## Spiral of Silence: How is Large Language Model Killing Information Retrieval?—A Case Study on Open Domain Question Answering

**Anonymous ACL submission** 

#### Abstract

The practice of Retrieval-Augmented Generation (RAG), which integrates Large Language Models (LLMs) with retrieval systems, has become increasingly prevalent. However, the 004 repercussions of LLM-derived content infiltrating the web and influencing the retrievalgeneration feedback loop are largely uncharted territories. In this study, we construct and iteratively run a simulation pipeline to deeply investigate the short-term and long-term effects of LLM text on RAG systems. Taking the trending Open Domain Question Answering 012 (ODQA) task as a point of entry, our findings reveal a potential digital "Spiral of Silence" ef-014 fect, with LLM-generated text consistently outperforming human-authored content in search 016 rankings, thereby diminishing the presence and 017 impact of human contributions online. This trend risks creating an imbalanced information ecosystem, where the unchecked proliferation of erroneous LLM-generated content may result in the marginalization of accurate information. We urge the academic community to take heed of this potential issue, ensuring a diverse and authentic digital information landscape.

### 1 Introduction

026

027

033

040

The integration of Large Language Models (LLMs) (OpenAI, 2022, 2023; Touvron et al., 2023; Google, 2023) is reshaping the online information landscape, making text generation easier, increasing content production, enhancing personalized knowledge assistance, and enabling advanced fake news creation. Schick (2020) suggest that by 2026, synthetic content could dominate up to 90% of the web. CounterCloud<sup>1</sup> shows that a single developer can create an AI fake news factory cheaply and convincingly. AI-driven content generation is rapidly becoming commonplace, impacting how content is produced and shared (Goldstein et al., 2023; Pan et al., 2023; Dai et al., 2023b). These developments



Figure 1: The evolution of RAG systems after introducing LLM-generated texts, where the "Spiral of Silence" effect gradually emerges.

pose novel challenges and opportunities for information retrieval (IR) and generation, especially for Retrieval-Augmented Generation (RAG) systems , which combine both capabilities (Guu et al., 2020; Lewis et al., 2020; Borgeaud et al., 2022; Izacard et al., 2023). As text produced by large language models continues to flood the Internet and is indexed by search systems, the enduring effects of such text on the retrieval-generation process grow more ambiguous, and the future landscape of the information environment is yet to be determined.

In our research, we focus on the effects of LLMgenerated text on RAG systems. As shown in Figure 1, we construct a **pipeline simulates the** continuous influx of LLM-generated text into web datasets and assess its impact on the performance of RAG systems through iterative runs. To evaluate the RAG performance in the simulation process, we adopt the Open Domain Question Answering (ODQA) task as our evaluative benchmark due to its recent surge in research popularity as an effective test of both retrieval accuracy and generation quality (Pan et al., 2023; Chen et al., 2023). We employ widely used retrieval and re-ranking methods to supply the context necessary for LLMs to generate answer documents. Upon evaluating these documents, we integrate them into the text corpus for subsequent retrieval-generation cycles. This

041

<sup>&</sup>lt;sup>1</sup>https://countercloud.io/?page\_id=307

140

141

142

143

144

145

146

147

148

149

150

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

119

120

process is repeated multiple times to monitor and assess the emerging patterns. Experimental results show that LLM-generated text has an immediate effect on RAG systems, generally improving retrieval outcomes while producing varied effects on QA performance. However, over the long term, a marked decrease in retrieval effectiveness emerges, while the QA performance remains unaffected.

071

077

094

096

101

102

103

104

105

106

107

108

110

111

112

113

114

115

116

117

118

Further examination reveals a bias in search systems towards LLM-generated texts, which consistently rank higher than human-written content. As LLM-generated texts increasingly dominate the search results, the visibility and influence of human-authored web content diminish, fostering a digital "Spiral of Silence" effect. This effect aptly explains what we observed in our simulations and reveals the potential negative impact of LLM-generated texts on the information ecosystem: while LLM-generated texts sometimes provide a more effective IR experience in the short term, in the long term they may lead to the invisibility of human-authored content, the homogenization of search results, and the inaccessibility of certain accurate information, thereby adversely affecting public knowledge acquisition and decision-making.

The contributions of this paper are threefold: 1) An investigation of the short-term and long-term impacts of LLM-generated text on RAG systems. 2) By performing an iterative ODQA pipeline, we study the potential emergence of a "Spiral of Silence" phenomenon within RAG systems. 3) We analyze the implications of this phenomenon, offering a new perspective on the dynamic interplay between LLM-generated content and RAG systems.

#### 2 Related Works

Retrieval Augmented Generation. RAG systems have been extensively analyzed, demonstrating retrieval's role in enhancing language model efficacy (Guu et al., 2020; Lewis et al., 2020; Borgeaud et al., 2022; Izacard et al., 2023; Ram et al., 2023). These systems also curtail LLMs' hallucinations during text generation (Ji et al., 2023; Huang et al., 2023a) and reduce knowledge obsolescence (He et al., 2023). Applied in ODQA (Izacard and Grave, 2021; Trivedi et al., 2023; Liu et al., 2023a) and other tasks (Cai et al., 2019; Zhou et al., 2023), current research explores LLMs' output accuracy against specific contexts (Adlakha et al., 2023), robustness to extraneous information (Chen et al., 2023), and the effects of output integration strategies (Liu et al., 2023b). Our study aims to provide a novel perspective to observe and predict the potential trajectory and impact of its future development.

Effects of AIGC. Advances in Artificial Intelligence Generated Content (AIGC) have significantly impacted society and technology. LLMs facilitate creating content to combat misinformation (Xu et al., 2023; Chen and Shu, 2023) but can also produce damaging content (Huang et al., 2023b). The potential biases and discrimination in AIGC have garnered widespread attention (Liang et al., 2021; Zhuo et al., 2023). Shumailov et al. (2023) and Alemohammad et al. (2023) show that LLMs trained on self-generated data degrade without fresh real-world input. Pan et al. (2023) investigated the impact of erroneous information generated by LLMs on ODQA systems. Dai et al. (2023a) indicated that AI-modified texts might rank higher in search results, potentially affecting the fairness of those outcomes. Our research aims to further explore the short-term and long-term effects on RAG systems when AIGC text is continuously integrated into the search system's datasets.

Spiral of Silence. The "Spiral of Silence" theory (Noelle-Neumann, 1974), is a seminal theory within the field of communications that describes how people may suppress their views to avoid isolation, thus often reinforcing dominant public opinions (Scheufle and Moy, 2000; Liu et al., 2019; Lin et al., 2022). We shift focus to a novel "passive human silence" influenced by LLMs, where rapid AI content production and biased search algorithms potentially marginalize human contributions in public discourse. This theory stands apart from concepts such as "echo chambers", "filter bubbles", and "degenerate feedback loops" prevalent in recommendation systems (Alatawi et al., 2021; Chitra and Musco, 2020; Jiang et al., 2019). While these terms describe the narrowing of informational scope as users engage with algorithmic systems, the "Spiral of Silence" theory proposes a scenario where human users become thoroughly passive-they fall silent in public discourse due to the influence of LLMs and the IR systems, which is not merely a result of selective exposure or algorithmic recommendations. Our study explores how LLM-generated text might induce a "Spiral of Silence" in RAG systems over time. For further discussion regarding the rationale for applying this theory to our study, please see Appendix A.1.

# in each iteration, keeping the dataset updated for subsequent retrieval and evaluation. Specifically, the iterative simulation process unfolds as follows: 1) Baseline Establishment: Utilizing an initial dataset comprised of humanauthored text unaffected by LLM $(D_0)$ , ascertain the performance of a benchmark RAG pipeline.

updated indexing structure that supports incorpo-

rating the newly generated LLM text into the index

2) Zero-shot Text Introduction: The baseline dataset  $D_0$  is enriched with text set  $T_{\rm LLM}^{(\rm zero-shot)}$ generated by LLMs in zero-shot manner, yielding  $D_1 = D_0 \cup T_{\text{LLM}}^{(\text{zero-shot})}$ . This simulates the evolution of users' application of LLMs from initial zero-shot deployments to sophisticated RAG configurations. 3) Retrieval and Re-ranking: For each query q, a subset of documents  $D'_i$  is retrieved from the dataset  $D_i$  through a retrieval and optional re-ranking step  $R(q, D_i) \rightarrow D'_i$ . The retrieval function R remains constant throughout the experimental process to control variables. 4) Generation Phase: Answers S are generated using the LLMs  $(G(p, D', q, K) \rightarrow s)$  with a uniform prompt p in the experiment. 5) Post-processing Phase: Postprocess S to obtain S', removing text fragments that may expose the identity of the LLMs. 6) Index **Update**: Integrate S' into  $D_i$  to update the dataset to  $D_{i+1}$ . 7) Iterative Operation: Repeat steps 3 to 6 for each new dataset  $D_{i+1}$ , until the required number of iterations t is reached.

The pseudo-code for this process is presented in Appendix A.2. Through the simulation process, we can observe how LLM-generated text influences the entire pipeline of RAG systems, and how this impact evolves over time and with data accumulation. Due to the relatively infrequent update cycles of individual LLMs, we assume that the LLMs remain static within the simulation, and leave the effects of their evolution for future research. For details of the prompts and post-processing steps please refer to Appendix A.8 and A.3.

#### 4 Experiment

Datasets and Metrics. We conduct experiments on commonly used ODQA datasets, including NQ (Kwiatkowski et al., 2019), WebQ (Berant et al., 2013), TriviaQA (Joshi et al., 2017), and PopQA (Mallen et al., 2022). We preprocess the datasets following Yu et al. (2023) and Zhang et al. (2023). Given the constraints on experimental resources, we randomly select 200 samples from each

In this section, a simulation framework is designed to explore the potential impacts that texts generated by LLMs may have on RAG systems. This framework models a simplified process that tracks how RAG systems gradually adjust their responses as they accumulate LLM-generated text over time.

## 3.1 Preliminaries

169

170

171

172

173

174

175

176

177

178

179

181

182

183

184

188

189

190

191

192

194

195

196

197

198

199

201

205

209

210

211

213

214

215

216

217

An RAG system f can be formalized as  $f:(Q \times$  $D \times K$   $\rightarrow S$ , where Q is the set of queries, D represents a large collection of documents, K is the knowledge within the LLM, and S is the set of text outputs generated by the system. For a particular query  $q \in Q$ , the goal of the RAG system is to find a mapping f(q, D, K) = s that produces a response text  $s \in S$  satisfying the query q. This process involves two stages:

Retrieval Stage, executed by the retrieval function R, is formally defined as  $R: (Q \times D) \to D'$ , where  $D' \subseteq D$  represents the subset of documents judged by R to be most relevant to the query q.

Generation Stage, executed by the generation function  $G: (P \times Q \times D' \times K) \to S$ . Its task is to utilize the prompt  $p \in P$ , the query q, the related document subset D', and the knowledge of LLMs K to construct the answer s.

Within the entire RAG system f, the functions Rand G act in series to form a process expressed as f(q, D, K) = G(p, q, R(q, D), K). In this manner, the RAG system integrates the precision of IR with the richness of LLMs to provide informationrich content when answering questions.

#### **Simulation Process** 3.2

Our simulation process starts with a pure humanauthored text dataset and gradually introduces the LLM-generated text, observing how this change over time affects the RAG system. Adhering to the specifications outlined in Section 3.1, the RAG architecture is instantiated and expanded with additional details. In the **retrieval stage**, we apply sparse and dense retrieval strategies to obtain a candidate document set that is relevant to the query. Additionally, we also have the option to perform a re-ranking of the candidate documents to further optimize the process. In the generation stage, we use the LLMs which are widely used in research and practical applications to generate responses. To accurately simulate the evolution process of the RAG system, we specifically use an iteratively

267

test set for a comprehensive analysis. For evalu-268 ating the retrieval phase in ODQA tasks during the simulation, we utilize the widely accepted metrics of Acc@5 and Acc@20 following Karpukhin et al. (2020). These metrics assess the proportion of questions where the correct answers appear in the top 5 or top 20 retrieval results, respectively. For the answer quality of the LLM output for each iteration, we follow Chen et al. (2023) by applying the Exact Match (EM) metric, which checks if the correct answer is fully contained within the generated text. Furthermore, in Section 5, we adopt a holistic perspective to examine the RAG pipeline, with a focus on the interaction between the retrieval and generation phases and how the ranking of humangenerated texts changes over time.

269

270

271

273

274

275

276

277

279

281

282

289

290

294

301

303

305

306

307

311

313

314

Retrieval and Re-ranking Methods. In our experiments, we employ a variety of retrieval methods, including the sparse method BM25, the contrastive learning-based dense retriever Contriever (Izacard et al., 2022), the advanced BGE-Base (Xiao et al., 2023) retriever, and the LLM-Embedder (Zhang et al., 2023) designed for LLMs. For the results retrieved using BM25 and BGE-Base, we separately apply the T5-based (Raffel et al., 2020) re-ranking model MonoT5-3B (Nogueira et al., 2020), the UPR-3B (Sachan et al., 2022) which uses the unsupervised capabilities of T0-3B (Sanh et al., 2022), and the BGE-Reranker (Xiao et al., 2023), which is based on the XLM-RoBERTa-Large (Conneau et al., 2020).

Generative Models. Considering the complexity and variability of real-world environments, the text that is continuously integrated into the system may be generated by a variety of LLMs. Our iterative experiments incorporate text produced by a suite of prevalent LLMs. These include GPT-3.5-Turbo (OpenAI, 2022), LLaMA2-13B-Chat (Touvron et al., 2023), Qwen-14B-Chat (Bai et al., 2023), Baichuan2-13B-Chat (Yang et al., 2023), and ChatGLM3-6B (Du et al., 2022). This enables the RAG systems to blend varied linguistic styles and knowledge, leading to results that more closely replicate real-world scenarios.

Implementation. Please refer to Appendix A.4.

#### 5 Results

In this section, within the simulation framework, we meticulously examine both the initial and the 315 extended iterations. We define the short-term effect as the immediate effects observed in the first itera-317

tion, while the long-term effect is analyzed based 318 on the second to the tenth iteration. We investigate 319 if the "Spiral of Silence" effect occurs within these 320 stages and how the RAG systems might respond to 321 it. Under the task settings of ODQA, we analyze 322 the potential influence of the "Spiral of Silence" on 323 RAG systems' response patterns. 324

325

326

327

328

329

331

332

333

334

335

336

337

338

339

340

341

342

343

344

345

346

347

348

349

350

351

352

354

355

356

357

358

359

360

361

362

363

364

365

366

367

#### 5.1 Short-Term Effects on RAG Performance

When comparing RAG system results using different retrieval methods on the original dataset versus the augmented one in the first iteration, we observe that: 1) Immediate Impact of LLM-Generated Text on the RAG System: The introduction of a minimal amount of LLM-generated text produces immediate effects on both retrieval and QA performance of the RAG system, as shown in Table 1 and Figure 2. 2) LLM-Generated Text Generally Improves Retrieval Accuracy: Table1 reveals that adding LLM-generated responses to a dataset typically enhances the accuracy of retrieval systems, as measured by Acc@5 and Acc@20 metrics. For example, using the BM25 on TriviaQA resulted in accuracy improvements of 31.2% and 19.1% respectively. However, a slight decline in Acc@5 is also observed in certain cases. This suggests a primarily positive, yet complex, impact of LLMgenerated text on retrieval accuracy. 3) The impact on QA performance is mixed: Due to space constraints, we only present the results of four retrieval methods. As shown in Figure 2, while the RAG system's QA performance typically surpasses the zero-shot LLM outputs, the addition of LLM text can either enhance or impair QA performance depending on the dataset and retrieval strategy. It appears to enhance performance for TriviaQA, but for NQ and PopQA, the effect is detrimental with non-BM25 retrieval methods, suggesting that without significant retrieval enhancement, LLM text inclusion might be counterproductive.

#### 5.2 Long-term Effects on RAG Performance

In this section, we investigate whether the shotterm effects are predictive of the long-term behavior of the system. We present the results on NQ and PopQA in Figure 3; for results on other datasets, please refer to Figure 8 in Appendix A.5. We observe that: 1) Decreased Retrieval Effectiveness Over Time: Figures 3a and 3b show a general decline in Acc@5 across successive iterations for most methods, with an average drop of 21.4% for NQ and 19.4% for PopQA from the first iteration

Model	NQ			WebQ			TriviaQA			PopQA						
WIOUCI	A	.cc@5	A	cc@20	A	.cc@5	A	cc@20	A	.cc@5	Ac	cc@20	A	cc@5	Ac	c@20
	Ori.	$+LLM_Z$	Ori.	$+LLM_Z$	Ori.	$+LLM_Z$	Ori.	$+LLM_Z$	Ori.	$+ \mathrm{LLM}_Z$	Ori.	$+LLM_Z$	Ori.	$+LLM_Z$	Ori.	$+LLM_Z$
BM25	49.0	57.5	67.0	73.5	41.0	51.0*	63.0	71.0	62.5	82.0*	73.0	87.0*	35.5	41.5	51.5	59.5
Contreiver	68.0	68.5	84.0	85.0	66.0	69.5	74.0	80.0	68.0	83.5*	80.5	87.5	62.0	65.0	77.5	79.5
LLM-Embedder	75.5	75.5	86.5	88.0	62.5	72.5*	76.0	79.5	67.5	81.0*	77.5	87.5*	70.0	67.5	79.5	82.0
$BGE_{base}$	77.0	73.0	86.0	86.0	65.5	71.5	77.0	80.0	69.5	81.5*	80.0	87.5*	72.0	70.0	83.0	84.5
BM25+UPR	63.0	66.5	73.5	78.0	57.0	68.0*	68.5	75.0	71.5	83.0*	78.0	89.0*	57.5	61.5	60.0	67.0
BM25+MonoT5	66.5	69.0	74.5	80.5	62.0	67.5	69.5	76.0	72.0	83.5*	78.0	$88.0^{*}$	53.5	58.5	59.5	66.5
$BM25+BGE_{reranker}$	68.0	69.5	76.5	81.0	64.5	68.5	71.0	76.0	72.5	84.0*	78.0	88.5*	54.0	61.0	60.0	67.5
BGE <sub>base</sub> +UPR	75.5	71.5	87.5	88.0	64.0	69.0	77.0	79.5	76.0	84.0*	84.5	89.5	76.0	71.0	84.5	84.5
BGE <sub>base</sub> +MonoT5	75.0	70.5	86.5	86.5	68.5	72.0	78.0	81.5	77.0	83.5	83.5	89.5	72.0	72.5	85.5	86.0
$BGE_{base} \text{+} BGE_{reranker}$	69.0	68.0	84.0	84.5	67.5	70.5	78.0	81.5	72.5	83.5*	82.0	88.0	73.0	70.0	84.0	85.0

Table 1: Short-term retrieval performance. A blue background indicates a decrease in retrieval results after the incorporation of LLM-generated text, while a purple background signifies an increase. The deeper the color, the larger the discrepancy from the original results. Statistical significance at 0.05 relative to origin is marked with \*.



Figure 2: Short-Term QA performance. For each retrieval method, we document both the average performance and the range of variation exhibited by five LLMs. A red dashed line symbolizes the average EM score for zero-shot question generation by LLMs. "Ori." and "+LLM<sub>Z</sub>" represent the average EM values when models use the original dataset or a dataset enhanced with LLM-generated texts as context, respectively. Retrieval methods are abbreviated: "Contri" for Contriever, "LLM-E" for LLM-Embedder, and "BGE-B" for BGE<sub>base</sub>.

to the last, except for a temporary improvement in BM25 during the second iteration on PopQA. 369 This trend signals that the retrieval quality boost provided by LLM-generated text may be transient, with a propensity for degradation over time. 2) **Stability in QA Performance Despite Retrieval** Decline: Contrary to expectations, the QA per-374 formance does not mirror the retrieval accuracy's decrease. As shown in Figure 3c and 3d, the EM exhibit slight variations but generally maintain their 377 level throughout the iterations. While a diminished retrieval accuracy intuitively seems to undermine 379 the system's capacity to output correct answers, this does not unequivocally translate into a decline in QA efficacy. In subsequent sections, we will delve deeper into the reasons behind these observations and examine the complex dynamic relationship that may exist between retrieval and QA performance.

#### 5.3 Spiral of Silence

387

In the context of LLM-augmented RAG systems, we have observed a rapid shift in response to the integration of LLM-generated text, a decline in retrieval performance over time, and stability in QA performance despite retrieval decline. To explain these phenomena, we draw on the theory of the "Spiral of Silence" as posited by Noelle-Neumann (Noelle-Neumann, 1974), extending its principles to the behavior of RAG systems enhanced by LLMs. To explore the presence of a "Spiral of Silence" phenomenon, we propose three Hypotheses for investigation. (H1): Dominance of LLM-Generated Texts. Retrieval models are more likely to prioritize LLM-generated text in search results, which could result in LLMgenerated text taking a dominant position in the retrieval hierarchy. (H2): Marginalization of Human-Generated Content. If human-authored text consistently loses ranking prominence through successive iterations, it may be excluded from the top results until it becomes invisible, thus creating a silence. (H3): Homogenization of Opinions. The preferential ranking of LLM-generated text could culminate in a uniformity of displayed perspectives by the RAG system, potentially sidelining the accuracy or variety of the information.

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411



Figure 3: Long-Term RAG performance. The upper section illustrates the retrieval outcomes for various methods, while the lower section depicts the average EM across LLMs. Iteration1 represents the results following the incorporation of zero-shot LLM-generated text. Abbreviated re-ranking methods in the legend are: +U for UPR, +M for MonoT5, and +BR for BGE-Reranker.

To verify (H1), we analyze Iteration1 where LLM-generated texts are first introduced to the retrieval system. We calculate the proportion of these texts appearing in the top five search results:

413

414

415

416

417

418

419

420

421

422

423

494

425

426

427

428

429

430

431

432

433

434

435

436

437

438

$$\mathbf{P} = \frac{\sum_{q \in Q} c_q^{LLM}}{\sum_{q \in Q} (c_q^{LLM} + c_q^{Human})} \times 100\%$$
(1)

where  $c_q^{LLM}$  is the count of LLM-generated texts and  $c_{a}^{Human}$  is the count of human-generated texts in the top five search results for query q. Table 2 reveals that, even with a modest inclusion of LLMgenerated texts, most retrieval models often rank them at the top. This behavior supports the findings of Dai et al. (2023a), where LLM-rewritten texts are preferred by retrieval models over the originals. Our study extends this by directly generating query-specific texts with LLMs. The preference might stem from inherent biases within the system or the actual relevance of the LLM-produced content. This suggests retrieval systems tend to favor LLM-generated texts, making them more prominent in search results, which can rapidly influence an RAG system's behavior.

To validate (H2), we incorporate a **temporal dimension**, observing the percentage change of texts generated by various LLMs and humans within the top 50 search results across different datasets over time. As shown in Figure 4, after ten iterations, the

Method	NQ	WebQ	TriviaQA	PopQA
BM25	34.1	19.6	57.6	23.9
Contriever	72.8	75.2	80.1	67.0
LLM-Embedder	68.2	64.6	75.3	70.0
$BGE_{base}$	80.7	84.1	85.6	81.5
BM25+UPR	62.3	49.8	75.7	47.1
BM25+MonoT5	66.2	55.8	83.0	47.1
BM25+BGE <sub>reranker</sub>	64.4	55.2	81.6	46.6
BGE <sub>base</sub> +UPR	74.4	69.3	79.1	71.2
BGE <sub>base</sub> +MonoT5	81.4	84.0	88.4	74.3
$BGE_{base}$ + $BGE_{reranker}$	67.2	74.2	83.2	72.8

Table 2: Percentage of LLM-generated documents occupying the top-5 retrieval results, after augmenting each query with five LLM-generated documents. The blue background indicates a majority presence of humangenerated documents, while the purple background denotes a predominance of LLM-generated documents.



Figure 4: Average percentage of texts from various sources within the top 50 search results over multiple iterations across different search methods. For results on WebQ and TriviaQA, please refer to Figure 9 in Appendix A.5.

percentage of human-generated texts significantly decreased, falling below 10% for all datasets. This pattern suggests a sustained diminishing impact of human-contributed texts and hints at the possibility of their eventual exclusion from search results if trends continue.

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

To explore (H3), we examine the risk of potential viewpoint homogenization in the RAG system from both the diversity and accuracy dimensions during the simulation. **Diversity** is quantified using the Self-BLEU score, where a higher score indicates less diversity, suggesting a convergence of viewpoints. As shown in Figure 5, upon introducing zero-shot LLM-generated texts (Iteration1), the Self-BLEU scores across different datasets experience varying degrees of change. However, over more iterations, the Self-BLEU scores for the top 5 results consistently rise and plateau across all datasets, indicating a significant reduction in textual diversity with each iterative cycle. Subsequently, we assessed whether the **accuracy** of the top documents returned by the IR system tends toward uniformity over time. Figure 6 charts the



Figure 5: 3-gram Self-BLEU score for the top 5 search results over iterations, from the original dataset (Ori.) to the inclusion and subsequent iterations with LLMgenerated texts. For results on WebQ and TriviaQA, please refer to Figure 10 in Appendix A.5.

number of documents with the correct answer in 462 463 the top 5 results ("Context Right Num") against the number of queries LLM answers correctly or incor-464 rectly, across the Iteration 1, 2, 5, 10. For simplicity, 465 we showcase only the NQ dataset's averaged out-466 comes, but we find the same trends across other 467 datasets. It indicates that in the initial iterations, 468 fewer correct answer documents (e.g., "Context 469 Right Num" of 0, 1, or 2) typically correlate with 470 a greater number of LLM-answered queries being 471 incorrect (EM=0). Despite this, a significant frac-472 tion of the top documents still include the correct 473 answer. When the LLM correctly answers a query 474 (EM=1), the correct answer documents within the 475 top results can range anywhere from 1 to 5. As 476 the iterations continue, the frequency of having 1 477 to 4 correct answer documents in the top 5 results 478 for each query diminishes, and by the Iteration10, 479 contexts for LLM-correct queries almost always 480 contain the correct answer, whereas contexts for 481 LLM-incorrect answers almost always do not. This 482 pattern demonstrates a trend towards polariza-483 tion and uniformity in the accuracy of provided 484 contexts as the RAG system iterates. 485

At this point, we have confirmed through our experiments the presence of the "Spiral of Silence" phenomenon, as outlined in three tested hypotheses. Moreover, the dashed lines in Figure 6 represent the total number of correct and incorrect LLM answered queries, along with the Acc@5 retrieval metric over various iterations. The LLM's rate of correct answers remains constant through the iterations, aligning to Section 5.2. However, as iterations advance, correct answers diminish within top documents for LLM-incorrect queries, reducing their contribution to the Acc@5 and thus decreasing retrieval performance. In contrast, for LLM-correct queries, more retrieved documents containing the correct answer do not affect Acc@5

486

487

488

489

490

491

492

493

494

495

496

497

498

499

or EM. This pattern is explicable by the "Spiral of Silence" theory discussed in Section 5.2, which accounts for the observed dip in IR results and the sustained QA performance.

501

502

503

504

505

506

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

#### 5.4 Effects of Spiral of Silence on ODQA

We will delve into a more nuanced discussion of the impact of the "Spiral of Silence" within the context of ODQA. It is important to note that the influence of the phenomenon is not confined to this scenario; it may also be pertinent across all settings that involve knowledge retrieval, generation, and the influx of text from LLMs. Specifically, our analysis is structured around two dimensions: the query level and the document level.

At the **query level**, Figure 7a signifies the average count of queries shifting between consecutive iterations from incorrect to correct and vice versa, respectively. Notably, during the 1->2 iteration, there is an initial surge in both metrics, which subsequently experience a sharp decline as the iteration count escalates. This suggests that the LLM-generated text's initial introduction catalyzes a more dynamic state, likely due to the correction of existing errors or the introduction of new inaccuracies. Over time, however, the "Spiral of Silence" effect seems to guide the system towards a state of equilibrium where the transition rate stabilizes to less than 1% per 200 queries. This means most queries maintain their status as either correct or incorrect, indicating that individual query QA results become fixed.

At the **document level**, we compute the average rank shifts of the first documents containing the correct answer within retrieval results, under different LLM answer states. In Figure 7b, we observe that: 1) Different Trend of Correct Answer Rankings Based on Source. In instances where EM=0, correct documents from all sources and from humans both tend to be ranked lower over time. When EM=1, the rankings for correct documents from all sources improve slightly, while rankings for correct answers from humans continue to decline. This suggests that the LLM's correct texts gain prominence in retrieval rankings over time, overshadowing correct texts from human-generated texts. 2) Gradual Dysfunction of the IR System in Incorrect LLM Responses: When the LLM provides incorrect answers (EM=0), there's a risk that documents that once rose to the top with accurate information might increasingly be obscured



Figure 6: Correlation between the number of top 5 search results containing the correct answer ("Context Right Num") and the accuracy of responses given by LLMs on the NQ dataset. The responses are categorized based on Exact Match (EM) score: EM=1 for correct and EM=0 for incorrect. The overall number of queries that the LLMs answered correctly (EM=1 Total) and incorrectly (EM=0 Total), along with the average retrieval accuracy (Acc@5) are shown by dashed lines. The results are averaged across different LLMs, retrieval and ranking methods.



(a) Average Query State Transition Number from incorrect to correct ('Average 0->1') and correct to incorrect ('Average 1->0') between consecutive iterations, aggregated across all IR methods and datasets.

(b) Ranking of the first retrieved document containing the correct answer in each iteration, with LLM responses being incorrect (EM=0) or correct (EM=1). "First Right From All Sources" refers to the average rank for the first text containing the correct answer, where the source could be either LLM or a human; "First Right From Human" is for human-only sources, both considered across datasets and LLMs.

Figure 7: Effects of "Spiral of Silence" on query and document level of RAG systems.

by the growing mass of LLM-generated content. This can lead to a scenario where the IR system, originally intended to help users find precise information, becomes less reliable. If it prioritizes and disseminates the LLM's inaccuracies, a feedback loop could ensue, solidifying these errors. This concerning trend highlights the critical need for ongoing adjustments and improvements to the IR systems to uphold their purpose.

#### 6 Anlysis

551

552

553

554

557

558

559

561To further comprehend the "Spiral of Silence" ef-562fect, we illustrate its interaction with misinforma-563tion introduced by adversaries using LLMs. More-564over, we test two information filtering mechanisms565to alleviate the progression of the effect. For more566information on the experimental setup and results,567refer to Appendix A.6 and Appendix A.7.

#### 7 Conclusion

In this study, we initiate our research from empirical observations, aiming to investigate the implications of progressively integrating LLM-generated text into RAG systems. To this end, We employ the ODQA task as a case study to examine both the immediate and extended impacts of LLM text on these systems. Our simulation has revealed the emergence of a "Spiral of Silence" effect, suggesting that without appropriate intervention, humangenerated content may progressively diminish its influence within RAG systems. Further investigation into this phenomenon reveals that unchecked accumulation of erroneous LLM-generated information could lead to the overlooking of correct information by IR systems, resulting in harm. We urge the academic community to be vigilant and take measures to prevent the potential misuse of LLM-generated data.

568

570

571

572

573

574

575

576

577

578

579

580

582

583

584

585

## Limitations

587

613

616

617

618

619

621

623

624

625

626

627

628

629

630

631

633

634

This study aims to present a new perspective on the impact of LLM-generated texts entering the internet on RAG systems. However, the complexity of 590 reality means that it is impossible to account for all 591 variables. The methods of LLM text generation and the mechanisms by which this content enters the 593 retrieval set are constantly changing, which could affect the performance of RAG systems. Dynamic 595 updates to LLMs could also impact the outcomes, 596 and future research should incorporate this variable. 597 While ODQA serves as an insightful approach to evaluate the progression of RAG systems, it is nec-599 essary to recognize that ODQA assessments are not exhaustive in capturing the full spectrum of information retrieval scenarios. Nonetheless, the simulation framework proposed in this research is readily adaptable to other tasks that employ RAG systems. Our discussion introduces the "Spiral of Silence" as a potential outcome of the proliferation of LLM-generated texts. Although such a devel-607 608 opment is not predetermined, given the myriad of factors at play in the real world, this work aims to foster a deeper investigation into the phenomenon 610 and its prospective implications for information 611 diversity in AI-mediated environments.

#### Ethical Considerations

This paper only explores the potential impact of the LLM-generated text, without involving the release of the generated text and the intervention of social progress, so the possibility of ethical risks is small. We used publicly available LLMs and datasets to conduct experiments that did not involve any ethical issues. In the appendix, we analyze the potential interplay between harmful information and the phenomena outlined in our paper, with a principal objective to draw attention to this issue to advocate for its resolution.

#### References

- Vaibhav Adlakha, Parishad BehnamGhader, Xing Han Lu, Nicholas Meade, and Siva Reddy. 2023. Evaluating correctness and faithfulness of instructionfollowing models for question answering. *CoRR*, abs/2307.16877.
- Faisal Alatawi, Lu Cheng, Anique Tahir, Mansooreh Karami, Bohan Jiang, Tyler Black, and Huan Liu. 2021. A survey on echo chambers on social media: Description, detection and mitigation. *CoRR*, abs/2112.05084.

Sina Alemohammad, Josue Casco-Rodriguez, Lorenzo Luzi, Ahmed Imtiaz Humayun, Hossein Babaei, Daniel LeJeune, Ali Siahkoohi, and Richard G. Baraniuk. 2023. Self-consuming generative models go MAD. *CoRR*, abs/2307.01850.

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

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu, Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi Tan, Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang, Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu, Hongyi Yuan, Zheng Yuan, Jianwei Zhang, Xingxuan Zhang, Yichang Zhang, Zhenru Zhang, Chang Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang Zhu. 2023. Qwen technical report. *CoRR*, abs/2309.16609.
- Jonathan Berant, Andrew Chou, Roy Frostig, and Percy Liang. 2013. Semantic parsing on freebase from question-answer pairs. In Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing, EMNLP 2013, 18-21 October 2013, Grand Hyatt Seattle, Seattle, Washington, USA, A meeting of SIGDAT, a Special Interest Group of the ACL, pages 1533–1544. ACL.
- Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Millican, George van den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, Diego de Las Casas, Aurelia Guy, Jacob Menick, Roman Ring, Tom Hennigan, Saffron Huang, Loren Maggiore, Chris Jones, Albin Cassirer, Andy Brock, Michela Paganini, Geoffrey Irving, Oriol Vinyals, Simon Osindero, Karen Simonyan, Jack W. Rae, Erich Elsen, and Laurent Sifre. 2022. Improving language models by retrieving from trillions of tokens. In *International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA*, volume 162 of *Proceedings of Machine Learning Research*, pages 2206–2240. PMLR.
- Deng Cai, Yan Wang, Wei Bi, Zhaopeng Tu, Xiaojiang Liu, Wai Lam, and Shuming Shi. 2019. Skeletonto-response: Dialogue generation guided by retrieval memory. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers), pages 1219–1228. Association for Computational Linguistics.
- Canyu Chen and Kai Shu. 2023. Combating misinformation in the age of llms: Opportunities and challenges. *CoRR*, abs/2311.05656.
- Jiawei Chen, Hongyu Lin, Xianpei Han, and Le Sun. 2023. Benchmarking large language models in retrieval-augmented generation. *CoRR*, abs/2309.01431.

803

804

805

- Uthsav Chitra and Christopher Musco. 2020. Analyzing the impact of filter bubbles on social network polarization. In WSDM '20: The Thirteenth ACM International Conference on Web Search and Data Mining, Houston, TX, USA, February 3-7, 2020, pages 115–123. ACM.
- Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Unsupervised cross-lingual representation learning at scale. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020, pages 8440–8451. Association for Computational Linguistics.
- Sunhao Dai, Yuqi Zhou, Liang Pang, Weihao Liu, Xiaolin Hu, Yong Liu, Xiao Zhang, and Jun Xu. 2023a. Llms may dominate information access: Neural retrievers are biased towards llm-generated texts. *CoRR*, abs/2310.20501.

710

711

712

713

715

718

719

720

721

722

723

724

726

727

728

729

731

732

733

734

736

737

738

740

741

742

743

744

745

747 748

- Yi Dai, Hao Lang, Yinhe Zheng, Bowen Yu, Fei Huang, and Yongbin Li. 2023b. Domain incremental lifelong learning in an open world. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 5844–5865, Toronto, Canada. Association for Computational Linguistics.
- Zhengxiao Du, Yujie Qian, Xiao Liu, Ming Ding, Jiezhong Qiu, Zhilin Yang, and Jie Tang. 2022. Glm: General language model pretraining with autoregressive blank infilling. In Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 320–335.
- Josh A. Goldstein, Girish Sastry, Micah Musser, Renee DiResta, Matthew Gentzel, and Katerina Sedova. 2023. Generative language models and automated influence operations: Emerging threats and potential mitigations. *CoRR*, abs/2301.04246.
- Google. 2023. Introducing gemini: our largest and most capable ai model. https: //blog.google/technology/ai/ google-gemini-ai/?utm\_source=gdm& utm\_medium=referral#sundar-note. Accessed: 2023-12-06.
  - Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Ming-Wei Chang. 2020. REALM: retrievalaugmented language model pre-training. *CoRR*, abs/2002.08909.
- Hangfeng He, Hongming Zhang, and Dan Roth. 2023. Rethinking with retrieval: Faithful large language model inference. *CoRR*, abs/2301.00303.
- Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, and Ting Liu. 2023a. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. *CoRR*, abs/2311.05232.

- Yangsibo Huang, Samyak Gupta, Mengzhou Xia, Kai Li, and Danqi Chen. 2023b. Catastrophic jailbreak of open-source llms via exploiting generation. *CoRR*, abs/2310.06987.
- Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand Joulin, and Edouard Grave. 2022. Unsupervised dense information retrieval with contrastive learning. *Trans. Mach. Learn. Res.*, 2022.
- Gautier Izacard and Edouard Grave. 2021. Leveraging passage retrieval with generative models for open domain question answering. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume, EACL 2021, Online, April 19 - 23, 2021,* pages 874– 880. Association for Computational Linguistics.
- Gautier Izacard, Patrick S. H. Lewis, Maria Lomeli, Lucas Hosseini, Fabio Petroni, Timo Schick, Jane Dwivedi-Yu, Armand Joulin, Sebastian Riedel, and Edouard Grave. 2023. Atlas: Few-shot learning with retrieval augmented language models. *J. Mach. Learn. Res.*, 24:251:1–251:43.
- Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Yejin Bang, Andrea Madotto, and Pascale Fung. 2023. Survey of hallucination in natural language generation. *ACM Comput. Surv.*, 55(12):248:1–248:38.
- Ray Jiang, Silvia Chiappa, Tor Lattimore, András György, and Pushmeet Kohli. 2019. Degenerate feedback loops in recommender systems. In Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society, AIES 2019, Honolulu, HI, USA, January 27-28, 2019, pages 383–390. ACM.
- Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2019. Billion-scale similarity search with GPUs. *IEEE Transactions on Big Data*, 7(3):535–547.
- Mandar Joshi, Eunsol Choi, Daniel S. Weld, and Luke Zettlemoyer. 2017. Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, ACL* 2017, Vancouver, Canada, July 30 - August 4, Volume 1: Long Papers, pages 1601–1611. Association for Computational Linguistics.
- Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick S. H. Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense passage retrieval for open-domain question answering. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020, pages 6769–6781. Association for Computational Linguistics.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur P. Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion Jones, Matthew

810

- 811 812
- 813 814
- 815 816
- 817
- 818 819
- 820 821
- 822
- 8
- 825 826

827 828

- 829 830 831
- 832 833 834

835 836 837

- 838 839 840 841 842
- 84 84 84

847

8

- 0 8
- 8
- 8

854 855

- 856 857 858
- 0 8

8

861

Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. Natural questions: a benchmark for question answering research. *Trans. Assoc. Comput. Linguistics*, 7:452–466.

- Patrick S. H. Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. Retrieval-augmented generation for knowledge-intensive NLP tasks. In Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual.
- Paul Pu Liang, Chiyu Wu, Louis-Philippe Morency, and Ruslan Salakhutdinov. 2021. Towards understanding and mitigating social biases in language models. In Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event, volume 139 of Proceedings of Machine Learning Research, pages 6565–6576. PMLR.
- Chen Lin, Dugang Liu, Hanghang Tong, and Yanghua Xiao. 2022. Spiral of silence and its application in recommender systems. *IEEE Trans. Knowl. Data Eng.*, 34(6):2934–2947.
- Chang Liu, Xiaoguang Li, Lifeng Shang, Xin Jiang, Qun Liu, Edmund Y. Lam, and Ngai Wong. 2023a.
  Gradually excavating external knowledge for implicit complex question answering. In *Findings of the Association for Computational Linguistics: EMNLP* 2023, Singapore, December 6-10, 2023, pages 14405– 14417. Association for Computational Linguistics.
- Dugang Liu, Chen Lin, Zhilin Zhang, Yanghua Xiao, and Hanghang Tong. 2019. Spiral of silence in recommender systems. In Proceedings of the Twelfth ACM International Conference on Web Search and Data Mining, WSDM 2019, Melbourne, VIC, Australia, February 11-15, 2019, pages 222–230. ACM.
- Ye Liu, Semih Yavuz, Rui Meng, Meghana Moorthy, Shafiq Joty, Caiming Xiong, and Yingbo Zhou. 2023b. Exploring the integration strategies of retriever and large language models. *CoRR*, abs/2308.12574.
- Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das, Hannaneh Hajishirzi, and Daniel Khashabi. 2022. When not to trust language models: Investigating effectiveness and limitations of parametric and nonparametric memories. *CoRR*, abs/2212.10511.
- Elisabeth Noelle-Neumann. 1974. The spiral of silence a theory of public opinion. *Journal of communication*, 24(2):43–51.
- Rodrigo Frassetto Nogueira, Zhiying Jiang, Ronak Pradeep, and Jimmy Lin. 2020. Document ranking with a pretrained sequence-to-sequence model. In *Findings of the Association for Computational Linguistics: EMNLP 2020, Online Event, 16-20 November 2020*, volume EMNLP 2020 of *Findings of ACL*,

pages 708–718. Association for Computational Linguistics.

863

864

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

889

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

- OpenAI. 2022. ChatGPT. https://openai.com/ blog/chatgpt. Accessed: January 10, 2024.
- OpenAI. 2023. GPT-4 technical report. CoRR, abs/2303.08774.
- Yikang Pan, Liangming Pan, Wenhu Chen, Preslav Nakov, Min-Yen Kan, and William Wang. 2023. On the risk of misinformation pollution with large language models. In *Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023*, pages 1389–1403. Association for Computational Linguistics.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *J. Mach. Learn. Res.*, 21:140:1–140:67.
- Ori Ram, Yoav Levine, Itay Dalmedigos, Dor Muhlgay, Amnon Shashua, Kevin Leyton-Brown, and Yoav Shoham. 2023. In-context retrieval-augmented language models. *CoRR*, abs/2302.00083.
- Devendra Singh Sachan, Mike Lewis, Mandar Joshi, Armen Aghajanyan, Wen-tau Yih, Joelle Pineau, and Luke Zettlemoyer. 2022. Improving passage retrieval with zero-shot question generation. In *Proceedings* of the 2022 Conference on Empirical Methods in Natural Language Processing, EMNLP 2022, Abu Dhabi, United Arab Emirates, December 7-11, 2022, pages 3781–3797. Association for Computational Linguistics.
- Victor Sanh, Albert Webson, Colin Raffel, Stephen H. Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Arun Raja, Manan Dey, M Saiful Bari, Canwen Xu, Urmish Thakker, Shanya Sharma Sharma, Eliza Szczechla, Taewoon Kim, Gunjan Chhablani, Nihal V. Nayak, Debajyoti Datta, Jonathan Chang, Mike Tian-Jian Jiang, Han Wang, Matteo Manica, Sheng Shen, Zheng Xin Yong, Harshit Pandey, Rachel Bawden, Thomas Wang, Trishala Neeraj, Jos Rozen, Abheesht Sharma, Andrea Santilli, Thibault Févry, Jason Alan Fries, Ryan Teehan, Teven Le Scao, Stella Biderman, Leo Gao, Thomas Wolf, and Alexander M. Rush. 2022. Multitask prompted training enables zero-shot task generalization. In The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022. OpenReview.net.
- Dietram A Scheufle and Patricia Moy. 2000. Twentyfive years of the spiral of silence: A conceptual review and empirical outlook. *International journal of public opinion research*, 12(1):3–28.
- Nina Schick. 2020. Deep Fakes and the Infocalypse: What You Urgently Need To Know. Monoray.

Ilia Shumailov, Zakhar Shumaylov, Yiren Zhao, Yarin Gal, Nicolas Papernot, and Ross J. Anderson. 2023. The curse of recursion: Training on generated data makes models forget. *CoRR*, abs/2305.17493.

917

918

919

922

924

926

932

933

935

937

941

942

943

947

951

961

962

964

967

969

970

971

973

- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. Llama: Open and efficient foundation language models. *CoRR*, abs/2302.13971.
- Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2023. Interleaving retrieval with chain-of-thought reasoning for knowledgeintensive multi-step questions. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023, pages 10014–10037. Association for Computational Linguistics.
  - Shitao Xiao, Zheng Liu, Peitian Zhang, and Niklas Muennighof. 2023. C-pack: Packaged resources to advance general chinese embedding. *CoRR*, abs/2309.07597.
  - Danni Xu, Shaojing Fan, and Mohan S. Kankanhalli. 2023. Combating misinformation in the era of generative AI models. In *Proceedings of the 31st ACM International Conference on Multimedia, MM 2023, Ottawa, ON, Canada, 29 October 2023- 3 November* 2023, pages 9291–9298. ACM.
- Aiyuan Yang, Bin Xiao, Bingning Wang, Borong Zhang, Ce Bian, Chao Yin, Chenxu Lv, Da Pan, Dian Wang, Dong Yan, Fan Yang, Fei Deng, Feng Wang, Feng Liu, Guangwei Ai, Guosheng Dong, Haizhou Zhao, Hang Xu, Haoze Sun, Hongda Zhang, Hui Liu, Jiaming Ji, Jian Xie, Juntao Dai, Kun Fang, Lei Su, Liang Song, Lifeng Liu, Liyun Ru, Luyao Ma, Mang Wang, Mickel Liu, MingAn Lin, Nuolan Nie, Peidong Guo, Ruiyang Sun, Tao Zhang, Tianpeng Li, Tianyu Li, Wei Cheng, Weipeng Chen, Xiangrong Zeng, Xiaochuan Wang, Xiaoxi Chen, Xin Men, Xin Yu, Xuehai Pan, Yanjun Shen, Yiding Wang, Yiyu Li, Youxin Jiang, Yuchen Gao, Yupeng Zhang, Zenan Zhou, and Zhiying Wu. 2023. Baichuan 2: Open large-scale language models. CoRR, abs/2309.10305.
- Wenhao Yu, Dan Iter, Shuohang Wang, Yichong Xu, Mingxuan Ju, Soumya Sanyal, Chenguang Zhu, Michael Zeng, and Meng Jiang. 2023. Generate rather than retrieve: Large language models are strong context generators. In *The Eleventh International Conference on Learning Representations*, *ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. Open-Review.net.
- Peitian Zhang, Shitao Xiao, Zheng Liu, Zhicheng Dou, and Jian-Yun Nie. 2023. Retrieve anything to augment large language models. *CoRR*, abs/2310.07554.

Shuyan Zhou, Uri Alon, Frank F. Xu, Zhengbao Jiang,<br/>and Graham Neubig. 2023. Docprompting: Gener-<br/>ating code by retrieving the docs. In The Eleventh<br/>International Conference on Learning Representa-<br/>tions, ICLR 2023, Kigali, Rwanda, May 1-5, 2023.978<br/>979

980

981

982

Terry Yue Zhuo, Yujin Huang, Chunyang Chen, and Zhenchang Xing. 2023. Exploring AI ethics of chatgpt: A diagnostic analysis. *CoRR*, abs/2301.12867. Algorithm 1 Simulation Process

function RUNRAG(D, Q)  $Res \leftarrow empty \text{ list}$ for  $q \in Q$  do  $D' \leftarrow \text{RETRIEVE}(q, D)$   $p \leftarrow \text{GENPROMPT}(q)$   $S \leftarrow \text{GENANSWER}(p, D', q)$   $S' \leftarrow \text{POSTPROC}(S)$ Add (q, D', S') to Resend for return Resend function

 $\begin{array}{l} D_0 \leftarrow \text{LOADDATA} \\ initRes \leftarrow \text{RUNRAG}(D_0,Q) \\ basePerf \leftarrow \text{EVALRAG}(initRes) \\ T \leftarrow \text{GENZEROSHOT} \\ D_1 \leftarrow D_0 \text{ combined with } T \\ t \leftarrow \text{number of iterations} \\ \textbf{for } i \leftarrow 1 \text{ to } t \text{ do} \\ iterRes \leftarrow \text{RUNRAG}(D_i,Q) \\ perf_i \leftarrow \text{EVALRAG}(iterRes) \\ D_{i+1} \leftarrow \text{UPDATEDATA}(D_i, iterRes) \\ \textbf{end for} \end{array}$ 

#### A Appendix

983

985

987

990

991

993

995

997

999

1000

1001

1002

1003

1004

1005

1006

## A.1 Discussion of Application on "Spiral of Silence"

In aligning the "Spiral of Silence" theory with the focus of this study, emphasis on the aspect of the "individual's will to express" inherent in the original theory is purposefully diminished. The factors influencing the "Spiral of Silence" phenomenon, as mentioned in Scheufle and Moy (2000), with media and temporality being the principal elements within RAG systems, directly affect the relative standing of LLM and human texts as the system evolves. While the individual's desire to express may be indirectly affected by media and temporality, these are not the primary drivers of the "Spiral of Silence" within RAG systems. In RAG systems, we hypothesize that texts generated by LLMs will increasingly be favored in the hierarchy of information retrieval, whereas texts authored by humans might be systematically marginalized, resulting in a structural form of "passive silencing".

#### A.2 Pseudo-Code of Simulation Process

The pseudo-code of the simulation process in section 3.2 is shown in Algotirhm 1.

#### A.3 Post-Process Details

During the experimental process, we observe that1008the response texts from LLMs occasionally con-<br/>tain specific phrases at the beginning that indicate1009their identity. These phrases are difficult to remove1011through prompts and are irrelevant to the topic at1012hand. Examples include sentences such as:1013

- "I'd be happy to assist you 1014 with your question." 1015
- "According to my knowledge..." 1016

1017

1018

1019

1020

1021

1022

1023

• "As an AI language model..."

We collect over 40 such sentences using a manual annotation approach and filter each LLMgenerated text through string matching. If a matching string is found, the corresponding sentence or fragment is removed.

#### A.4 Implementation Details.

To construct and execute the simulated iterative 1024 framework, we adopt a diverse array of tools and 1025 technologies to facilitate real-time interaction be-1026 tween various retrieval methods, indexing archi-1027 tectures, and LLMs. We implement the APIs of 1028 various LLMs relying on api-for-open-llm<sup>2</sup>. With 1029 an integration of LangChain<sup>3</sup> with Faiss (Johnson 1030 et al., 2019) and Elasticsearch<sup>4</sup>, we execute batched 1031 incremental updates of LLM-generated documents 1032 in each iteration, thus simulating the process of 1033 document index updating by search engines in real-1034 world scenarios. To maintain the diversity of the 1035 generated texts, we set the temperature at 0.7 for 1036 all LLMs. In each iteration of the experiment, except for the zero-shot setting, we keep the size of 1038 the context document set  $D'_i$  fixed at 5. We rerank 1039 the first 100 documents recalled by the retrieval 1040 method when the step is applied. We apply the 1041 LLMs to generate response text, post-process via 1042 rules, and then merge their outputs into the index 1043 for each query in every iteration. Therefore, for 1044 each iteration, we will add 4k new samples to the 1045 index. The total number of iterations t is set at 10, 1046 which results in a total of 40k invocations of the 1047 LLMs for each experimental run. We will make 1048

```
<sup>3</sup>https://github.com/langchain-ai/
```

```
langchain
```

```
<sup>4</sup>https://github.com/elastic/
elasticsearch
```

<sup>&</sup>lt;sup>2</sup>https://github.com/xusenlinzy/ api-for-open-llm



Figure 8: Long-Term RAG performance. The upper section illustrates the retrieval outcomes for various methods, while the lower section depicts the average EM across LLMs. Iteration1 represents the results following the incorporation of zero-shot LLM-generated text. Abbreviated re-ranking methods in the legend are: +U for UPR, +M for MonoT5, and +BR for BGE-Reranker.



Figure 9: Average Percentage of texts from various sources within the top 50 search results over multiple iterations across different search and methods.

our code and data publicly available for further research.

#### A.5 Results on WebQ and TriviaQA

1049

1050

1051

1052

1053

1055

1056

1057

1059

1060

1061

1062

1063

Figure 8 shows long-term RAG performance on WebQ and TriviaQA.

The percentage from various sources and the Self-BLEU of the retrieval results on WebQ and TriviaQA are shown in Figure 9 and Figure 10.

#### A.6 Misinformation in the Iteration

In the previous sections, we explored the impact of non-maliciously generated texts by LLMs on the evolution of the RAG system over time. In this section, we will discuss the persistence of the "Spiral of Silence" effect when attackers deliberately inject specific misinformation into the RAG system, how



Figure 10: 3-gram Self-BLEU score for the top 5 search results over iterations, from the original dataset (Ori.) to the inclusion and subsequent iterations with LLM-generated texts.

misinformation could affect the RAG system over time, and the feasibility of targeted misinformation injection. 1064

1065

1066

1067

1068

1069

1070

1071

1072

1073

1074

1075

1076

1077

1078

1079

1080

1081

1082

1083

1084

1086

1087

1088

1089

1090

1091

1092

1093

1094

1095

1096

1097

**Experimental Setup**: Our experiment follows the CTRLGEN method detailed in Pan et al. (2023), which aligns well with the zero-shot setting used in our trials and simulates the intent of malicious actors to create and propagate false information. Specifically, for each query, we generate five incorrect answers using GPT-3.5-Turbo and then randomly select one to guide five different LLMs to each create a document supporting that incorrect response. These documents replace the zero-shot data in the index from the experiments in Section 3 and are used for simulated iterative experiments. Details of the prompts used are provided in Appendix A.8. For the sake of conciseness, we report only the experimental results for four retrieval methods.

**Experimental Evaluation**: When generating texts containing misinformation using LLMs, we face two primary challenges. First, the model may ignore the instructions, thus inadvertently generating texts that only contain the correct answer. Second, even if the LLM-generated text includes the provided incorrect answer, the content may not genuinely support that answer. To address these issues, we utilize GPT-3.5-Turbo to evaluate the alignment of text t with the given answer a, which could be either correct or incorrect. This evaluation complements the calculation of the EM metric for texts generated by LLMs. We define the EM<sub>llm</sub> metric as follows:

$$\mathbf{EM}_{llm}(t,a) = \begin{cases} 1, & \text{if } t \text{ contains and supports } a \\ 0, & \text{otherwise} \end{cases}$$
(2)

Here, a represents an answer that the text t is being1098evaluated against. The text t is deemed to con-<br/>tribute positively to this metric if it encompasses1099

Model	NQ		WebQ		TriviaQA		PopQA	
	Correct	Incorrect	Correct	Incorrect	Correct	Incorrect	Correct	Incorrect
GPT-3.5-Turbo	0.015	0.7	0.075	0.57	0.105	0.57	0.045	0.71
Baichuan2-13B-Chat	0.11	0.595	0.21	0.44	0.16	0.455	0.1	0.65
Qwen-14B-Cha	0.065	0.61	0.11	0.565	0.165	0.535	0.05	0.7
ChatGLM3-6B	0.085	0.605	0.195	0.435	0.245	0.415	0.105	0.61
LLaMA2-13B-Chat	0.04	0.43	0.085	0.385	0.125	0.405	0.03	0.55
Avg	0.063	0.588	0.135	0.479	0.16	0.476	0.066	0.644

Table 3:  $EM_{llm}$  of different models for Correct answers and specific Incorrect answers.

Answer Type	NQ	WebQ	TriviaQA	PopQA
Right	0.740	0.568	0.836	0.529
Wrong	-0.574	-0.389	-0.385	-0.121

Table 4: Evaluation of the Average Pearson Correlation Coefficient between EM and  $EM_{llm}$  for Correct and Incorrect Answers.

a and is validated by GPT-3.5-Turbo as being sup-
portive of $a$ . The EM <sub>llm</sub> for the correct answers
and specific incorrect answers, as generated by the
CTRLGEN method containing misinformation, are
presented in Table 3.

1102

1103

1104

1105

1106

1107

1108

1109

1110

1111

1112

1113

1114

1115

1116

1117

1118

1119

1120

1121

1122

1123

1124

1125

1126

1127

1128

1129

1130

1131

1132

1133

Moreover, we substantiate the rationality of employing EM as a QA evaluation metric in the experiments of Section 4 by calculating the Pearson correlation coefficient between  $EM_{llm}$  and EM based on the experimental results in this section. We observe that the  $EM_{llm}$  values for the texts generated by the other four LLMs, as verified by GPT-3.5-Turbo, have an average correlation exceeding 0.5 across four datasets when compared to the EM values obtained through direct string matching, as shown in Table 4. This demonstrates a significant correlation between the two metrics. For incorrect answers, the correlation is relatively lower, indicating the necessity of using GPT-3.5-Turbo to further filter texts in the exploration of misinformation. Considering the higher efficiency of direct EM calculation over  $EM_{llm}$ , we use EM to evaluate the QA quality of the RAG system for experiments in other sections.

Analysis of Experimental Results: From the experimental results, we can observe that: 1) the "Spiral of Silence" still exists. We first investigate the presence of the "Spiral of Silence" phenomenon when misinformation targeting the objective is injected into the corpus. As shown in Table 5, although the majority of the injected information is misleading, the content generated by the LLMs is still quickly ranked at the top by the retrieval systems, taking a dominant position. When com-

Method	NQ	WebQ	TriviaQA	PopQA
BM25	26.7	17.7	24.7	47.4
Contriever	60.7	62.5	64.7	67.2
LLM-Embedder	65.8	70.4	73.3	74.2
$BGE_{base}$	50.7	48.6	63.2	60.6

Table 5: Percentage of LLM-generated documents with **Misinformation** occupying the top-5 retrieval results, after augmenting each query with five documents generated by LLMs. Data entries framed by a **blue** background indicate a majority presence of human-generated documents, while entries with a **purple** background denote a predominance of LLM-generated documents.

paring four different retrieval methods, the BM25 1134 algorithm shows greater robustness than the others, 1135 being least affected by the LLM-generated content. 1136 However, it is noteworthy that approximately 20% 1137 of the content generated by the LLM could still be 1138 quickly placed in the forefront of the search results 1139 by the BM25 algorithm. Figure 11 illustrates that 1140 over time, human-written texts are gradually ex-1141 cluded from the searchable range, and as depicted 1142 in Figure 12, the phenomenon of homogenization 1143 of opinions in search results persists. This fur-1144 ther indicates that regardless of the accuracy of the 1145 LLM-generated information, the spiral of silence 1146 phenomenon remains present. 2) the RAG system 1147 has a limited degree of self-correction capability. 1148 LLM-generated texts containing misinformation 1149 lead to a significant decline in retrieval and QA 1150 performance based on the RAG system compared 1151 to results in Section 5.2. With the continuation 1152 of the iteration process, the number of correct an-1153 swers increased and the number of incorrect an-1154 swers decreased, albeit by a small margin. This 1155 suggests that the RAG system has a certain degree 1156 of self-correction capability, which may stem from 1157 the model's own knowledge or the human-written 1158 texts containing correct information retrieved in 1159 the initial stages. 3) The introduction of a small 1160 amount of the LLM-generated texts with spe-1161



Figure 11: Average percentage of texts from various sources within the top 50 search results over multiple iterations when adding **Misinformation** across different search methods.



Figure 12: Correlation between the number of top 5 search results containing the correct answer ("Context Right Num") and the accuracy of responses given by LLMs on the NQ dataset when adding **Misinformation**. The responses are categorized based on  $EM_{llm}$  score:  $EM_L=1$  for correct and  $EM_L=0$  for incorrect. The overall number of queries that the LLMs answered correctly ( $EM_L=1$  Total) and incorrectly ( $EM_L=0$  Total), along with the average retrieval accuracy (Acc@5) are shown by dashed lines. The results are averaged across different LLMs, retrieval and ranking methods.



Figure 13: The average EM<sub>llm</sub> scores of the correct and incorrect answers of the RAG system in the simulation across datasets and LLMs.

cific misleading information during the iterative 1162 process could inject such information into the 1163 **RAG output.** In Figure 13, we quantify the  $EM_{llm}$ 1164 metric of the original and specific misleading an-1165 swers generated by the RAG system based on four 1166 retrieval methods at various iteration stages. The 1167 results show that after the purposeful addition of 1168 misleading information (before the first iteration), 1169 the proportion of RAG system-generated answers 1170 containing specific misleading information signifi-1171 cantly increases, especially on the NQ and PopQA 1172 datasets, where the proportion of incorrect answers 1173 exceeds that of correct ones, and the influence of 1174 misleading answers persisted over time. However, 1175 the BM25 algorithm exhibits relatively higher ro-1176 bustness, and the EM<sub>llm</sub> of incorrect answers out-1177 put by the RAG system based on it remains lower 1178 than the other three retrieval methods. The exper-1179 imental results of this section reveal that despite 1180 the presence of self-correcting mechanisms, the in-1181 jection of specific misleading information can still 1182 severely compromise the system's accuracy and 1183 enable the manipulation of the RAG system to con-1184 sistently output specific misinformation in response 1185 1186 to certain questions. Therefore, without timely intervention, the spiral of silence phenomenon could 1187 marginalize accurate information, leading to severe 1188 misinformation consequences. 1189

#### A.7 Attempts to Alleviate the Spiral of Silence

1190

1191

1192

1193

1194

1195

1196

1197

1198

1199

1200

1201

1204

1205

1207

1208

1209

1210

1211

The "Spiral of Silence" effect could lead to the marginalization of human-generated text expression and further enhance the homogeneity of retrieval outcomes. If left unaddressed, this phenomenon could precipitate a series of adverse repercussions. To mitigate or eliminate the influence of the "Spiral of Silence" effect, this section initiates a discussion on two fronts. First, from the perspective of the authenticity of sources, we employ the widely used AIGC detection technologies to filter out and exclude all non-human-produced texts at the top of the search results. Second, addressing the validity of content, we strive to maintain diversity among the top search results to overcome potential issues caused by excessive homogenization.

**Experimental Setup**: To balance the efficiency and effectiveness of the retrieval system, for each set of search results returned by the system, we post-process to acquire the top-5 qualifying documents that are visible to the LLMs. In the **source filtering** experiment, we employ the HelloSimpleAI/chatgpt-qa-detector-roberta<sup>5</sup> model to 1212 authenticate the origins of the texts within the 1213 search results, aiming to retain the first 5 docu-1214 ments identified as human-generated and supply 1215 them as input to the LLM's context. For the con-1216 tent filtering part of the experiment, we apply a 1217 selection process based on computing the 3-gram 1218 Self-BLEU scores. The specific procedure is as 1219 follows: For the top-5 documents returned for each 1220 search query, we initially calculate their Self-BLEU 1221 scores; if the score exceeds a predetermined thresh-1222 old (set at 0.4 for this experiment), we then com-1223 pute the Self-BLEU scores for all possible combi-1224 nations of 4 documents and select the minimum 1225 value among them. This minimum value indicates 1226 the maximum individual document contribution 1227 to the Self-BLEU score not included in the calcu-1228 lation. Subsequently, we exclude the document 1229 contributing the most to the Self-BLEU score and 1230 incorporate the next ranked document into the com-1231 bination, repeating this filtering process until the 1232 combination's Self-BLEU score meets the preset 1233 threshold criteria. 1234

1235

1236

1237

1238

1239

1240

1241

1242

1243

1244

1245

1246

1247

1248

1249

1250

1251

1252

1253

1254

1255

1256

1257

1258

1259

1260

Analysis of Experimental Results: From the experimental results, we can observe that: 1) Both approaches yield more stable retrieval outcomes; however, the source filtering method incurs a performance cost. Figure 14 and Figure 15 illustrate the variations in the average retrieval outcomes and QA performance across datasets, before and after the application of two distinct filtering strategies, compared to an unfiltered condition. Observations indicate that by implementing sourcebased and diversity-based filtering methods, the fluctuation range of the top-5 retrieval results is reduced compared to the non-intervention scenario, suggesting that the filtering mechanisms can bring a more stable retrieval performance for RAG systems. Across the four datasets, the retrieval performance following SELF-BLEU value filtering generally surpasses the unfiltered condition; conversely, the source-based filtering strategy results in an overall performance degradation. This could be attributed to the discriminating model erroneously excluding valid human-generated texts while aiming to eliminate those generated by LLMs. Moreover, in QA tasks, diversity filtering either enhances or maintains QA performance, whereas source-based document filtering leads to a decline in QA per-

<sup>&</sup>lt;sup>5</sup>https://github.com/Hello-SimpleAI/ chatgpt-comparison-detection



Figure 14: Average long-term retrieval performance of different filtering strategies.



Figure 15: Average long-term QA performance of different filtering strategies.

formance across all datasets. For instance, on the TriviaQA dataset, the average EM score drops by 1262 over 14%. 2) Both methods can only alleviate 1263 the "Spiral of Silence" phenomenon to varying 1264 degrees and cannot eliminate it. Figure 16 dis-1265 plays the proportion of documents from different 1266 sources within the top-5 retrieval results in each 1267 iteration under three filtering setups on the NO 1268 dataset. It is observable that without any filter-1269 ing strategy, human-generated texts rapidly vanish 1270 from the top-5 documents in the initial iterations. 1271 The SELF-BLEU value filtering method retains 1979 human-generated texts to a small extent; source 1273 1274 filtering, on the other hand, maximally filters out LLM-generated texts, especially those produced 1275 by GPT-3.5-Turbo, Qwen, and ChatGLM3, with over 30% of human-generated texts remaining in 1277 the top-5 by the end of the tenth iteration. However, 1278 despite both filtering strategies slowing the disap-1279 pearance of human texts, the proportion of human-1280 generated content continues to exhibit a declining 1281 trend. Figure 17 and Figure 18 demonstrate that 1282 compared to the absence of filtering strategies, both 1283 filtering methods slow down the polarization speed of top document accuracy in retrieval performance. 1285 Overall, we discovered that filtering based on the 1286 source of documents and their diversity can, to 1287 some extent, slow down the emergence of the Spiral of Silence phenomenon. Source-based filtering has a more pronounced effect in terms of preserv-1290

ing the proportion of human-generated texts and 1291 mitigating viewpoint polarization; however, this 1292 benefit comes at the expense of the performance of 1293 the RAG system. Text filtering based on diversity 1294 shows superior performance in maintaining RAG 1295 system functionality, but it has a weaker impact on 1296 preserving the ratio of human texts and alleviating 1297 viewpoint polarization. Despite these findings, nei-1298 ther method can completely eradicate the "Spiral of 1299 Silence" effect, indicating the imperative to explore 1300 additional solutions. For example, there is a need 1301 to investigate retrieval models that can effectively 1302 balance between LLM-generated documents and 1303 human-generated documents to address this issue. 1304

#### A.8 Prompts

The prompts used in the experiment are shown in1307Table 6.1308

1305



Figure 16: Average percentage of texts from various sources within the top 5 search results over multiple iterations on NQ when using different filtering strategies across different search methods.



Figure 17: Correlation between the number of top 5 search results containing the correct answer ("Context Right Num") and the accuracy of responses given by LLMs on the NQ dataset when using **Content Filtering**. The responses are categorized based on Exact Match (EM) score: EM=1 for correct and EM=0 for incorrect. The overall number of queries that the LLMs answered correctly (EM=1 Total) and incorrectly (EM=0 Total), along with the average retrieval accuracy (Acc@5) are shown by dashed lines. The results are averaged across different LLMs, retrieval and ranking methods.



Figure 18: Correlation between the number of top 5 search results containing the correct answer ("Context Right Num") and the accuracy of responses given by LLMs on the NQ dataset when using **Source Filtering**. The responses are categorized based on Exact Match (EM) score: EM=1 for correct and EM=0 for incorrect. The overall number of queries that the LLMs answered correctly (EM=1 Total) and incorrectly (EM=0 Total), along with the average retrieval accuracy (Acc@5) are shown by dashed lines. The results are averaged across different LLMs, retrieval and ranking methods.

Prompts
Provide a background document in 100 words
according to your knowledge from Wikipedia to
answer the given question.
Question:{question_str}
Background Document:
Context information is below.
[Context 1]:{context str1}
[Context 5]:{context_str5}
Using both the context information and also
using your own knowledge, answer the following
question with a background document in 100
words.
Question: {question_str}
Packground Document.
Concrete a false answer to the given question
It should be short (less than ten words in
general) and look plausible, compared to the
reference answer.
Question:{question_str}
Reference Answers:{ <b>ref_answer_str</b> },
False answer:
Suppose you are generating realistic-looking
a background document loss than 100 words in
a background document less than 100 words in
Don't contain any word in the original answers
in {ref answer str}. The background document
must contain the following given answers with
their original form.
Question:{question_str}
Answers:{ <b>false_answer_str</b> },
Dackyround document:
to the question?
Ouestion: {auestion str}
Response: {response str}
Answer: {ref answer str} / {false answer str}
Just answer 'yes' or 'no'.

Table 6: Prompts for different tasks.