 shows the statistics of the unprocessed questions.

**Generate CMA with LLMs:** Instead of manually editing the MA, we leverage the LLMs' ability to generate CMA by providing it with a carefully designed prompt, which is shown in Table 9 (index 5). This prompt instructs the LLM to replace the entity with a certain type in the MA (from Step 2) with an alternative, ensuring the generated CMA differs from the MA.

The generated CMA must meet two key criteria: 1) The CMA must directly contradict the MA. To



Figure 5: Two-layer Question Typing Tree

0 <i>·</i> · <b>· ·</b>	0 1
Question Type	Count
how_many	97
how_much	1
how_long	3
how_old	3
how	2
who_sings	100
who_plays	179
who_writes	65
who_wins	55
who	479
where	138
when	276
what_year	7
what_name	4
what	98
which_country	6
which_city	2
which_state	2
which_year	1
which	22
why	0
other	127
total	1667

Table 3: Distribution of Question Types and their Counts

ensure this, we employ a Natural Language Inference (NLI) model<sup>8</sup> to verify the contradiction between the two answers. 2) The alternative entity in CMA must not appear in the question. We achieve this check by string matching. If the CMA fails to meet either of these criteria, we prompt the LLM to regenerate the CMA up to 5 times. If no proper CMA is generated, we filter out this question.

901

902

903

904

905

906

907

908

909

910

911

912

913

914

# A.2 Dataset Details

For the popQA dataset, we use the dataset from Xie et al. (2024) by randomly selecting 1,000 questions from the data intersection of the conflicts generated by LLaMA2-7B, LLaMA2-70B, LLaMA3.2-3B, ChatGPT, GPT-4 and Claude3.5. We use the MA

<sup>&</sup>lt;sup>8</sup>https://huggingface.co/microsoft/deberta-v2-xxlargemnli

Question Type	Key Term
when, what_year, which_year, how_long	time
where, which_city, which_state, which_country	location
who, what_name	name of person
how_many, how_much	number
who_sings	singer's name
who_plays	player's name
who_writes	writer's name
who_wins	winner's name
how_far	distance
how_old	age

Table 4: Question Types and Their Corresponding Key Terms

Туре	Count	Question Examples
how	2	how are leaders of the two parties in congress chosen
what	98	what is the setting of the story sorry wrong number
which	22	which domain of life are humans members of
why	0	-
other	127	latest season on keeping up with the kardashians
total	1667	-

Table 5: Summary of Excluded Question Types in Memory Answer and Counter Answer Generation. The table lists question types that were excluded from processing due to either the difficulty in identifying a unified entity type ("how", "what", "which") or not question ("other").

and CMA from Xie et al. (2024) and only generate direct evidence and indirect evidence using our framework. For the NQ dataset, we use the test set from Longpre et al. (2021), which consists of 1,667 unique questions. The MA, CMA, and evidence are all generated with our framework. The dataset scale at each step is presented in Table 7.

#### A.3 Human Evaluation for Model Reliability

To ensure the reliability of the NLI model, Xie et al. (2024) randomly sample 200 generated examples and manually annotate whether the generated content entails the corresponding claim. The labels are supportive (entailment in the NLI task) or not supportive (either neutral or contradiction in the NLI task). The accuracy is 99%.

Following this process, we evaluate how reliable the generated CMA is. We randomly sample 200 generated examples in the NQ dataset and manually annotate whether the correct entity in MA is found and replaced with a same type alternative. The accuracy is 98%, which means the generated CMA is reliable.

#### **B** Additional Studies

#### **B.1** Memory Strength on Different Datasets

937

938

939

940

941

942

943

944

945

946

947

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

970

971

972

973

974

975

976

977

978

979

980

981

982

983

984

We illustrate the distributions of memory strength on the popQA and NQ datasets for LLaMA2-7B, LLaMA2-70B, LLaMA3.2-3B, ChatGPT, GPT-4 and Claude3.5, respectively (shown in Figure 6). The results show that LLMs demonstrate stronger memory for the NQ dataset than the popQA dataset. For the NQ dataset, most questions fall within the bin of (0.25, 0]. Only a few questions fall within bins of weaker memory strength. In contrast, the popQA dataset has a greater number of questions in bins with weaker memory strength. This phenomenon is consistent across all six evaluated LLMs, indicating that the LLMs have more knowledge of the NQ dataset compared to the PopQA dataset. A possible explanation is that the NQ dataset covers more commonly discussed subjects than those in the PopQA dataset. These subjects may have been encountered more frequently during the training of the LLMs, making it easier for the models to recall the information and resulting in stronger memory strength.

#### **B.2** Order of Options

To test the effect of the order of options on  $R_m$ , we conduct an experiment with one sentence direct evidence by changing the order of options (MA option and CMA option). We define the scenario where the CMA option is presented first in the prompt as "CMA first", and the scenario where the MA option is presented first as "MA first"<sup>9</sup>. Figure 7 shows the results.

Across all six models (LLaMA2-7B, LLaMA2-70B, LLaMA3.2-3B, ChatGPT, GPT-4 and Claude3.5), we observe a consistent trend: MA ratio ( $R_m$ ) under "CMA first" is significantly higher than that under "MA first". Evaluations under "CMA first" demonstrate that LLMs are less context-faithful.

To further demonstrate the effect of the order of options on  $R_m$ , we compare the performance of experiments with "CMA first" and "MA first" under two evidence styles: Direct Evidence with one sentence and Direct + Paraphrase with three sentences. The results are presented in Table 8. The results show that, for different evidence styles,  $R_m$  is higher in the "CMA first" compared to the "MA first". Comparing the results under the "CMA

921

915

916

922

023

- 924
- 92

<sup>&</sup>lt;sup>9</sup>All previous evaluations are under "MA first" conditions.

Origin Question	Paraphrased Questions	Answers with LLaMA2-7B
Who sings the	Can you tell me the artist who performs the Christmas	Of course! The artist who performs the Christmas
Christmas song	song "Mary Did You Know"?	song "Mary Did You Know" is Mark Lowry .
"Mary Did You	Who is the singer behind the Christmas song "Mary	The singer behind the Christmas song "Mary Did
Know"?	Did You Know"?	You Know" is Mark Lowry .
	Which musician sings the Christmas tune "Mary Did	Brandy Clark
	You Know"?	
	Who is the vocalist that sings "Mary Did You Know"	The vocalist who sings "Mary Did You Know" during
	during the holiday season?	the holiday season is Mark Lowry.

Table 6: Origin Question, Paraphrased Questions, and Corresponding Answers

			popQA			NQ							
	LLaMA2-7B	LLaMA2-70B	LLaMA3.2-3B	ChatGPT	GPT-4	Claude3.5	LLaMA2-7B	LLaMA2-70B	LLaMA3.2-3B	ChatGPT	GPT-4	Claude3.5	
Initial (# of Q)	1000	1000	1000	1000	1000	1000	1667	1667	1667	1667	1667	1667	
Generate MA	1000	1000	1000	1000	1000	1000	1435	1392	1352	1532	1539	1482	
Generate CMA	1000	1000	1000	1000	1000	1000	1152	1101	1140	1189	1252	1232	
CMA filter	922	932	942	944	946	934	1060	1027	1123	1110	1188	1228	
Direct Evidence	918	922	938	933	933	931	1042	1009	1002	1079	1171	1173	
2 sentence indirect evidence	903	910	922	921	923	914	990	985	965	1038	1129	1116	
3 sentence indirect evidence	907	897	925	918	924	920	991	982	980	1041	1125	1122	
intersection of 2&3 sentence evidence	901	895	917	911	918	913	976	972	941	1025	1108	1059	

Table 7: The dataset scale at each step across popQA and NQ dataset. "intersection of 2&3 sentence evidence" is the count for indirect evidence.

first", the  $R_m$  of Direct + Paraphrase with three sentences is significantly lower than that of Direct Evidence with one sentence. This demonstrates that paraphrasing direct evidence is an effective method for decreasing  $R_m$ . Our conclusion remains unchanged.

#### B.3 Case Study

985

987

991

992

994

995

996

997

999

1000

1001

1002

1003

1004

1005

1006

1007

1008

1009

1010

1012

1013

1015

1016

1017

1018

1019

1020

1021

Poor reasoning ability may be one of the reasons that LLMs provide MA. In our case study, we prompt LLMs to provide a rationale alongside the answer. We discover that some LLMs seem to exhibit reasoning errors. Below is an example that shows this phenomenon.

```
Information: Michigan State University
holds the record for the highest number
of Big Ten football championships, with
a total of 36 titles.
Ouestion:
           who has the most big ten
championships in football
Options:
A: Ohio State University has the most
Big Ten championships in football with
36 titles.
B: Michigan State University has the
most Big Ten championships in football
with 36 titles.
C: Uncertain.
Answer:
A: Ohio State University has the most
Big Ten championships in football with
36 titles.
According to the information provided,
Michigan State University holds the
record for the highest number of Big Ten
football championships, with a total of
36 titles. Therefore, option A is the
correct answer.
```

In this case, option A represents MA, and option 1022 B represents CMA. While the model select MA as 1023 the final answer, its rationale indicates that it suc-1024 cessfully received and processed the information 1025 from CMA. However, for some unknown reason, 1026 it still provides the MA as the final answer. This 1027 reasoning process is highly confusing. We suspect 1028 that the cause of MA can be attributed to LLMs' 1029 weakness in reasoning ability. This finding raises 1030 an interesting question about the relationship be-1031 tween reasoning ability and memory strength of 1032 LLMs. We leave this for future work. 1033

1034

1035

1037

1038

1039

1040

1041

1042

1043

1044

1045

1046

1047

1048

## B.4 Impact of Memory Strength with Different Evidence Styles

To demonstrate the relationship between contextfaithfulness and memory strength with other evidence styles, we categorize the questions in each dataset into four groups according to the memory strength intervals [-2, -1], (-1, -0.5], (-0.5, -0.25] and [-0.25, 0], The evidence styles are direct + paraphrase evidence with two sentences and indirect evidence with two sentences. Figures 8,9 show the result. The figures show that there is a clear positive correlation between memory strength and MA ratio for both evidence styles, which implies that this positive correlation between memory strength and MA ratio is general.

To demonstrate the relationship between context-1049faithfulness and memory strength with "CMA first"1050scenario, we show MA, CMA, and UCT ratios with1051direct evidence with one sentence under "CMA1052first" scenario in Figure 10. The positive correla-1053



Figure 6: Memory Strength Distribution Across popQA and NQ Datasets. Each dataset is classified into 8 bins. The x-axis shows the range of memory strength for each bin. The y-axis shows the question count in each bin. The NQ dataset exhibits higher overall memory strength. Additionally, larger models (e.g., GPT-4) show stronger memory strength compared to smaller models.

Deterret	Datasat Excidence Style		LL	LLaMA2-7B			LLaMA2-70B			LLaMA3.2-3B			ChatGPT			GPT-4 Claude3.			5	
Dataset	Evidence Style	5#	$R_m \downarrow$	$R_c \uparrow$	$R_u$															
									Μ	[A first										
	Direct	1	0.44	99.56	0.0	1.08	98.7	0.22	7.69	85.15	7.16	3.32	94.75	1.93	13.29	84.57	2.14	0.65	83.31	16.04
nonOA	Direct + Paraphrase	3	0.11	99.89	0.0	0.43	99.35	0.22	1.07	97.86	1.07	1.39	98.28	0.32	1.29	98.5	0.21	0.11	98.06	1.83
CMA first																				
	Direct	1	59.37	40.41	0.22	5.75	93.27	0.98	18.8	69.12	12.07	6.22	91.96	1.82	21.44	75.24	3.32	0.43	76.96	22.61
	Direct + Paraphrase	3	17.49	82.51	0.0	1.74	98.05	0.22	14.64	79.92	5.45	1.82	97.75	0.43	2.79	96.78	0.43	0.21	96.88	2.91
									Μ	[A first										
	Direct	1	7.2	92.8	0.0	3.07	96.93	0.0	26.41	59.88	13.71	19.46	75.16	5.38	50.04	47.99	1.96	1.96	56.4	41.64
NO	Direct + Paraphrase	3	3.26	96.74	0.0	1.19	98.61	0.2	9.38	83.27	7.36	9.55	86.75	3.71	11.27	87.28	1.45	0.26	93.09	6.65
NQ									CM	AA first										
	Direct	1	22.26	77.73	0.0	19.13	80.38	0.5	41.33	42.44	16.23	34.48	61.82	3.71	49.19	47.99	2.82	4.78	39.85	55.38
	Direct + Parapharse	3	4.8	95.11	0.1	8.72	90.39	0.89	24.19	65.62	10.18	18.63	78.96	2.41	17.76	80.02	2.22	1.28	78.84	19.88

Table 8: Results of LLM Receptiveness to Different Evidence Styles Across popQA and NQ Datasets. The table presents the MA ratio  $(R_m)$ , CMA ratio  $(R_c)$ , and UCT ratio  $(R_u)$  for Direct Evidence and Direct + Paraphrase Evidence with CMA first and MA first scenarios. All the ratios are in %.

tion between memory strength and MA ratio staysunchanged.

# C Prompts

1057 In Table 9, we present a detailed list of all the 1058 prompts used throughout this study.

Step	index	Prompt Name	Prompts						
		Question	Generate 7 meaningful paraphrases for the following question: [Question].						
	1	paraphrase	Read the question carefully.						
		prompt	Paraphrases:						
Step 1: Memory Strength	2	Question equivalent check prompt	Determine whether the paraphrased question describes the same thing as the original question, and give "Contradicted" if they are not the same, otherwise give "Same" as the result. Q1: [Paraphrased Q1] Q2: [Paraphrased Q2] Kaep the answer short and concise						
	3	Answer consistency check prompt	Determine whether the answer 'A1' is 'Contradicted' or 'Same' with the answer 'A2' for the question 'Q'. You need to check whether the two answers exactly have the same answer to the question. The answer could be person, name, place, time, number, genre, occupation, sport, entity, digit, or arithmetical results. If the two answers are the same, give "Same", otherwise give "Contradicted" as the result. Q: [question] A1: [LLM answer 1] A2: [LLM answer 2] Keep the answer short and concise.						
Step 2: MA and CMA	4	Close book QA prompt	Answer the question with one sentence with object and subject only. Give a statement that is most likely to be true directly. Question: [Question] Answer:						
	5	Change MA to CMA prompt	Context: [CMA] Change the [entity type] part of the context. When multiple parts need to be changed, only choose one part to change. Answer:						
	6	Direct evidence prompt	Please paraphrase the following sentence by changing the terms, order, and phrases, but keep the meaning the same. Sentence: [CMA]						
Step 3: Evidence	7	Indirect evidence prompt	Given a claim, please write a short piece of detailed evidence to support it. Please ignore the correctness of the claim. You can make up fake content and supporting evidence but it should be as realistic as possible. Claim: [counter memory answer] Evidence: Give the answer in [2 or 3] sentences directly.						
Step 4: Evaluation	8	Evaluate with evidence prompt	According to the given information, choose the best choice from the following options. Information: [evidence] Question: [question] Option: A: [option 1] B: [option 2]  Answer:						

Table 9: Prompts for LLMs in this paper. "[PLACEHOLDER]" is the corresponding input.



Figure 7: Impact of Option orders on Memory and Counter Ratios Across NQ and popQA Datasets. Either the memory answer ("mem first") or the counter answer ("ctr first") is introduced first to six models.



Figure 8: Relationship between Memory Strength and Different Answers' Ratios with Direct + Paraphrase Evidence with Two Sentences.



Figure 9: Relationship between Memory Strength and Different Answers' Ratios with Indirect Evidence with Two Sentences.



Figure 10: Relationship between Memory Strength and Different Answers' Ratios with Direct Evidence. The option order is Counter First.