



RAG AI Agent: A Local, Privacy-Preserving Document Q&A System Using Retrieval-Augmented Generation

1st Author - Mr. Rahul Baburao Virale, 2nd Author - Mr. Omkar Bajarang Bhosale,
3rd Author - Mr. Pratik Zunjarrao Patil, 4th Author - Mr. Hanuman Sandeepan Wandare, 5th Author -
Mr. Umesh Anandrao Patil.

Guide of Research paper - Prof. U. A. Patil

Department of Computer Science Engineering

D.Y. Patil Technical Campus Talsande, Kolhapur, Maharashtra, India

Abstract — Large language models can write fluently but they cannot tell you what is actually in your PDF. They hallucinate, confabulate, and confidently miss the point. This paper describes a Retrieval-Augmented Generation (RAG) system built to close that gap.

Objective: Large language models cannot read documents they have never seen — they hallucinate when asked, or simply refuse. This project set out to build a local, privacy-preserving document question-answering system where every response is grounded in text retrieved directly from the user's own PDFs. The specific objectives were: (1) implement a full Retrieval-Augmented Generation pipeline running entirely on local hardware with no cloud API dependency; (2) support multi-turn conversational memory so users can ask follow-up questions naturally; (3) add abstractive document summarization; and (4) measure faithfulness and retrieval quality against a plain LLM and a BM25 keyword-search baseline.

Methodology: Uploaded PDFs are parsed with PyPDF Loader, split into 1,000-character overlapping chunks, and embedded with Nomic-embed-text (768 dimensions) via Ollama. Vectors are stored persistently in ChromaDB. At query time, the user's question is embedded with the same model and a cosine similarity search retrieves the top-5 most relevant chunks, which are injected into a structured prompt for Mistral 7B. LangChain's RetrievalQA chain orchestrates the pipeline; ConversationBufferMemory persists session history for multi-turn dialogue. A MapReduceDocumentsChain handles long-document summarization. The frontend is built with Next.js 15 and React 19, styled with Tailwind CSS and Shadcn UI components, communicating with a Python Flask REST API backend.

Results: Evaluated on 50 annotated query-document pairs across five document types, the system achieved a RAGAS faithfulness score of 0.91, Precision@5 of 0.82, and a hallucination rate of approximately 7% — compared to 42% for a plain Mistral baseline and 18% for BM25 retrieval. Out-of-scope queries were correctly refused in 94% of cases. Multi-turn memory resolved 93% of pronoun-reference follow-up questions (vs. 33% without memory). Average query latency was 5.2 seconds for cached documents. The system handled up to 5 concurrent users with no errors; degradation became significant beyond 10 users due to sequential GPU inference.

Conclusions: A locally-hosted RAG system using open-source models achieves production-quality faithfulness for document question answering without sending any user data to external servers. The 57-percentage-point reduction in hallucination over a plain LLM baseline confirms that retrieval-augmented prompting, combined with an explicit non-fabrication instruction, reliably grounds model responses in document content. The architecture is modular and extensible — the clearest next steps are OCR support for scanned PDFs, hybrid dense-sparse retrieval, and batched inference for multi-user deployments.

Keywords: Retrieval-Augmented Generation, RAG, Large Language Models, Vector Database, ChromaDB, Mistral, LangChain, Document Question Answering, Semantic Search, Nomic Embeddings, Flask, Next.js

I. INTRODUCTION

1.1 Problem Statement

Most people interact with large language models expecting them to know things. They do, to a point. But ask a model about the contents of a document it has never seen and you get either a hallucination or a polite admission of ignorance. Neither is useful. For anyone who works with PDFs — researchers, students, legal professionals, engineers — this is a hard ceiling.

The core issue is that LLMs store knowledge implicitly in their weights, fixed at training time. They cannot read new documents. Feeding a PDF into a chat prompt is a workaround, not a solution: most models have context windows too small for long documents, and even those with large windows tend to lose precision when the relevant passage is buried deep in a 60-page file.

1.2 Why Existing Tools Fall Short

Commercial tools like ChatGPT with file upload or Claude's document mode solve part of the problem but require sending potentially sensitive documents to third-party servers. For institutional, legal, or personal data, that trade-off is often unacceptable. Open-source alternatives exist but are either too complex to deploy or rely on external API calls that reintroduce the privacy problem.

There is a gap between what users need (ask questions about my documents, privately, on my own machine) and what current tools reliably provide.

1.3 Retrieval-Augmented Generation as the Answer

RAG [1] addresses this by separating knowledge storage from language generation. Documents are chunked, embedded into vectors, and stored in a vector database. At query time, the system retrieves the most semantically relevant chunks and injects them directly into the model's prompt. The model never has to recall anything from its weights — it reads the context provided to it.

This approach has two practical consequences. First, the model can answer questions about documents it was never trained on. Second, because the answer is derived from retrieved text, it is far easier to verify: the source chunks can be shown alongside the response.

1.4 Research Gaps Addressed

Three specific gaps shaped the design of this system:

- Gap 1 — Local deployment: Most published RAG implementations rely on OpenAI or Cohere APIs. Running Mistral 7B and Nomic-embed-text fully through Ollama eliminates cloud dependencies and keeps all document data on the user's machine.
- Gap 2 — Conversational context: Most RAG demos treat each query independently. Real users ask follow-up questions. This system persists conversation history across a session via LangChain's ConversationBufferMemory.
- Gap 3 — Accessible evaluation: RAG papers rarely report faithfulness alongside latency in a unified evaluation. This work uses the RAGAS framework to report both retrieval quality and generation faithfulness against a plain LLM baseline.

1.5 Contributions

- A fully local, open-source RAG pipeline using Mistral 7B and Nomic-embed-text via Ollama
- Multi-turn conversational memory integrated into the retrieval-generation chain

- An abstractive summarization module using MapReduce chaining in LangChain
- A responsive Next.js/React frontend with Shadcn UI components for document management and chat
- Evaluation against plain LLM and BM25 baselines, showing a 57% improvement in faithfulness

1.6 Paper Organization

Section II surveys prior work. Section III states the research questions. Section IV describes the system architecture and methodology. Section V presents evaluation results. Section VI discusses findings. Sections VII and VIII cover future scope and conclusions.

II. LITERATURE REVIEW

2.1 Retrieval-Augmented Generation

Lewis et al. [1] introduced RAG as a framework combining parametric memory (the language model) with non-parametric memory (a retrieved document store). Their original implementation used a Dense Passage Retrieval (DPR) encoder and BART for generation. The key insight was that retrieval at inference time, rather than at training time, allows a fixed model to remain current and domain-adaptive. Subsequent work by Izacard and Grave [2] showed that re-ranking retrieved passages before generation further improves factual precision.

More recent surveys [3] categorize RAG systems along three axes: naive RAG (retrieve-then-generate), advanced RAG (with re-ranking and query rewriting), and modular RAG (with tool-augmented generation). Our system sits in the advanced category, using cosine similarity retrieval from a persistent vector store with overlapping chunk windows.

2.2 Dense Retrieval and Vector Stores

Karpukhin et al. [4] demonstrated that dense bi-encoder retrieval consistently outperforms BM25 for open-domain question answering. The gap is especially pronounced on paraphrase-heavy queries where keyword overlap is low. Embedding models like all-MiniLM-L6-v2 [5] and, more recently, Nomic-embed-text [6] provide sentence-level representations that capture semantic similarity effectively on consumer hardware.

ChromaDB [7], an open-source vector database, provides persistent storage and cosine similarity search with simple Python bindings. Unlike FAISS, ChromaDB is designed for persistent use and does not require rebuilding the index on each startup, which matters for a document Q&A application where the corpus changes incrementally.

2.3 LLMs for Document Understanding

Mistral 7B [8] achieves strong performance on reading comprehension benchmarks at a fraction of the memory footprint of GPT-4 class models. It runs comfortably on consumer GPUs (8GB+ VRAM) via Ollama, making local deployment practical. Recent work by Touvron et al. [9] on Llama 2 and by Jiang et al. [8] on Mistral both suggest that smaller, instruction-tuned models can match larger models on focused, context-provided tasks — exactly the setting of a RAG system.

2.4 Hallucination and Faithfulness in LLMs

Ji et al. [10] provide a comprehensive taxonomy of hallucination in neural text generation. They distinguish between intrinsic hallucination (contradicting the source) and extrinsic hallucination (adding unverifiable information). RAG reduces both by anchoring generation to retrieved text, but does not eliminate them. The model can still misread its context or selectively ignore contradictory chunks.

The RAGAS evaluation framework [11] was developed specifically to assess RAG pipelines on four dimensions: answer relevance, faithfulness, context recall, and context precision. We use faithfulness and answer relevance in our evaluation.

2.5 Conversational Memory in LLMs

LangChain [12] provides several memory primitives for multi-turn dialogue. ConversationBufferMemory stores the full conversation history and prepends it to each new prompt. For long sessions, this grows unwieldy, but for the typical document Q&A use case (5-10 turns), it is sufficient. Alternatives like ConversationSummaryMemory [12] periodically compress the history using the LLM itself, trading accuracy for token efficiency.

2.6 Related Systems

Gao et al. [13] describe a survey-level comparison of RAG vs. fine-tuning vs. in-context learning for domain adaptation. Their conclusion is that RAG is preferable when the knowledge domain changes frequently and fine-tuning is impractical. PrivateGPT [14] and LlamaIndex [15] both implement local RAG pipelines but target developer use cases rather than end-users. Our system prioritizes the end-user experience, wrapping the pipeline in a conversational chat interface.

III. RESEARCH OBJECTIVES AND RESEARCH QUESTIONS

Six research questions guide this work:

RQ1 (Retrieval Effectiveness): What Precision@5 and MRR can dense cosine similarity retrieval from ChromaDB achieve on a domain-general document corpus?

RQ2 (Answer Faithfulness): How does the faithfulness and hallucination rate of a RAG system compare to a plain LLM and a BM25 baseline on the same query set?

RQ3 (Conversational Context): Does maintaining a session-level conversation buffer produce measurably more coherent follow-up responses than stateless querying?

RQ4 (Latency): What end-to-end response latency does the system achieve, broken down by ingestion, retrieval, and generation stages?

RQ5 (Scalability): How does response time degrade under concurrent user load?

RQ6 (Summarization Quality): Can a MapReduce summarization chain produce coherent abstractive summaries of multi-page PDFs using a local 7B model?

IV. METHODOLOGY

4.1 System Architecture

The system is organized into three layers: a Next.js frontend, a Python Flask backend, and a data storage layer. Figure 1 shows the high-level flow.

The frontend provides a document upload panel and a conversational chat interface. Users drag-and-drop a PDF, which is sent to the Flask backend via a multipart POST request. Once ingested, the document appears in the document list and users can begin querying.

The backend handles two primary operations: document ingestion and query execution. These are separated deliberately, so the embedding computation happens offline (at upload time) and does not add to query latency.

The storage layer consists of two components: the local filesystem for raw PDFs and ChromaDB for vector embeddings and chunk metadata.

4.2 PDF Ingestion Pipeline (Module 3)

When a PDF is uploaded, the following steps execute in sequence:

- Text extraction: PyPDF Loader reads the document page by page and concatenates raw text.
- Chunking: LangChain's RecursiveCharacterTextSplitter divides the text into chunks of 1,000 characters with a 200-character overlap. Overlap prevents context loss at chunk boundaries.
- Embedding: Each chunk is passed to Nomic-embed-text via the Ollama API. This produces a 768-dimensional dense vector per chunk.
- Storage: The vectors, chunk text, and document metadata are stored in ChromaDB under a collection keyed by document ID.

4.3 RAG Query Pipeline (Module 4)

Query execution proceeds as follows. The user's query text is embedded using the same Nomic-embed-text model to produce a 768-dimensional query vector. ChromaDB performs cosine similarity search against the document collection and returns the top-5 most similar chunks. These chunks are assembled into a context string and injected into a structured prompt template that instructs Mistral to answer only

from the provided context and to acknowledge when the answer is not present. LangChain's RetrievalQA chain manages this orchestration.

The prompt template includes an explicit instruction: 'If the answer is not contained in the following context, say so. Do not invent information.' This single instruction accounts for a large part of the faithfulness improvement over the plain LLM baseline.

4.4 Conversational Memory (Module 5)

LangChain's ConversationBufferMemory stores the full message history for a session. On each new query, the memory is serialized and injected into the prompt as a 'Chat History' block. This allows the model to handle pronoun references ('what about the second point?') and follow-up questions without the user restating context.

Sessions are keyed by a UUID generated at page load. Different browser sessions maintain independent histories. Memory is cleared when the user uploads a new document.

4.5 Summarization (Module 6)

Document summarization uses LangChain's MapReduceDocumentsChain. The Map step sends each chunk to Mistral with a one-sentence summary instruction. The Reduce step takes the collection of per-chunk summaries and asks Mistral to synthesize them into a coherent 3-5 paragraph abstract. This avoids the context-window constraint that would prevent processing a 100-page document in a single prompt.

4.6 Frontend Architecture

The user interface is built with Next.js 15.4.2 using the App Router and React 19.1.0. TypeScript is used throughout. The Shadcn UI component library provides the visual primitives — cards, buttons, badges, scroll areas — styled with Tailwind CSS utility classes. The upload panel uses the HTML5 File API with drag-and-drop. Chat messages are rendered with markdown support to handle code blocks and lists in LLM responses.

Communication with the Flask backend uses the native Fetch API. The chat interface updates optimistically: the user's message appears immediately, and a loading indicator shows while the backend processes the query.

4.7 Technology Stack

Layer	Technology	Version	Role
Frontend	Next.js + React	15.4.2 / 19.1.0	UI and routing
Frontend	TypeScript + Tailwind CSS	—	Type safety and styling
Frontend	Shadcn UI	Latest	Component library
Backend	Python + Flask	3.10 / 3.0	REST API gateway
Backend	LangChain	0.2.x	RAG orchestration
AI Model	Mistral 7B (via Ollama)	0.1	Text generation
AI Model	Nomic-embed-text (via Ollama)	1.5	Dense embeddings
Storage	ChromaDB	0.5.x	Vector store
Storage	PyPDF Loader	—	PDF text extraction
Runtime	Ollama	0.1.x	Local model serving

V. RESULTS AND EVALUATION

5.1 RQ1: Retrieval Effectiveness

Retrieval was evaluated on a hand-annotated set of 50 query-document pairs drawn from five document types: research papers, technical manuals, policy reports, news articles, and student textbooks. Ground-truth relevance was determined by two annotators with agreement computed using Cohen's kappa ($\kappa = 0.81$, strong agreement).

Metric	RAG System	BM25 Baseline	Difference
Precision@5	0.82	0.61	+0.21
Recall@5	0.74	0.59	+0.15
Mean Reciprocal Rank	0.88	0.71	+0.17
NDCG@5	0.79	0.64	+0.15

The dense retrieval system outperforms BM25 on every metric. The gap is largest on Precision@5 (+21 percentage points), confirming the advantage of semantic over keyword matching, particularly on queries where the user's phrasing differs from the document's phrasing.

5.2 RQ2: Answer Faithfulness and Hallucination

Answer quality was assessed using two RAGAS metrics — faithfulness and answer relevance — plus a manual hallucination count by the annotators.

Evaluation Dimension	RAG System	LLM-Only (Mistral)	BM25 + Mistral
Faithfulness (RAGAS)	0.91	0.34	0.71
Answer Relevance (RAGAS)	0.85	0.51	0.67
Hallucination Rate	~7%	~42%	~18%
Out-of-context refusals (correct)	94%	N/A	79%

The RAG system reduces hallucination from 42% (plain LLM) to 7%, a 35-percentage-point drop. The explicit 'do not invent information' instruction in the prompt, combined with the retrieved context, is primarily responsible for this result. When queries fell outside the document's scope, the system correctly refused to answer in 94% of cases.

5.3 RQ3: Conversational Context

We designed 15 two-turn conversation scenarios where the second query used a pronoun or partial reference that required understanding the first. Examples: 'Who wrote this paper?' followed by 'What institution are they from?'. The RAG system with memory resolved 14 of 15 correctly (93%). Without memory (stateless), the system correctly resolved 5 of 15 (33%). This confirms that session-level memory is essential for natural multi-turn interaction.

5.4 RQ4: Response Latency

Stage / Document Size	Small (< 10 pages)	Medium (10–50 pages)	Large (50–100 pages)
PDF ingestion + chunking	0.4 sec	1.8 sec	3.9 sec
Embedding generation	0.8 sec	3.6 sec	7.8 sec
Retrieval from ChromaDB	0.3 sec	0.4 sec	0.5 sec
Mistral generation	4.8 sec	5.1 sec	5.3 sec
Total (query only, cached)	5.1 sec	5.5 sec	5.8 sec
Total (first upload + query)	6.3 sec	10.9 sec	17.5 sec

Retrieval from ChromaDB is consistently fast regardless of corpus size (< 0.5 sec). The ingestion bottleneck is embedding generation, which scales linearly with document length. Mistral generation time is stable at roughly 5 seconds regardless of document size, because the context provided to the model is always the same size (top-5 chunks of 1,000 characters each).

5.5 RQ5: Scalability Under Load

Concurrent Users	Avg. Response Time	95th Percentile	Error Rate
1	5.2 sec	6.1 sec	0%
5	7.4 sec	9.8 sec	0%
10	13.6 sec	18.2 sec	2%
25	31.4 sec	44.7 sec	8%

The system handles up to 5 concurrent users acceptably. Beyond 10, latency grows super-linearly because Ollama processes inference requests sequentially on the GPU. This is a hardware constraint, not a software architecture problem, but it defines the practical deployment limit of a single-machine setup.

5.6 RQ6: Summarization Quality

Abstractive summaries were generated for 10 documents ranging from 5 to 80 pages. Three evaluators rated each summary on coherence (1–5), coverage (1–5), and conciseness (1–5). Mean scores: coherence 4.1, coverage 3.8, conciseness 4.3. The main failure mode was incomplete coverage: for very long documents, the MapReduce chain occasionally dropped minority topics that appeared in only a few chunks. The overall quality was judged acceptable for orientation and navigation purposes, if not for exhaustive reference.

VI. DISCUSSION

6.1 What the Results Actually Tell Us

The 91% faithfulness score is the headline number, but the more practically important finding is the out-of-context refusal rate: 94% of queries about topics not in the document were correctly refused. This is harder to achieve than it sounds. A model given irrelevant chunks and told to answer only from them will sometimes hallucinate anyway, especially if the question is adjacent to something the model knows from training. The explicit refusal instruction, combined with the top-5 retrieval mechanism, prevents most of this.

The 7% residual hallucination rate deserves scrutiny. Manual review of the failing cases shows three patterns: (1) the model correctly retrieves relevant chunks but misreads a number or date; (2) the query is a comparative question ('which of these is larger?') where the answer requires reasoning across multiple

chunks rather than reading a single passage; (3) the relevant information is in a table or list and PyPDF extracts it with formatting artifacts that confuse the chunker.

Pattern (1) is a model-level limitation. Pattern (2) is a retrieval architecture limitation: the current setup treats chunks as independent, not as a graph. Pattern (3) is fixable in the preprocessing pipeline.

6.2 Comparison with Existing Systems

Against PrivateGPT [14] and LlamaIndex-based demos [15], this system offers a comparable faithfulness score with a more accessible user interface. Most existing local RAG demos expose the system through a command-line or Jupyter interface. The Next.js frontend makes the system practical for non-technical users.

The BM25 comparison is informative. BM25 retrieval is cheap, well-understood, and already widely deployed. The 21-point Precision@5 gap justifies the additional complexity of running an embedding model, but only if the user's queries are semantically diverse. For a narrow-domain corpus where users ask straightforward keyword-matchable questions, BM25 may be sufficient.

6.3 Limitations Worth Being Honest About

The system does not handle scanned PDFs. PyPDF Loader extracts text from digitally-typed documents only. A scanned PDF returns an empty string, and the system currently shows a warning rather than attempting OCR. This excludes a large class of real-world documents.

The single-GPU deployment is not production-ready. Five concurrent users is a ceiling, not a target. Anyone expecting to deploy this for a team needs to address the inference bottleneck — either with batched inference via vLLM or with a request queue and multiple model instances.

Evaluation was done on a hand-curated dataset of 50 queries. That is enough to show trends but not enough to characterize performance on rare document types (handbooks, contracts, code documentation). Larger-scale evaluation with standardized benchmarks like BEIR [16] would strengthen the claims.

6.4 Practical Takeaways

For a student or researcher who wants a private, local document assistant, this system works well today. Upload a paper, ask questions about it, follow up. The latency is noticeable (5 seconds) but not disruptive. The answers are generally accurate and cite their source.

For institutional deployment, the gap between this prototype and a production system is real but bridgeable. The architecture is sound. The main engineering work is containerization, request queuing, user authentication, and OCR integration.

VII. FUTURE SCOPE

Several directions are worth pursuing:

- OCR integration: Adding Tesseract or PaddleOCR would allow the system to process scanned and image-heavy PDFs, which currently fail silently. This is probably the highest-impact near-term improvement for real-world usability.
- Hybrid retrieval: Combining dense (ChromaDB cosine) and sparse (BM25) retrieval with Reciprocal Rank Fusion [17] would likely close the 18% residual hallucination gap on comparative and numerical queries where dense retrieval alone underperforms.
- Cross-encoder re-ranking: Adding a cross-encoder (e.g., ms-marco-MiniLM-L-12-v2) to rerank the top-5 retrieved chunks before passing them to Mistral has been shown to improve faithfulness by 5-10 percentage points in similar systems [2].
- Multi-document queries: The current system queries a single document per session. Extending ChromaDB to handle a user's entire document library and routing queries across it would substantially increase utility.

- Agentic extensions: LangChain Agents and the ReAct framework [18] would allow the system to decompose complex queries into sub-queries, call external tools (calculators, web search), and synthesize results — moving from document Q&A toward a more general research assistant.
- Production deployment: Containerizing with Docker, switching from Ollama to vLLM for batched inference, and adding Celery-backed task queues would make the system viable for team-scale deployment.
- Fine-tuning: For specialized domains (legal, medical), fine-tuning the embedding model on domain-specific text pairs would likely improve retrieval precision more cost-effectively than scaling the model size.

VIII. CONCLUSION

We built a local RAG system that lets users ask questions about their own PDFs, privately, without sending data to any external service. The system uses Mistral 7B and Nomic-embed-text via Ollama, stores vectors in ChromaDB, and wraps everything in a Next.js chat interface.

The evaluation results are straightforward: faithfulness of 0.91, hallucination rate of 7%, and Precision@5 of 0.82. Compared to a plain LLM, that is a 57-point improvement in faithfulness. Compared to BM25 retrieval, it is a 21-point improvement in retrieval precision. Multi-turn memory resolves 93% of pronoun-reference follow-up queries correctly.

The main limitations are known: no OCR support, single-GPU scalability ceiling, and evaluation on a modest 50-query dataset. None of these are fundamental architectural problems.

The more interesting finding, for us at least, is how much of the faithfulness improvement comes from the prompt rather than from the retrieval. Telling the model not to invent answers, and giving it retrieved text to work from, accounts for most of the gap. The retrieval system's job is to make sure the right text is available. The prompt's job is to make sure the model reads it honestly. Both matter, and neither alone is sufficient.

IX. REFERENCES (IEEE FORMAT)

- [1] P. Lewis et al., "Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks," in *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 33, pp. 9459–9474, 2020.
- [2] G. Izacard and E. Grave, "Leveraging Passage Retrieval with Generative Models for Open Domain Question Answering," in *Proc. 16th Conference of the European Chapter of the ACL (EACL)*, pp. 874–880, 2021.
- [3] Y. Gao et al., "Retrieval-Augmented Generation for Large Language Models: A Survey," *arXiv preprint arXiv:2312.10997*, 2023.
- [4] V. Karpukhin et al., "Dense Passage Retrieval for Open-Domain Question Answering," in *Proc. 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 6769–6781, 2020.
- [5] N. Reimers and I. Gurevych, "Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks," in *Proc. 2019 Conference on EMNLP*, pp. 3982–3992, 2019.
- [6] Z. Nussbaum et al., "Nomic Embed: Training a Reproducible Long Context Text Embedder," *arXiv preprint arXiv:2402.01613*, 2024.
- [7] J. Anton et al., "Chroma: Open-Source Embedding Database," *GitHub repository*, <https://github.com/chroma-core/chroma>, 2023.
- [8] A. Q. Jiang et al., "Mistral 7B," *arXiv preprint arXiv:2310.06825*, 2023.
- [9] H. Touvron et al., "Llama 2: Open Foundation and Fine-Tuned Chat Models," *arXiv preprint arXiv:2307.09288*, 2023.
- [10] Z. Ji et al., "Survey of Hallucination in Natural Language Generation," *ACM Computing Surveys*, vol. 55, no. 12, pp. 1–38, 2023.
- [11] S. Es et al., "RAGAS: Automated Evaluation of Retrieval Augmented Generation," *arXiv preprint arXiv:2309.15217*, 2023.

- [12] H. Chase, "LangChain," GitHub repository, <https://github.com/langchain-ai/langchain>, 2022.
- [13] J. Xu et al., "Retrieval Meets Long Context Large Language Models," arXiv preprint arXiv:2310.03025, 2023.
- [14] I. Zaccardi, "PrivateGPT: Interact Privately with Your Documents," GitHub repository, <https://github.com/imartinez/privateGPT>, 2023.
- [15] J. Liu, "LlamaIndex: A Data Framework for LLM Applications," GitHub repository, https://github.com/jerryliu/llama_index, 2022.
- [16] N. Thakur et al., "BEIR: A Heterogeneous Benchmark for Zero-Shot Evaluation of Information Retrieval Models," in Proc. 35th NeurIPS Datasets and Benchmarks Track, 2021.
- [17] G. V. Cormack, C. L. A. Clarke, and S. Buettcher, "Reciprocal Rank Fusion Outperforms Condorcet and Individual Rank Learning Methods," in Proc. 32nd ACM SIGIR Conference, pp. 758–759, 2009.
- [18] S. Yao et al., "ReAct: Synergizing Reasoning and Acting in Language Models," in Proc. 11th International Conference on Learning Representations (ICLR), 2023.

