

CMP6200

Individual Undergraduate Project

2024 - 2025

A2 - Literature Review and Methods

University Artificially Intelligent Assistant



Course: Computer & Data Science
Student Name: Lewis Higgins
Student Number: 22133848
Supervisor Name: Dr. Atif Azad

Contents

1	Report Introduction	2
1.1	Aims and Objectives	2
1.2	Literature Search Methodology	3
2	Literature Review	5
2.1	Themes	5
2.2	Review of Literature	7
2.2.1	Artificial Intelligence (AI)	7
2.2.2	Natural language processing (NLP)	8
2.2.3	Large language models	9
2.2.4	Retrieval-Augmented Generation	10
2.2.5	Chatbots / Conversational Agents	11
2.3	Summary	12
3	Appendix	13
3.1	Gantt Chart	13
	References	18
	Bibliography	19

Report Introduction

1.1 Aims and Objectives

This project aims to leverage to aid new and existing students alike while they are attending university with helpful information about university itself, such as university societies, locations/campuses, and policies through the medium of a digital chatbot companion to converse with. The project's objectives are to:

- Conduct a thorough literature review on the surrounding topics, namely AI, LLMs and NLP.
- Create effective documentation for all stages of development, highlighting challenges faced during the process.
- Leverage Retrieval-Augmented Generation alongside a cloud-based LLM to query a vector database of university-related data.
- Develop a chatbot capable of accurately answering user queries related to university buildings, policies, and societies with a minimum 80% accuracy rate.
- Evaluate the effectiveness of an AI assistant on university student acclimatization.

1.2 Literature Search Methodology

My literature search will be performed using multiple reputable databases for academic papers, including:

- IEEE Xplore
- Scopus / Elsevier
- Google Scholar
- arXiv
- BCU Online Library

By using multiple different databases to source my information from, I can ensure that any potentially relevant literature will be found. Figure 1.1 depicts how in a search for 1685 articles about employee retention strategies and turnover, only 582 (25.7%) appeared in multiple databases (Wanyama et al., 2022), meaning that the remaining 74.3% of articles were exclusive to the single database in which they were found, emphasising the importance of searching multiple databases.

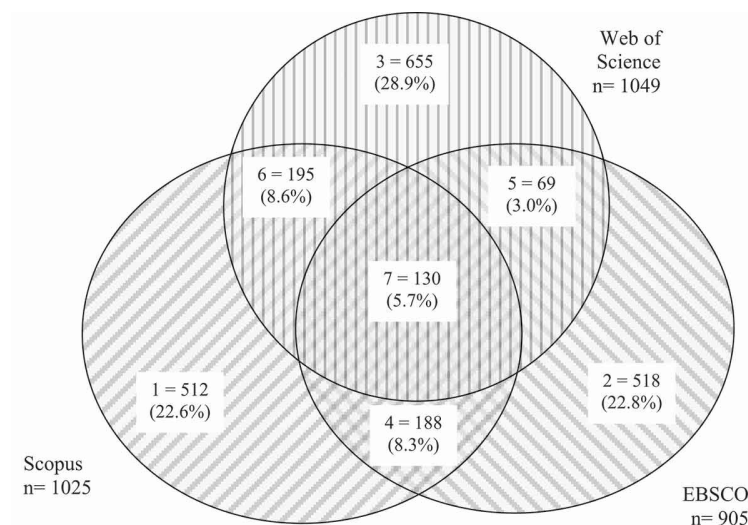


Figure 1.1: Distribution of searched articles across databases. (Wanyama et al., 2022)

All searches performed for recent literature will have a heavy preference to more recent literature, due to the constantly evolving fields my project is based on. The search terms I will use to retrieve the data I will be studying are:

- Artificial Intelligence / AI
- Natural Language Processing / NLP
- Large Language Models / LLMs

- Chatbots / Conversational Agents
- Retrieval-Augmented Generation / RAG

By using these specific terms that are directly relevant to the core themes of my project, I will be ensuring that I only retrieve literature that will be of crucial use in its development.

Literature Review

2.1 Themes

To develop the artefact and conduct thorough background research on relevant literature to further my knowledge of the subject areas, key general themes of the project were identified in Table 2.1. From these themes, further keywords to be used in the literature search were derived to ensure that retrieved literature is directly relevant to my research and development of the final artefact. Due to the constantly evolving fields the project focuses on, it will be necessary to limit the results to primarily those written in recent years as there are frequent new developments in the subject areas.

Theme	Description	Keywords
AI	A field of computing dedicated to allowing computers to simulate human learning by training them on large amounts of data so that they can recognise patterns to classify or predict unknown data. AI can only be as good as the data it is trained upon, and can develop biases if it is fed too much data of a certain type.	Generative AI, Human-Centred AI, Explainable AI, AI Ethics, AI Bias
Natural Language Processing	NLP refers to the use of machine learning to encode and process text to understand it in a similar way to humans, which can be used to allow direct two-way conversation between users and computers.	Embedding models, Vectorisation Semantic search, Entity linking
LLMs	Large Language Models are a type of machine learning model dedicated to the recognition and generation of text. As suggested by their name, they are trained on enormous amounts of text data, which allows them to have active conversations with users. There are many different LLMs, and as their size and complexity increases, so too does the necessary processing power.	Fine-tuning, Prompt engineering, Impact on industry, GPT4o, LLaMA, Gemini, Evaluation
Retrieval-Augmented Generation (RAG)	The optimisation of the generated text output of an LLM, incorporating an external data source to enhance its contextual knowledge and enhancing the subject relevancy of outputs.	Embedding, Vector databases, Document retrieval, Prompt engineering
Chatbot Conversational Agent	Software that simulates a natural conversation between the computer and end user. Many chatbots, including the one to be produced in this project, utilise recent developments such as Generative AI and natural language processing (NLP) to interpret and respond to user queries. (IBM, 2024d)	NLP, Digital assistant, ChatGPT, Risks, Impact on industry

Table 2.1: The themes and keywords used in the literature search.

2.2 Review of Literature

2.2.1 Artificial Intelligence (AI)

Researchers have always wanted to harness the processing power of computers to act in a similar manner indistinguishable from that of humans, most notably from as long ago as 1950, where the question was posed 'Can machines think?' (Turing, 1950). Ever since, constant innovations were made in computer intelligence and machine learning, from playing games of checkers at a better level than human players (Samuel, 1959) to classifying the contents of millions of images using convolutional neural networks (Krizhevsky et al., 2012).

Recently, AI is used across many disciplines and for different purposes to complete tasks faster than, and in some cases better than, human workers. Wirtz et al. (2018) write that 'service robots' ¹ can complete a variety of tangible or intangible actions, such as reading and sending text as a chatbot.

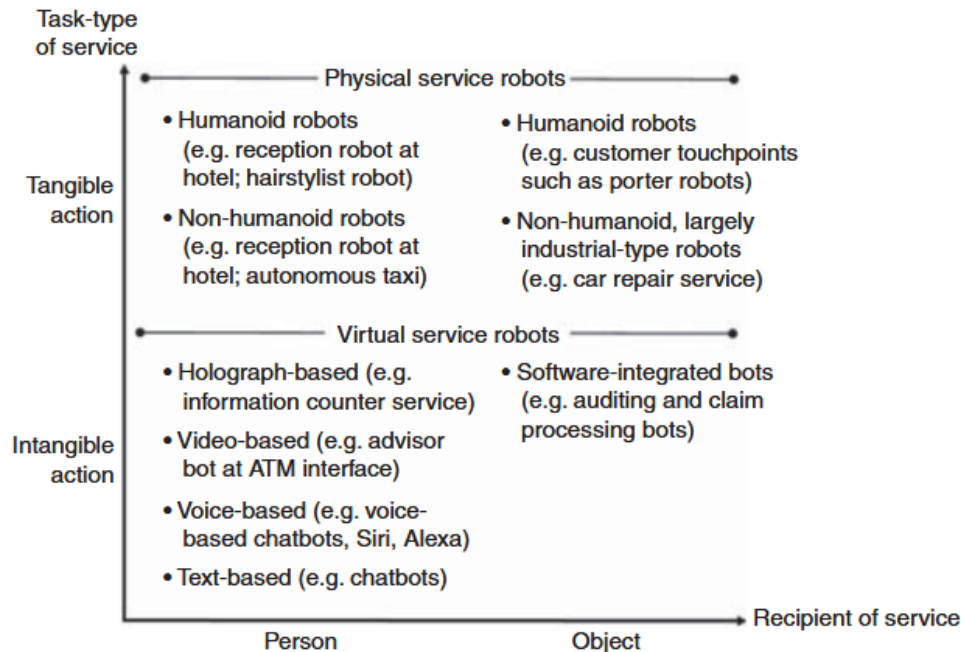


Figure 2.1: Service robots categorization by task-type and recipient of service (Wirtz et al., 2018).

Today, AI is still a constantly evolving field that is seeing bleeding-edge developments on a highly frequent basis, and more recently, is becoming instrumental in many people's work and private lives with the introduction of large language models (LLMs) (Maedche et al., 2019). However, when developing a project that utilises AI, it is important that they are ethical and human-centred in the development process, which is known as Human-Centred AI (HCAI). Another issue is the "black-box problem" - the inability to know an AI's reasoning, meaning that eXplainable AI (XAI) is a growing necessity (Miró-Nicolau et al., 2025).

¹Defined as "system-based autonomous and adaptable interfaces that interact, communicate and deliver service to an organization's customers" (Wirtz et al., 2018, p.909)

In focusing on HCAI and XAI, the focus shifts from the machine executing the algorithms, and instead to the user and their experience using the AI. Shneiderman (2020) strongly advocates for the promotion of HCAI for the benefit of both companies and their users, which is a commonly accepted idea due to the ethical risks of using AI.

Because AI calculates outcomes from its training data rather than understanding social norms and perspectives, the use of it in sociotechnical systems poses serious risks due to the 'traps' it can fall into, because it cannot account for every possibility such as the personal tendencies and biases of its users (Selbst et al., 2019), and therefore developers require a shift in focus - from the final product at the end of development to the development process itself and end users, which also echoes Shneiderman's views.

2.2.2 Natural language processing (NLP)

The ability for a computer to interpret and understand human language greatly enhances the scale of their capabilities. This was recognised during the 1950s, where machine translation from Russian to English was demonstrated for the first time, albeit in a basic form (Jones, 1994). Ever since, NLP has persistently been a key topic in computing, and even more so become in recent years, with its applications becoming very wide in scope with modern processing power.

One of the key advancements in NLP is vectorisation, a process where data is embedded into a numerical equivalent that a computer can interpret, enabling Natural Language Understanding (NLU) and the identification of semantic similarities between words through the use of an embedding model like Word2Vec (Mikolov et al., 2013) without the need to manually label data. Word2Vec was a key innovation in NLP, and Mikolov and Le went on to improve it further with Doc2Vec (Le and Mikolov, 2014), which could embed entire documents into semantically searchable vectorised forms.

Embedding models have further improved since, most notably with Vaswani et al. (2017)'s Transformer architecture enhancing models such as BERT (Devlin et al., 2019), which is able to analyse context through analysing multiple neighbours of a word rather than reading from left to right, gaining a higher understanding of the text it processes. Many embedding models have since been developed, though one of the most reputable is OpenAI's recent text-embeddings-3 model (OpenAI, 2024b), which can be used in the development of the chatbot at a low cost.

2.2.3 Large language models

LLMs are colossal machine learning models that leverage NLP to generate text, and have become widely used across industries in place of technical support and human resources (Vrontis et al., 2022). The training data required for an LLM is immense, reaching 45 terabytes of text data for ChatGPT in 2023 (Dwivedi et al., 2023).

This data is harvested from the web (Dubey et al., 2024) and social media due to it being one of the largest repositories of opinionated text data (Z. Wang et al., 2016), such as posts on platforms like Facebook and X. However, meticulous care is taken into the specific sources used to remove Personally Identifiable Information (PII) to minimise privacy and ethical concerns (Dubey et al., 2024).

The previously mentioned Transformer from Vaswani et al. (2017) became a staple in LLMs due to the major reduction in necessary processing power to produce higher-quality results, and it continues to underpin many LLMs today, including ChatGPT (Brown et al., 2020). Even with these enhancements, LLMs are still extremely performance intensive, requiring more than 8 top-range server-grade GPUs to run some of the most powerful high-parameter open-source models like LLaMA 3.1's 405 billion parameter model (Dubey et al., 2024), and many therefore use cloud API solutions to access LLMs.

The amount of parameters in a model does not entirely account for the quality of its responses, as studied by Ouyang et al. (2022) in Figure 2.2 wherein their surveys revealed their fine-tuned LLM "InstructGPT" with over 100x less parameters than a 175 billion parameter GPT3 model would often give answers preferred by its human assessors, which reveals that the fine-tuning and prompt engineering of an LLM is as vitally important to the quality of its responses as the amount of parameters.

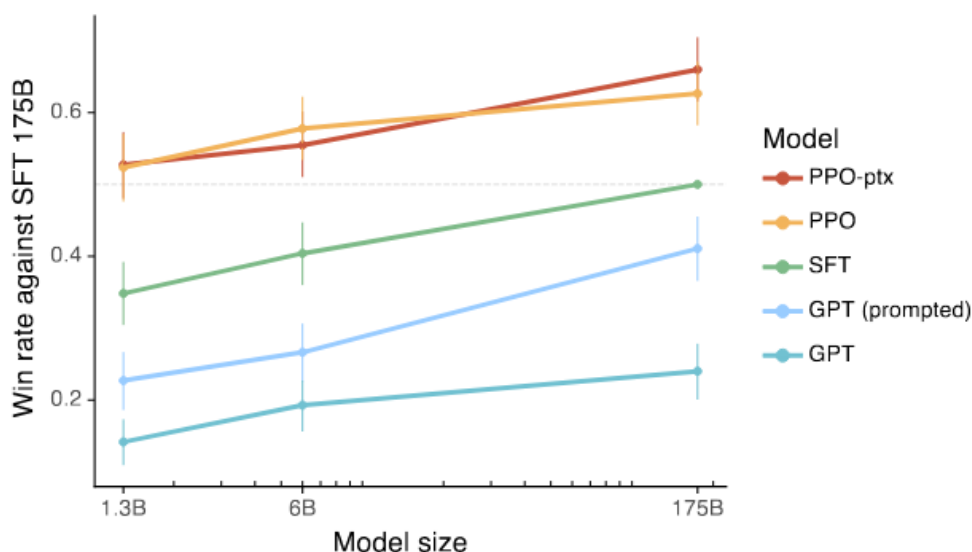


Figure 2.2: Human evaluations of the GPT models produced by Ouyang et al. (2022). PPO and PPO-ptx are their models.

The simplest way to measure the accuracy and quality of an LLM's responses is through human evaluation surveys such as that conducted by Ouyang et al. (2022), though software

approaches such as the open-source DeepEval can be used. DeepEval offers 14 metrics to test LLM outputs with (DeepEval, 2024), with a notable metric being "G-Eval", originally introduced by Liu et al. (2023), which uses an "LLM-as-a-judge" approach where an LLM will evaluate and grade the quality of the output.

2.2.4 Retrieval-Augmented Generation

While LLMs are highly useful tools across many industries, they are not without limitations. The most notable of these limitations are hallucinations (P. Lewis et al., 2021), where the LLM will fabricate information that conflicts with user input, earlier conversation context or true facts (Zhang et al., 2023). This occurs as a direct result of the LLM's parametric memory² being overfitted or biased, which can be counteracted through introducing an external knowledge source, known as non-parametric memory (Komeili et al. (2022), Siriwardhana et al. (2023)).

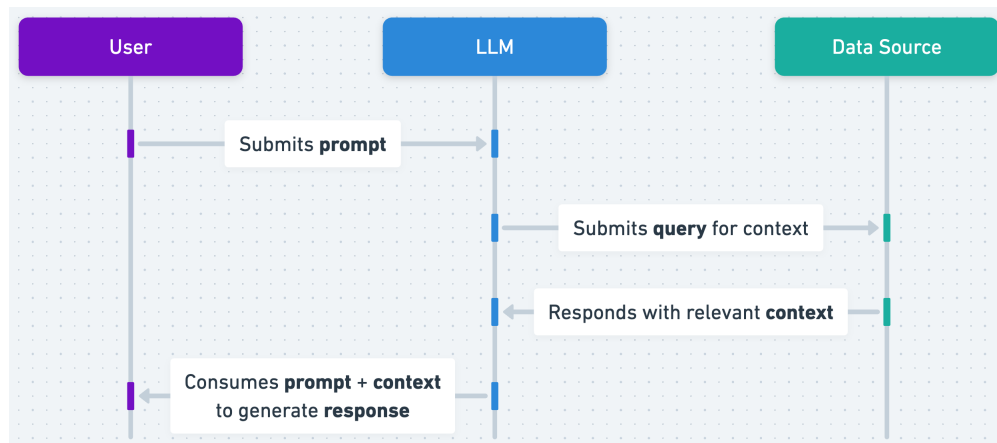


Figure 2.3: A basic overview of a RAG workflow (OpenAI, 2024a).

Siriwardhana et al. (2023) expanded upon the earlier works of Karpukhin et al. (2020) and M. Lewis et al. (2020) by creating "RAG-end2end", which explored the capabilities of RAG on a dynamically updating knowledge store, meaning the LLM itself would not have to be retrained every time the data updates, saving enormous amounts of processing power.

RAG is dependent upon external knowledge stores such as vector databases, which store and process vectorised data for non-parametric memory (Li, 2023), which makes them an essential part of the backend of a RAG-enabled chatbot as studied by Odede and Frommholz (2024).

Many software options exist for vector databases, such as Milvus (J. Wang et al., 2021), Pinecone (Pinecone, 2024) and Chroma³ (Chroma, 2024). A study by Xie et al. (2023) compared these three, citing Pinecone's 'robust distributed computing capabilities and scalability', and its common usage in real-time searching scenarios.

Pinecone was also used in chatbots by Odede and Frommholz (2024) and Singer et al. (2024), showcasing its potential as a vector database solution for chatbots.

²Knowledge that the LLM has from its training data (Siriwardhana et al., 2023).

³Also known as ChromaDB.

LangChain (LangChain, 2024) is a popular open-source framework for RAG pipelines that can be used to connect backend elements together, as described by Singer et al. (2024) when they used it to chunk their text data and connect to their Pinecone vector database to store the embedded data.

2.2.5 Chatbots / Conversational Agents

Conversational agents, better known as chatbots, leverage MLP in order to simulate a conversational flow between a user and machine, and have become mainstream products in recent years (Liao et al., 2018), though have existed as far back as 1966 with the creation of "ELIZA" for the IBM-7094 (Weizenbaum, 1966). As time has passed, advancements in chatbots have occurred in "waves", where each new wave has brought a major innovation (Schöbel et al., 2024).

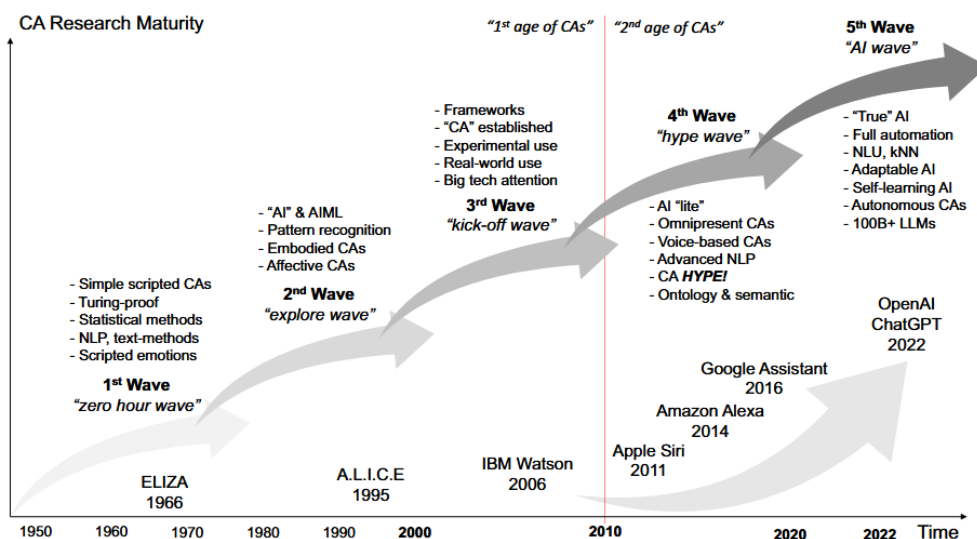


Figure 2.4: The five waves of conversational agent research (Schöbel et al., 2024).

As a product of the considerable developments in the field, chatbots are now widely used across industries such as education (Kuhail et al., 2023). However, the use of the latest wave of chatbots based on LLMs, especially in educational settings, poses significant risks as studied by Neumann et al. (2024) due to the risk of hallucinations being interpreted as absolute fact, although Shuster et al. (2021) argued that this risk can be greatly reduced through introducing RAG to the backend LLM, which is further backed by the RAG-based chatbot created by Ge et al. (2023), which they found to also give superior answers in their medical field of study to those of a general-purpose chatbot without RAG.

Many platforms exist to aid chatbot development, though they are typically aimed at users from non-IT backgrounds (Srivastava and Prabhakar, 2020). Popular platforms include IBM's watsonx Assistant (IBM, 2024a), Google's Dialogflow (Google, 2024) and Microsoft's Bot Framework (Microsoft, 2024). However, these are primarily targeted at enterprise clients which is reflected in their pricing. Instead of using these, the chatbot can be manually developed using LangChain as its framework.

2.3 Summary

In conclusion, this literature review has revealed multiple key areas of focus for the development of the chatbot. The overall design of the chatbot must be iterative and human-centred, and user feedback should be obtained at every possible opportunity to ensure the resultant product is high quality. A deep exploration into AI, specifically in its applications in NLP and LLMs, has revealed that the best option for the chatbot will be to leverage a pre-existing cloud-based LLM's RAG capabilities, as to do so on a local machine would require an infeasible amount of processing power. The non-parametric memory accessed through RAG would be a vector database created with Pinecone storing embeddings generated by OpenAI's text-embeddings-3-small model, and the overall framework will be LangChain. This will keep the cost of the project low while maintaining a tolerable level of quality in the bot's responses.

Appendix

3.1 Gantt Chart

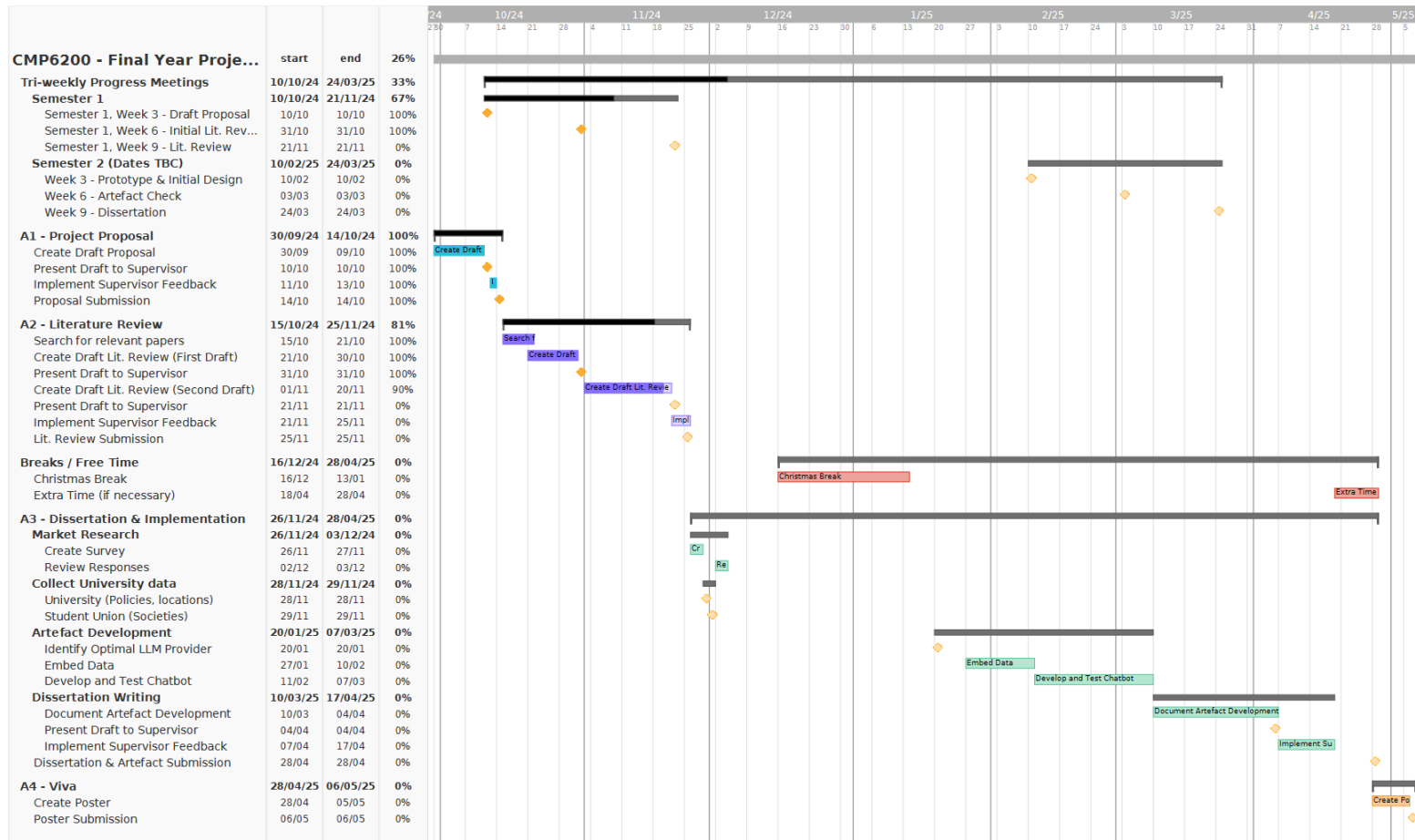


Figure 3.1: The updated Gantt Chart for the development timeline.

References

- Chroma (2024). *Chroma*. URL: <https://www.trychroma.com/> (visited on 11/24/2024).
- DeepEval (Nov. 22, 2024). *Introduction | DeepEval - The Open-Source LLM Evaluation Framework*. URL: <https://docs.confident-ai.com/docs/metrics-introduction> (visited on 11/24/2024).
- Devlin, Jacob, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova (June 2019). “BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding”. In: *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. NAACL-HLT 2019. Ed. by Jill Burstein, Christy Doran, and Thamar Solorio. Minneapolis, Minnesota: Association for Computational Linguistics, pp. 4171–4186. DOI: [10.18653/v1/N19-1423](https://doi.org/10.18653/v1/N19-1423).
- Dubey, Abhimanyu et al. (Aug. 15, 2024). *The Llama 3 Herd of Models*. DOI: [10.48550/arXiv.2407.21783](https://doi.org/10.48550/arXiv.2407.21783). arXiv: [2407.21783](https://arxiv.org/abs/2407.21783).
- Dwivedi, Y.K., N. Kshetri, L. Hughes, E.L. Slade, A. Jeyaraj, A.K. Kar, A.M. Baabdullah, A. Koohang, V. Raghavan, M. Ahuja, H. Albanna, M.A. Albashrawi, A.S. Al-Busaidi, J. Balakrishnan, Y. Barlette, S. Basu, I. Bose, L. Brooks, D. Buhalis, L. Carter, S. Chowdhury, T. Crick, S.W. Cunningham, G.H. Davies, R.M. Davison, R. Dé, D. Dennehy, Y. Duan, R. Dubey, R. Dwivedi, J.S. Edwards, C. Flavián, R. Gault, V. Grover, M.-C. Hu, M. Janssen, P. Jones, I. Junglas, S. Khorana, S. Kraus, K.R. Larsen, P. Latreille, S. Laumer, F.T. Malik, A. Mardani, M. Mariani, S. Mithas, E. Mogaji, J.H. Nord, S. O’Connor, F. Okumus, M. Pagani, N. Pandey, S. Papagiannidis, I.O. Pappas, N. Pathak, J. Pries-Heje, R. Raman, N.P. Rana, S.-V. Rehm, S. Ribeiro-Navarrete, A. Richter, F. Rowe, S. Sarker, B.C. Stahl, M.K. Tiwari, W. van der Aalst, V. Venkatesh, G. Viglia, M. Wade, P. Walton, J. Wirtz, and R. Wright (2023). ““So what if ChatGPT wrote it?” Multidisciplinary perspectives on opportunities, challenges and implications of generative conversational AI for research, practice and policy”. In: *International Journal of Information Management* 71. DOI: [10.1016/j.ijinfomgt.2023.102642](https://doi.org/10.1016/j.ijinfomgt.2023.102642).
- Ge, Jin, Steve Sun, Joseph Owens, Victor Galvez, Oksana Gologorskaya, Jennifer C Lai, Mark J Pletcher, and Ki Lai (Nov. 1, 2023). “Development of a Liver Disease-Specific Large Language Model Chat Interface using Retrieval Augmented Generation”. In: *medRxiv*, p. 2023.11.10.23298364. DOI: [10.1101/2023.11.10.23298364](https://doi.org/10.1101/2023.11.10.23298364).
- Google (2024). *Conversational Agents and Dialogflow*. Google Cloud. URL: <https://cloud.google.com/products/conversational-agents> (visited on 11/24/2024).
- IBM (Apr. 3, 2024a). *IBM watsonx Assistant Virtual Agent*. URL: <https://www.ibm.com/products/watsonx-assistant> (visited on 11/24/2024).
- Karpukhin, Vladimir, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih (Nov. 2020). “Dense Passage Retrieval for Open-Domain Question Answering”. In: *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. EMNLP 2020. Ed. by Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu. Online: Association for Computational Linguistics, pp. 6769–6781. DOI: [10.18653/v1/2020.emnlp-main.550](https://doi.org/10.18653/v1/2020.emnlp-main.550).
- Komeili, Mojtaba, Kurt Shuster, and Jason Weston (May 2022). “Internet-Augmented Dialogue Generation”. In: *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. ACL 2022. Ed. by Smaranda Mure-

- san, Preslav Nakov, and Aline Villavicencio. Dublin, Ireland: Association for Computational Linguistics, pp. 8460–8478. DOI: [10.18653/v1/2022.acl-long.579](https://doi.org/10.18653/v1/2022.acl-long.579).
- Krizhevsky, Alex, Ilya Sutskever, and Geoffrey E Hinton (2012). “ImageNet Classification with Deep Convolutional Neural Networks”. In: *Advances in Neural Information Processing Systems*. Vol. 25. Curran Associates, Inc. DOI: [10.1145/3065386](https://doi.org/10.1145/3065386).
- Kuhail, M.A., N. Alturki, S. Alramlawi, and K. Alhejori (2023). “Interacting with educational chatbots: A systematic review”. In: *Education and Information Technologies* 28 (1), pp. 973–1018. DOI: [10.1007/s10639-022-11177-3](https://doi.org/10.1007/s10639-022-11177-3).
- LangChain (2024). *Introduction / LangChain*. URL: <https://python.langchain.com/docs/introduction/> (visited on 11/24/2024).
- Le, Quoc V. and Tomas Mikolov (May 22, 2014). *Distributed Representations of Sentences and Documents*. DOI: [10.48550/arXiv.1405.4053](https://doi.org/10.48550/arXiv.1405.4053). arXiv: [1405.4053](https://arxiv.org/abs/1405.4053).
- Lewis, Mike, Marjan Ghazvininejad, Gargi Ghosh, Armen Aghajanyan, Sida Wang, and Luke Zettlemoyer (June 26, 2020). *Pre-training via Paraphrasing*. DOI: [10.48550/arXiv.2006.15020](https://doi.org/10.48550/arXiv.2006.15020). arXiv: [2006.15020](https://arxiv.org/abs/2006.15020).
- Lewis, Patrick, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela (Apr. 12, 2021). *Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks*. DOI: [10.48550/arXiv.2005.11401](https://doi.org/10.48550/arXiv.2005.11401). arXiv: [2005.11401](https://arxiv.org/abs/2005.11401).
- Li, Feifei (Aug. 1, 2023). “Modernization of Databases in the Cloud Era: Building Databases that Run Like Legos”. In: *Proc. VLDB Endow.* 16 (12), pp. 4140–4151. ISSN: 2150-8097. DOI: [10.14778/3611540.3611639](https://doi.org/10.14778/3611540.3611639).
- Liao, Q. Vera, Muhammed Mas-ud Hussain, Praveen Chandar, Matthew Davis, Yasaman Khazaeni, Marco Patricio Crasso, Dakuo Wang, Michael Muller, N. Sadat Shami, and Werner Geyer (Apr. 19, 2018). “All Work and No Play?” In: *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*. CHI '18. New York, NY, USA: Association for Computing Machinery, pp. 1–13. ISBN: 978-1-4503-5620-6. DOI: [10.1145/3173574.3173577](https://doi.org/10.1145/3173574.3173577).
- Liu, Yang, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu (May 23, 2023). *G-Eval: NLG Evaluation using GPT-4 with Better Human Alignment*. DOI: [10.48550/arXiv.2303.16634](https://doi.org/10.48550/arXiv.2303.16634). arXiv: [2303.16634](https://arxiv.org/abs/2303.16634).
- Maedche, Alexander, Christine Legner, Alexander Benlian, Benedikt Berger, Henner Gimpel, Thomas Hess, Oliver Hinz, Stefan Morana, and Matthias Söllner (Aug. 1, 2019). “AI-Based Digital Assistants”. In: *Business & Information Systems Engineering* 61 (4), pp. 535–544. ISSN: 1867-0202. DOI: [10.1007/s12599-019-00600-8](https://doi.org/10.1007/s12599-019-00600-8).
- Microsoft (2024). *Microsoft Bot Framework*. URL: <https://dev.botframework.com/> (visited on 11/24/2024).
- Mikolov, Tomas, Kai Chen, Greg Corrado, and Jeffrey Dean (Sept. 7, 2013). *Efficient Estimation of Word Representations in Vector Space*. DOI: [10.48550/arXiv.1301.3781](https://doi.org/10.48550/arXiv.1301.3781). arXiv: [1301.3781](https://arxiv.org/abs/1301.3781).
- Miró-Nicolau, M., A. Jaume-i-Capó, and G. Moyà-Alcover (2025). “A comprehensive study on fidelity metrics for XAI”. In: *Information Processing and Management* 62 (1). DOI: [10.1016/j.ipm.2024.103900](https://doi.org/10.1016/j.ipm.2024.103900).
- Neumann, Alexander Tobias, Yue Yin, Sulayman Sowe, Stefan Decker, and Matthias Jarke (2024). “An LLM-Driven Chatbot in Higher Education for Databases and Information

- Systems”. In: *IEEE Transactions on Education*. Conference Name: IEEE Transactions on Education, pp. 1–14. ISSN: 1557-9638. DOI: [10.1109/TE.2024.3467912](https://doi.org/10.1109/TE.2024.3467912).
- Odede, Julius and Ingo Frommholz (Mar. 10, 2024). “JayBot – Aiding University Students and Admission with an LLM-based Chatbot”. In: *Proceedings of the 2024 Conference on Human Information Interaction and Retrieval*. CHIIR ’24. New York, NY, USA: Association for Computing Machinery, pp. 391–395. ISBN: 9798400704345. DOI: [10.1145/3627508.3638293](https://doi.org/10.1145/3627508.3638293).
- OpenAI (2024a). *Retrieval Augmented Generation (RAG) and Semantic Search for GPTs / OpenAI Help Center*. URL: <https://help.openai.com/en/articles/8868588-retrieval-augmented-generation-rag-and-semantic-search-for-gpts> (visited on 11/23/2024).
- Ouyang, Long, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe (Mar. 4, 2022). *Training language models to follow instructions with human feedback*. DOI: [10.48550/arXiv.2203.02155](https://doi.org/10.48550/arXiv.2203.02155). arXiv: [2203.02155](https://arxiv.org/abs/2203.02155).
- Pinecone (2024). *Pinecone Documentation*. Pinecone Docs. URL: <https://docs.pinecone.io/guides/get-started/overview> (visited on 11/24/2024).
- Samuel, A. L. (July 1959). “Some Studies in Machine Learning Using the Game of Checkers”. In: *IBM Journal of Research and Development* 3 (3). Conference Name: IBM Journal of Research and Development, pp. 210–229. ISSN: 0018-8646. DOI: [10.1147/rd.33.0210](https://doi.org/10.1147/rd.33.0210).
- Schöbel, Sofia, Anuschka Schmitt, Dennis Benner, Mohammed Saqr, Andreas Janson, and Jan Marco Leimeister (Apr. 1, 2024). “Charting the Evolution and Future of Conversational Agents: A Research Agenda Along Five Waves and New Frontiers”. In: *Information Systems Frontiers* 26 (2), pp. 729–754. ISSN: 1572-9419. DOI: [10.1007/s10796-023-10375-9](https://doi.org/10.1007/s10796-023-10375-9).
- Selbst, Andrew D., Danah Boyd, Sorelle A. Friedler, Suresh Venkatasubramanian, and Janet Vertesi (Jan. 29, 2019). “Fairness and Abstraction in Sociotechnical Systems”. In: *Proceedings of the Conference on Fairness, Accountability, and Transparency*. FAT* ’19. New York, NY, USA: Association for Computing Machinery, pp. 59–68. ISBN: 978-1-4503-6125-5. DOI: [10.1145/3287560.3287598](https://doi.org/10.1145/3287560.3287598).
- Shneiderman, Ben (Oct. 16, 2020). “Bridging the Gap Between Ethics and Practice: Guidelines for Reliable, Safe, and Trustworthy Human-centered AI Systems”. In: *ACM Trans. Interact. Intell. Syst.* 10 (4), 26:1–26:31. ISSN: 2160-6455. DOI: [10.1145/3419764](https://doi.org/10.1145/3419764).
- Shuster, Kurt, Spencer Poff, Moya Chen, Douwe Kiela, and Jason Weston (Nov. 2021). “Retrieval Augmentation Reduces Hallucination in Conversation”. In: *Findings of the Association for Computational Linguistics: EMNLP 2021*. Findings 2021. Ed. by Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih. Punta Cana, Dominican Republic: Association for Computational Linguistics, pp. 3784–3803. DOI: [10.18653/v1/2021.findings-emnlp.320](https://doi.org/10.18653/v1/2021.findings-emnlp.320).
- Singer, Maxwell B., Julia J. Fu, Jessica Chow, and Christopher C. Teng (Mar. 1, 2024). “Development and Evaluation of Aeyeconsult: A Novel Ophthalmology Chatbot Leveraging Verified Textbook Knowledge and GPT-4”. In: *Journal of Surgical Education* 81 (3), pp. 438–443. ISSN: 1931-7204. DOI: [10.1016/j.jsurg.2023.11.019](https://doi.org/10.1016/j.jsurg.2023.11.019).

- Siriwardhana, Shamane, Rivindu Weerasekera, Elliott Wen, Tharindu Kaluarachchi, Rajib Rana, and Suranga Nanayakkara (Jan. 12, 2023). “Improving the Domain Adaptation of Retrieval Augmented Generation (RAG) Models for Open Domain Question Answering”. In: *Transactions of the Association for Computational Linguistics* 11, pp. 1–17. ISSN: 2307-387X. DOI: [10.1162/tac1_a_00530](https://doi.org/10.1162/tac1_a_00530).
- Srivastava, Saurabh and T.V. Prabhakar (Sept. 2020). “Desirable Features of a Chatbot-building Platform”. In: *2020 IEEE International Conference on Humanized Computing and Communication with Artificial Intelligence (HCCAI)*. 2020 IEEE International Conference on Humanized Computing and Communication with Artificial Intelligence (HCCAI), pp. 61–64. DOI: [10.1109/HCCAI49649.2020.00016](https://doi.org/10.1109/HCCAI49649.2020.00016).
- Turing, A. M. (Oct. 1, 1950). “I.—COMPUTING MACHINERY AND INTELLIGENCE”. In: *Mind* LIX (236), pp. 433–460. ISSN: 1460-2113, 0026-4423. DOI: [10.1093/mind/LIX.236.433](https://doi.org/10.1093/mind/LIX.236.433).
- Vaswani, Ashish, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin (Dec. 4, 2017). “Attention is all you need”. In: *Proceedings of the 31st International Conference on Neural Information Processing Systems*. NIPS’17. Red Hook, NY, USA: Curran Associates Inc., pp. 6000–6010. ISBN: 978-1-5108-6096-4. DOI: [10.48550/arXiv.1706.03762](https://doi.org/10.48550/arXiv.1706.03762).
- Vrontis, Demetris, Michael Christofi, Vijay Pereira, Shlomo Tarba, Anna Makrides, and Eleni Trichina (Mar. 26, 2022). “Artificial intelligence, robotics, advanced technologies and human resource management: a systematic review”. In: *The International Journal of Human Resource Management* 33 (6). Publisher: Routledge, pp. 1237–1266. ISSN: 0958-5192. DOI: [10.1080/09585192.2020.1871398](https://doi.org/10.1080/09585192.2020.1871398).
- Wang, Jianguo, Xiaomeng Yi, Rentong Guo, Hai Jin, Peng Xu, Shengjun Li, Xiangyu Wang, Xiangzhou Guo, Chengming Li, Xiaohai Xu, Kun Yu, Yuxing Yuan, Yinghao Zou, Jiquan Long, Yudong Cai, Zhenxiang Li, Zhifeng Zhang, Yihua Mo, Jun Gu, Ruiyi Jiang, Yi Wei, and Charles Xie (June 18, 2021). “Milvus: A Purpose-Built Vector Data Management System”. In: *Proceedings of the 2021 International Conference on Management of Data*. SIGMOD ’21. New York, NY, USA: Association for Computing Machinery, pp. 2614–2627. ISBN: 978-1-4503-8343-1. DOI: [10.1145/3448016.3457550](https://doi.org/10.1145/3448016.3457550).
- Wang, Zhaoxia, Chee Seng Chong, Landy Lan, Yinping Yang, Seng Beng Ho, and Joo Chuan Tong (Dec. 2016). “Fine-grained sentiment analysis of social media with emotion sensing”. In: *2016 Future Technologies Conference (FTC)*. 2016 Future Technologies Conference (FTC), pp. 1361–1364. DOI: [10.1109/FTC.2016.7821783](https://doi.org/10.1109/FTC.2016.7821783).
- Wanyama, Seperia B., Ronald W. McQuaid, and Markus Kittler (May 4, 2022). “Where you search determines what you find: the effects of bibliographic databases on systematic reviews”. In: *International Journal of Social Research Methodology* 25 (3), pp. 409–422. ISSN: 1364-5579. DOI: [10.1080/13645579.2021.1892378](https://doi.org/10.1080/13645579.2021.1892378).
- Weizenbaum, Joseph (Jan. 1, 1966). “ELIZA—a computer program for the study of natural language communication between man and machine”. In: *Commun. ACM* 9 (1), pp. 36–45. ISSN: 0001-0782. DOI: [10.1145/365153.365168](https://doi.org/10.1145/365153.365168).
- Wirtz, J., P.G. Patterson, W.H. Kunz, T. Gruber, V.N. Lu, S. Paluch, and A. Martins (2018). “Brave new world: service robots in the frontline”. In: *Journal of Service Management* 29 (5), pp. 907–931. DOI: [10.1108/JOSM-04-2018-0119](https://doi.org/10.1108/JOSM-04-2018-0119).

- Xie, Xingrui, Han Liu, Wenzhe Hou, and Hongbin Huang (Dec. 2023). “A Brief Survey of Vector Databases”. In: *2023 9th International Conference on Big Data and Information Analytics (BigDIA)*. 2023 9th International Conference on Big Data and Information Analytics (BigDIA). ISSN: 2771-6902, pp. 364–371. DOI: [10.1109/BigDIA60676.2023.10429609](https://doi.org/10.1109/BigDIA60676.2023.10429609).
- Zhang, Yue, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu, Tingchen Fu, Xinting Huang, Enbo Zhao, Yu Zhang, Yulong Chen, Longyue Wang, Anh Tuan Luu, Wei Bi, Freda Shi, and Shuming Shi (Sept. 24, 2023). *Siren’s Song in the AI Ocean: A Survey on Hallucination in Large Language Models*. DOI: [10.48550/arXiv.2309.01219](https://doi.org/10.48550/arXiv.2309.01219). arXiv: [2309.01219](https://arxiv.org/abs/2309.01219).

Bibliography

Sources consulted but not directly cited

- AWS (2024). *What is RAG? - Retrieval-Augmented Generation AI Explained* - AWS. Amazon Web Services, Inc. URL: <https://aws.amazon.com/what-is/retrieval-augmented-generation/> (visited on 11/23/2024).
- Cloudflare (2024). *What is a large language model (LLM)?* URL: <https://www.cloudflare.com/en-gb/learning/ai/what-is-large-language-model/> (visited on 10/28/2024).
- Confident AI (2024). *LLM Evaluation Metrics: The Ultimate LLM Evaluation Guide - Confident AI*. URL: <https://www.confident-ai.com/blog/llm-evaluation-metrics-everything-you-need-for-llm-evaluation> (visited on 11/24/2024).
- Databricks (Oct. 18, 2023). *Retrieval Augmented Generation*. Databricks. URL: <https://www.databricks.com/glossary/retrieval-augmented-generation-rag> (visited on 11/23/2024).
- Elastic (2024). *What is Semantic Search? | A Comprehensive Semantic Search Guide*. URL: <https://www.elastic.co/what-is/semantic-search> (visited on 11/24/2024).
- IBM (Aug. 16, 2024b). *What is AI?* URL: <https://www.ibm.com/topics/artificial-intelligence> (visited on 10/28/2024).
- IBM (Aug. 11, 2024c). *What Is NLP (Natural Language Processing)? | IBM*. URL: <https://www.ibm.com/topics/natural-language-processing> (visited on 11/04/2024).
- IBM (2024d). *What is a chatbot?* URL: <https://www.ibm.com/topics/chatbots> (visited on 10/28/2024).
- IBM (2024e). *What is generative AI?* URL: <https://research.ibm.com/blog/what-is-generative-AI> (visited on 10/28/2024).
- ICO (2024). *Definitions*. URL: <https://ico.org.uk/for-organisations/uk-gdpr-guidance-and-resources/artificial-intelligence/explaining-decisions-made-with-artificial-intelligence/part-1-the-basics-of-explaining-ai/definitions/> (visited on 10/28/2024).
- MIT (Nov. 9, 2023). *Explained: Generative AI*. URL: <https://news.mit.edu/2023/explained-generative-ai-1109> (visited on 10/28/2024).