

Comparative Analysis of RAG vs. Fine-Tuning for Unanswerable Question Detection

Anonymous ACL submission

Abstract

Lore ipsum dolor sit amet, consectetur adipisciing elit. Cras ullamcorper elementum metus, eget dignissim felis posuere vitae. Nullam bibendum tellus sed bibendum gravida. Fusce posuere laoreet tortor, at facilisis magna imperdier vel. Donec finibus, ipsum a finibus ornare, leo risus dictum felis, sit amet porta diam diam sit amet lectus. Aenean aliquam condimentum velit. Quisque sodales volutpat ligula, non gravida metus tincidunt in. Sed fermentum eros ac libero gravida, eu hendrerit enim porttitor. Fusce vehicula dignissim vehicula. Aenean tincidunt posuere lobortis. Nam sodales nunc id ipsum ullamcorper elementum. Nam molestie velit quis elit blandit placerat.

1 Introduction

Question answering (QA) represents a primary use of large language models (LLMs), especially in domains that require accurate information. However, QA systems face a fundamental challenge that becomes particularly significant in high-stakes applications: they often hallucinate answers to unanswerable questions instead of abstaining (Kalai et al., 2025).

Healthcare and legal domains are particularly critical. QA systems for consumer health queries, clinical decision support, and drug interaction checking must never fabricate information, as hallucinated medical advice can have life-threatening consequences (Pal et al., 2023). Similarly, systems for legal research, compliance checking, and contract analysis can cause penalties, liability, or rights violations through fabricated answers. Prominent scandals have occurred where AI systems invented non-existent legal cases, resulting in attorney sanctions (Dahl et al., 2024). These domains particularly need robust unanswerable question detection because knowledge bases have clear boundaries yet cannot cover every scenario. Medical guidelines

cannot address every rare disease combination; legal databases cannot cover every novel situation. When questions fall outside a system's knowledge scope, it must recognize this rather than fabricate answers.

Two prominent paradigms address the hallucination challenge through different approaches to knowledge integration. Retrieval-Augmented Generation (RAG) retrieves relevant passages from a knowledge base before generation, representing an external, non-parametric approach to grounding responses (Lewis et al., 2020; Shuster et al., 2021). Fine-tuning directly encodes knowledge into model parameters through training on domain-specific question-answer pairs, including unanswerable examples, representing an internal, parametric approach (Roberts et al., 2020). While other mitigation techniques exist - including prompt engineering, chain-of-thought reasoning (Wei et al., 2022), and human-in-the-loop verification (Ouyang et al., 2022) - RAG and fine-tuning represent the two fundamental paradigms for integrating domain knowledge into QA systems.

Organizations deploying QA systems must choose between RAG and fine-tuning for knowledge integration, yet lack empirical guidance on which approach better handles unanswerable questions. Specifically, we address: *Do models more reliably abstain from answering when knowledge is provided externally through retrieval (RAG) or encoded internally through fine-tuning? How does each paradigm balance answering answerable questions correctly while recognizing when questions fall outside the knowledge scope?* These questions are critical for high-stakes deployments where incorrect answers carry significant consequences.

While prior work has compared RAG and fine-tuning on standard QA metrics like exact match and F1 scores, these comparisons focus primarily on accuracy when answers exist. Little empirical research examines how these paradigms differ in their

ability to detect unanswerable questions and appropriately abstain (Soudani et al., 2024; Balaguer et al., 2024). Furthermore, existing hallucination detection work focuses on identifying fabricated content after generation, rather than comparing prevention strategies through different knowledge integration approaches (Sadat et al., 2023; Farquhar et al., 2024). Finally, studies on SQuAD 2.0 - a dataset explicitly designed to test abstention behavior - concentrate on architectural improvements to extractive QA models rather than comparing fundamental knowledge integration paradigms (Rajpurkar et al., 2018). Our work fills this gap through systematic evaluation of RAG versus fine-tuning specifically on unanswerable question detection.

We present a systematic comparison of RAG and fine-tuning paradigms specifically focused on unanswerable question detection using SQuAD 2.0. Beyond standard metrics, we introduce an answer attribution score that measures whether generated answers are grounded in provided/retrieved context versus hallucinated from parametric memory. This metric directly captures the key distinction between external and internal knowledge integration. We implement three systems using Llama models: (1) zero-shot baseline with no retrieval or fine-tuning, (2) RAG system with dense retrieval from the SQuAD 2.0 corpus, and (3) Llama fine-tuned on SQuAD 2.0 using QLoRA. Our controlled experimental design isolates the effect of knowledge integration paradigm on abstention behavior while maintaining comparable model capacity and training data.

Our evaluation reveals distinct tradeoffs between RAG and fine-tuning for unanswerable question detection...

2 Background

Question answering systems have evolved significantly with the introduction of large-scale benchmarks. Rajpurkar et al. (2016) introduced the Stanford Question Answering Dataset (SQuAD 1.0), a reading comprehension benchmark containing over 100,000 questions where all answers could be extracted from provided Wikipedia passages. However, SQuAD 1.0 had a critical limitation: every question was guaranteed to be answerable from the given context, meaning systems never needed to recognize when they lacked sufficient information. To address this, Rajpurkar et al. (2018) released SQuAD 2.0, augmenting the dataset with

over 50,000 unanswerable questions. This established the paradigm of evaluating not just answer accuracy, but also the crucial ability to recognize knowledge boundaries—a capability we leverage to evaluate how different knowledge integration paradigms handle unanswerable questions.

Early approaches to unanswerable question detection in SQuAD 2.0 focused on architectural modifications to extractive models. These methods typically added answer verification components, threshold-based mechanisms, or modified BERT-based architectures to output “no answer” predictions when appropriate. However, this line of work concentrated on improving specific model architectures for extractive QA rather than examining fundamental differences in how knowledge is integrated into systems—the core question we address for modern generative LLMs.

3 Methods

Nunc dolor eros, lobortis egestas risus quis, facilisis ornare mi. Morbi a lacinia sapien. Etiam risus purus, lobortis vitae purus at, commodo condimentum tortor. Nulla diam metus, egestas vitae interdum ut, pretium sit amet nisl. Nam eget malesuada dui. Cras ultrices, felis aliquam ultricies iaculis, massa dolor semper turpis, quis placerat erat tortor et velit. Vestibulum posuere aliquam neque ut dignissim. Quisque mauris justo, euismod at convallis eget, aliquet eu lacus. Sed ut luctus neque. Praesent vel mauris sit amet odio ornare commodo. Maecenas lacinia iaculis magna, id porttitor quam tincidunt vitae. Vivamus et vehicula quam. Mauris malesuada euismod turpis, ut pellentesque tortor condimentum porttitor. Maecenas viverra lorem libero, in tristique ex iaculis non. Duis tempus lobortis lorem sed auctor.

4 Results and Discussion

Donec ut posuere quam. Nam eget dapibus erat, gravida viverra felis. Fusce eleifend urna sed risus suscipit, in vestibulum elit luctus. Vivamus ut leo convallis, interdum magna et, accumsan odio. Aenean vulputate purus at est sollicitudin, nec egestas justo iaculis. Fusce convallis sit amet nisi ac pharetra. Aenean maximus aliquam nibh eu ornare. Duis sit amet dictum magna, id malesuada urna. Curabitur cursus gravida odio. Donec eleifend iaculis vestibulum. Aliquam suscipit, lectus id malesuada ultricies, lorem ipsum malesuada nisi, sed venenatis magna metus non ante. Nulla a porttiti-

181 tor mauris, et pretium justo. Nulla dignissim mi
182 non pharetra lacinia. Pellentesque habitant morbi
183 tristique senectus et netus et malesuada fames ac
184 turpis egestas. Interdum et malesuada fames ac
185 ante ipsum primis in faucibus.

186 5 Conclusion

187 Cras cursus condimentum semper. Donec rhoncus
188 ante ac mi aliquam rutrum. Nullam pulvinar felis
189 vitae tellus porttitor, non vehicula ante consectetur.
190 Fusce accumsan tortor at nunc finibus, non male-
191 suada felis cursus. Nulla malesuada mattis mollis.
192 Donec finibus consequat tincidunt. Quisque me-
193 tus risus, blandit dictum mollis ut, molestie eget
194 ex. Suspendisse mauris erat, ultrices vitae euismod
195 ac, varius tempor nulla. Vivamus pellentesque nisl
196 eu venenatis aliquam. Donec tincidunt ut lectus
197 nec interdum. Nam eu libero luctus, dapibus odio
198 ac, gravida quam. Quisque blandit quis velit vel
199 consequat.

200 6 Author’s Contributions

201 Fusce pretium magna mauris. Cras mi elit, ve-
202 nenatis id turpis quis, vehicula ornare turpis. Vi-
203 vamus molestie aliquet efficitur. Vestibulum ac
204 massa justo. Sed vulputate, sem eu sollicitudin
205 volutpat, erat metus tempor ante, in hendrerit ante
206 diam a massa. Proin in fermentum libero. Donec
207 nec quam a nulla tincidunt bibendum. Proin vitae
208 lobortis odio. Vestibulum ut justo est. Pellentesque
209 habitant morbi tristique senectus et netus et male-
210 suada fames ac turpis egestas. Suspendisse vitae
211 tincidunt nunc. Integer a urna sed massa pulvinar
212 bibendum vitae ultricies est. Ut lacinia dolor purus,
213 eget accumsan mauris tempor posuere. Morbi blan-
214 dit ipsum elementum, fringilla ligula ac, viverra
215 nibh.

216 References

217 Angels Balaguer, Vinamra Benara, Renato Luiz de Fre-
218 itas Cunha, Todd Hendry, Daniel Holstein, Jen-
219 nifer Marsman, Nick Mecklenburg, Sara Malvar,
220 Leonardo O Nunes, Rafael Padilha, and 1 others.
221 2024. Rag vs fine-tuning: pipelines, tradeoffs,
222 and a case study on agriculture. *arXiv preprint arXiv:2401.08406*.

224 Matthew Dahl, Varun Magesh, Mirac Suzgun, and
225 Daniel E Ho. 2024. Large legal fictions: Profiling le-
226 gal hallucinations in large language models. *Journal
227 of Legal Analysis*, 16(1):64–93.

- 228 Sebastian Farquhar, Jannik Kossen, Lorenz Kuhn, and
229 Yarin Gal. 2024. Detecting hallucinations in large
230 language models using semantic entropy. *Nature*,
231 630(8017):625–630.
- 232 Adam Tauman Kalai, Ofir Nachum, Santosh S. Vem-
233 pala, and Edwin Zhang. 2025. Why language models
234 hallucinate. *ArXiv*, abs/2509.04664.
- 235 Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio
236 Petroni, Vladimir Karpukhin, Naman Goyal, Hein-
237 rich Küttler, Mike Lewis, Wen-tau Yih, Tim Rock-
238 täschel, Sebastian Riedel, and Douwe Kiela. 2020.
239 Retrieval-augmented generation for knowledge-
240 intensive nlp tasks. In *Advances in Neural Infor-
241 mation Processing Systems*, volume 33, pages 9459–
242 9474. Curran Associates, Inc.
- 243 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida,
244 Carroll Wainwright, Pamela Mishkin, Chong Zhang,
245 Sandhini Agarwal, Katarina Slama, Alex Ray, John
246 Schulman, Jacob Hilton, Fraser Kelton, Luke Miller,
247 Maddie Simens, Amanda Askell, Peter Welinder,
248 Paul F Christiano, Jan Leike, and Ryan Lowe. 2022.
249 Training language models to follow instructions with
250 human feedback. In *Advances in Neural Information
251 Processing Systems*, volume 35, pages 27730–27744.
252 Curran Associates, Inc.
- 253 Ankit Pal, Logesh Kumar Umapathi, and Malaikannan
254 Sankarasubbu. 2023. Med-HALT: Medical domain
255 hallucination test for large language models. In *Pro-
256 ceedings of the 27th Conference on Computational
257 Natural Language Learning (CoNLL)*, pages 314–
258 334, Singapore. Association for Computational Lin-
259 guistics.
- 260 Pranav Rajpurkar, Robin Jia, and Percy Liang. 2018.
261 Know what you don’t know: Unanswerable ques-
262 tions for SQuAD. In *Proceedings of the 56th Annual
263 Meeting of the Association for Computational Lin-
264 guistics (Volume 2: Short Papers)*, pages 784–789,
265 Melbourne, Australia. Association for Computational
266 Linguistics.
- 267 Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and
268 Percy Liang. 2016. SQuAD: 100,000+ questions for
269 machine comprehension of text. In *Proceedings of
270 the 2016 Conference on Empirical Methods in Natu-
271 ral Language Processing*, pages 2383–2392, Austin,
272 Texas. Association for Computational Linguistics.
- 273 Adam Roberts, Colin Raffel, and Noam Shazeer. 2020.
274 How much knowledge can you pack into the parame-
275 ters of a language model? In *Proceedings of the
276 2020 Conference on Empirical Methods in Natural
277 Language Processing (EMNLP)*, pages 5418–5426,
278 Online. Association for Computational Linguistics.
- 279 Mobashir Sadat, Zhengyu Zhou, Lukas Lange, Jun
280 Araki, Arsalan Gundroo, Bingqing Wang, Rakesh
281 Menon, Md Parvez, and Zhe Feng. 2023. Delu-
282 cionQA: Detecting hallucinations in domain-specific
283 question answering. In *Findings of the Association
284 for Computational Linguistics: EMNLP 2023*, pages

285 822–835, Singapore. Association for Computational
286 Linguistics.

287 Kurt Shuster, Spencer Poff, Moya Chen, Douwe Kiela,
288 and Jason Weston. 2021. [Retrieval augmentation](#)
289 [reduces hallucination in conversation](#). In *Findings*
290 [of the Association for Computational Linguistics:](#)
291 *EMNLP 2021*, pages 3784–3803, Punta Cana, Do-
292 minican Republic. Association for Computational
293 Linguistics.

294 Heydar Soudani, Evangelos Kanoulas, and Faegheh Ha-
295 sibi. 2024. [Fine tuning vs. retrieval augmented gener-](#)
296 [ation for less popular knowledge](#). In *Proceedings of*
297 [the 2024 Annual International ACM SIGIR Confer-](#)
298 [ence on Research and Development in Information](#)
299 [Retrieval in the Asia Pacific Region, SIGIR-AP 2024](#),
300 page 12–22, New York, NY, USA. Association for
301 Computing Machinery.

302 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten
303 Bosma, brian ichter, Fei Xia, Ed Chi, Quoc V Le,
304 and Denny Zhou. 2022. [Chain-of-thought prompt-](#)
305 [ing elicits reasoning in large language models](#). In
306 *Advances in Neural Information Processing Systems*,
307 volume 35, pages 24824–24837. Curran Associates,
308 Inc.

309 **A First Appendix Title**

310 This is appendix A content.

311 **B Second Appendix Title**

312 This is appendix B content.

313 **C Third Appendix Title**

314 This is appendix C content.