

Heidelberg University
Institut of Computer Science
Data Science Group

Master's Thesis

**Building Conversational Question
Answering Systems**

Name: Stephan Lenert
Matriculation Number: 3734583
Supervisor: Prof. Dr. Michael Gertz
Submission Date: March 15, 2024

Hiermit versichere ich, dass ich die Arbeit selbst verfasst und keine anderen als die angegebenen Quellen und Hilfsmittel benutzt und wörtlich oder inhaltlich aus fremden Werken Übernommenes als fremd kenntlich gemacht habe. Ferner versichere ich, dass die übermittelte elektronische Version in Inhalt und Wortlaut mit der gedruckten Version meiner Arbeit vollständig übereinstimmt. Ich bin einverstanden, dass diese elektronische Fassung universitätsintern anhand einer Plagiatssoftware auf Plagiate überprüft wird.

Abgabedatum: March 15, 2024

Zusammenfassung

Das Interesse an der Entwicklung von Systemen für die konversationelle Informationssuche hat durch den Boom von ChatGPT und Large Language Modellen im generellen stark zugenommen. Jedoch ist eine alleinige Verwendung von ChatGPT aufgrund von Herausforderungen wie der Aktualisierung und Veränderung von Wissen, Halluzinationen oder der Interpretierbarkeit von KI nicht ausreichend. Die Bewältigung dieser Hürden erfordert die Entwicklung von Retriever-Reader-Systemen. Allerdings ist die Konstruktion solcher Systeme für reale Anwendungsfälle in der aktuellen Forschung noch nicht ausreichend abgedeckt. Diese Arbeit schlägt ein umfassendes Rahmenwerk für Retrieval-Augmented Generation (RAG)-basierte konversationelle Frage-Antwort-Systeme vor, das vier Komponenten umfasst: Extractor, Contextual Query Understanding (CQU), Retriever und Reader. Es werden das Problemfeld, die Herausforderungen und Lösungsansätze für jede Komponente beschrieben, wobei der Fokus darauf liegt, bestehende Retriever- oder Reader-Modelle zu nutzen, anstatt neue Architekturen zu entwickeln. Für die Herausforderungen im Bereich des Retrievers werden Lösungen vorgeschlagen, die auf synthetischen Daten für Probleme im Zusammenhang mit dem Datensatz im generellen und verschiedenen Methoden wie Mixture-of-Experts für Probleme im Zusammenhang mit dem Evidence Set basieren. Für die Extract-Komponente wird eine pipelinefähige Kombination von Operationen vorgeschlagen. Die Aufgabe des Readers wird in mehrere kleinere Herausforderungen unterteilt, für die Lösungsansätze basierend auf Fine-Tuning oder Nachbearbeitung des Evidence Sets entwickelt werden. Unter Verwendung dieses Frameworks wird ein beispielhafter Frage-Antwort-Chatbot für die Prüfungsordnung der Universität Heidelberg entwickelt. Die Evaluation dieses Chatbots zeigt Engpässe in den Komponenten Retriever und Extractor auf, hebt die Nachteile kleinerer Sprachmodelle wie Llama2-7B-chat im Vergleich zu größeren Modellen wie gpt-3.5-turbo hervor und unterstreicht die Grenzen synthetischer Daten für die automatische Evaluation, die vor allem nur für faktoide Fragen anwendbar sind.

Abstract

The resurgence of interest in conversational information seeking, reignited by the boom of systems like ChatGPT and large language models in general, faces challenges such as knowledge extension, hallucination, and AI interpretability. Overcoming these hurdles necessitates the development of retriever-reader systems, yet constructing such systems from scratch for real-world use cases remains underexplored in current research. This thesis proposes a comprehensive framework for retrieval-augmented generation (RAG)-based conversational question-answering, comprising four components: Extractor, Contextual Query Understanding (CQU), Retriever, and Reader. We delineate the problem field, challenges, and solutions for each component, focusing on leveraging existing Retriever or Reader models rather than reinventing them. We address challenges in the Retriever, proposing solutions based on synthetic data for knowledge source challenges and various methods such as Mixture-of-Experts for evidence set challenges. For the Extract component, we propose a pipeline combination of operations. The Reader's task is broken down into multiple micro challenges and proposes solutions for them based on fine-tuning or post-processing of the evidence set. Applying this framework, we develop an exemplary question-answering chatbot for Heidelberg University's examination regulations. Evaluation uncovers bottlenecks in the Retriever and Extractor components, highlights the drawbacks of smaller language models like Llama2-7B-chat compared to larger models like gpt-3.5-turbo, and emphasizes the limitations of synthetic data for automatic evaluation, primarily applicable to factoid questions.

Contents

1. Introduction	1
1.1. Motivation	1
1.2. Objectives and Contributions	2
1.3. Thesis Structure	3
2. Background and Related Work	4
2.1. Question Answering	4
2.1.1. Basics	5
2.1.2. Information Retrieval Architectures	9
2.1.3. Extraction Approaches	11
2.1.4. Retrieval Approaches	13
2.1.5. Reader Approaches	17
2.1.6. Limitations	20
2.2. Conversational Question Answering	22
2.2.1. Basics	22
2.2.2. Contextual Query Understanding	24
2.2.3. Initiative	26
2.3. Efficient Large Language Models	27
2.3.1. Fine-Tuning	28
2.3.2. Compression	31
2.4. Related Work	35
2.4.1. Question Answering based on PDFs	35
2.4.2. Open-domain Conversational Question Answering	37
3. Open-domain QA Chatbot	39
3.1. Overview and Objective	39
3.2. Document Enquiry Model	41
3.3. Conversational Retrieval-Augmented Generation	43
3.3.1. Extract	44
3.3.2. Contextual Query Understanding	48
3.3.3. Retriever	49

Contents

3.3.4. Reader	55
3.4. Summary	59
4. Experimental Evaluation	61
4.1. Data	61
4.2. Evaluation Approach	65
4.2.1. Retrieval Evaluation Metrics	67
4.2.2. Reader Evaluation Metrics	68
4.2.3. End-to-End RAG Evaluation	70
4.3. Experimental Setup and Implementation	72
4.3.1. Extraction	74
4.3.2. Retriever	78
4.3.3. Reader	79
4.3.4. CQU	81
4.4. Experimental Results	82
4.4.1. Evaluation Limitations	83
4.4.2. Retrieval Results	84
4.4.3. Reader Results	85
4.4.4. Conversational Question Answering Results	88
5. Conclusions and Future Work	93
A. Appendix A	96
A.1. Implementation Trees	96
A.1.1. Example	96
A.2. Data Augmentation Implementation	97
A.2.1. Examples for Few-Shot Prompt	97
A.3. Evaluation	98
A.3.1. Model Configurations	98
A.3.2. Retriever	98

List of Acronyms

AI Artificial Intelligence.

AutoEval automatic evaluation.

BERT Bidirectional Encoder Representations from Transformers.

BLEU Bilingual Evaluation Understudy.

CAsT Conversational Assistance Track.

CIS Conversational Information-System.

ConRAG Conversational Retrieval-Augmented Generation.

Conv QA Conversational Question Answering.

CoQA Conversational Question Answering.

CoT chain-of-thought.

CQU Contextual Query Understanding.

DL Deep Learning.

DocVQA Document Visual Question Answering.

DPR Dense Passage Retrieval.

ER Examination Regulation.

FiD Fusion-in-Decoder.

HR Hit-Ratio.

IR Information Retrieval.

KB Knowledge Base.

LaPraDoR Large-scale Pretrained Dense Zero-shot Retriever.

LeoLM Linguistisch Erweitertes Offenes Language Model.

Llama Large Language Models with Adapters.

LLM Large Language Model.

LoRa Low-Rank Adaptation.

MaHuEval manual human evaluation.

MIT Massachusetts Institute of Technology.

ML Machine Learning.

MRC Machine Reading Comprehension.

MRR Mean Reciprocal Rank.

NN Neural Network.

OBQ Optimal Brain Quantization.

OCR Optical Character Recognition.

ODCQA Open-Domain Conversational Question Answering.

ODQA Open-Domain Question Answering.

OOD Out-of-Domain.

PEFT Parameter-efficient Fine-tuning.

PoC proof-of-concept.

PrLM Transformer-based Pre-trained Language Model.

QA Question Answering.

QG Automatic Question Generation.

QuReTeC Query Resolution by Term Classification.

List of Acronyms

R2-D2 Rank Twice, Read Twice.

RAG Retrieval-Augmented Generation.

REALM Retrieval-Augmented Language Model pre-training.

ROGUE Recall-Oriented Understudy for Gisting Evaluation.

seq-2-seq Sequence-to-Sequence.

TREC Text REtrieval Conference.

List of Figures

2.1. Adjusted QA Framework Classification by Farea et al. [Farea et al., 2022]	6
2.2. Reader-Retriever-System Architecture for QA by Zhu et al. [Zhu et al., 2021].	9
2.3. Types of Dense Retriever by Zhu et al. [Zhu et al., 2021].	15
2.4. Adjusted Graphic of the Extractive Reader by Jurafsky et al. [Jurafsky and Martin, 2023]	18
2.5. Overview of RAG by Lewis et al. [Lewis et al., 2020b]	19
2.6. Overview of FiD by Izacard et al. [Izacard and Grave, 2021]	20
2.7. Concepts of a Conversation in regards to a CIS	24
2.8. General System Architecture of a Conv QA System by Gao et al. [Gao et al., 2023]	25
2.9. Adapted Stages of Efficiency Improvement for LLMs by Treviso et al. [Treviso et al., 2023]	28
2.10. Adapted Fine-Tuning Approaches for Large Language Model (LLM) by Treviso et al. [Treviso et al., 2023]	29
2.11. General Idea of LoRa by Hu et al. [Hu et al., 2022]	31
2.12. Adapted Compression Approaches for LLM by Treviso et al. [Treviso et al., 2023]	32
3.1. Overview of the Example Use-Case	40
3.2. Overview of the System Architecture	41
3.3. Overview of the System Architecture in context of the sub-tasks of M	45
3.4. Fine-Tuning Process for Retriever	53
3.5. Scale of Passage implementations	56
4.1. Documents per Faculty of the German ER Dataset	62
4.2. Documents per Faculty of the English ER Dataset	63
4.3. Passages per Document of the German ER Dataset	63
4.4. Passage Length Distribution of the German ER Dataset	64
4.5. Passages per Document of the English ER Dataset	65
4.6. Passage Length Distribution of the English ER Dataset	65

List of Figures

4.7.	Adapted Evaluation Aspects and Components by TruLens [TruLens, 2024]	66
4.8.	Possible Implementations of the Extract Pipeline	73
4.9.	All Components of the System Architecture	73
4.10.	Implemented Extraction Pipeline	74
4.11.	Passage Length Distribution of the German ER Dataset before Filtering	76
4.12.	Implemented Retrievers	79
4.13.	Implemented Readers	81
4.14.	Implemented CQU	81
4.15.	Multiple Retriever Performance Comparisons	86
4.16.	Multiple Reader Performance Comparisons	88
4.17.	Error Type Distribution per Turn	90
4.18.	Error Distribution per Model	91
A.1.	Example Implementation Zero-Shot Baseline	96

List of Tables

4.1. Human-based Reader Evaluation for context length of 3 and no shuffled contexts	87
4.2. Answer Relevance per Model, Language and Question Type	89
4.3. Error-free turns per Model and Language	89
4.4. Average Context Relevance per Question Type	91
4.5. Multiple Metrics per Model and Language	92
A.1. Parameters used for Models in End-to-End Evaluation	98
A.2. Retriever comparison BM25 and DPR over indexEnglish and indexGerman	99
A.3. Retriever comparison BM25 and DPR over translated IndexEnglish and translated IndexGerman	100
A.4. Retriever comparison BM25-CE versions and DPR over IndexEnglish and IndexGerman	101
A.5. Retriever comparison BM25-CE versions and DPR over translated IndexEnglish and translated IndexGerman	102

1. Introduction

1.1. Motivation

The journey of exploring natural language querying dates back to as early as 1961, with researchers embarking on projects like Baseball [Green et al., 1961], a program designed to respond to users' natural language queries within the domain of baseball. In 1999, the Text REtrieval Conference (TREC) initiated the TREC-8 Question Answering track, marking "the first large-scale evaluation of domain-independent question-answering systems" [Voorhees, 1999]. A more renowned Question Answering (QA) system is IBM's *Watson*, an open-domain QA system that famously triumphed on the television game show Jeopardy! in 2011 [Ferrucci, 2012]. The advent of conversational assistants such as Alexa, Siri, and Cortana further fueled interest in conversational information seeking. In 2019, the establishment of the TREC Conversational Assistance Track (CAsT) aimed to nurture Conversational Information-System (CIS) as an active research field and provide large-scale reusable test beds. However, the true surge in user engagement came with the release of ChatGPT by OpenAI in November 2022, which achieved an astounding one million users within five days [Shewale, 2024].

Since the widespread adoption of ChatGPT, interest in conversational information seeking has surged once again. However, using a generative LLM only has certain drawbacks, including the following:

1. Parameterized Knowledge: Expansion or updating of knowledge is challenging.
2. Explainability: Understanding the rationale behind predictions is difficult.
3. Hallucination: The model may invent seemingly factual information not present in the underlying knowledge source.

These limitations render the sole use of generative models inadequate for practical conversational information-seeking use cases. Consequently, researchers have explored various approaches to address these issues and develop conversational information-seeking systems that mimic human-like conversations [Ferrucci, 2012, Guu et al., 2020, Lewis

1. *Introduction*

et al., 2020b, Nakano et al., 2022]. Among these approaches, Retrieval-Augmented Generation (RAG) [Lewis et al., 2020b] has emerged as a leading solution, evident in trending frameworks like Langchain [Langchain, 2023] and even ChatGPT itself [OpenAI, 2023].

However, despite advancements in research, there is a notable absence of papers detailing the adoption of these techniques in real-world use cases. To our knowledge, only two such papers exist [Feng, 2021, Gholami and Noori, 2021], which do not comprehensively address the entire problem domain associated with this task, such as the extraction step, and fail to leverage modern possibilities with LLMs.

1.2. Objectives and Contributions

This thesis aims to bridge this gap by addressing the question of how to build a faithful conversational question-answering system for real-world use cases, a topic that has not been adequately explored in scientific literature to date. This thesis resolves this issue by firstly laying out the fundamental concepts and techniques necessary, secondly developing a holistic framework for constructing such a system based on RAG, and thirdly implementing a proof-of-concept (PoC) based on the examination regulations of Heidelberg University. Finally, we evaluate the PoC using synthetic data and human evaluation. The core contributions of this thesis are as follows:

- Providing a holistic overview of the challenges in conversational question-answering systems tailored to a specific document collection and developing a comprehensive framework for constructing such systems based on the four components: Retriever, Reader, CQU and Extract.
- Breaking down the extraction of passages from a document collection into pipeline operations.
- Identifying challenges faced by the Retriever component towards both the knowledge base and the evidence set and how to resolve them using synthetic-data-based fine-tuning and Mixture-of-Experts.
- Decomposing the task of the Reader component into micro-challenges and proposing solutions for them.
- Offering multiple approaches for evaluating conversational question-answering systems, with the major insight, that the component-wise evaluation for non-factoid questions needs new approaches and currently at least in the case of this thesis, the human end-to-end evaluation was the most robust approach.

1. Introduction

- Highlighting the limitations of synthetic data for system evaluation, as those can only be used for factoid question generation due to the used pattern of question-context-answer tuples.
- Providing insights into the current bottlenecks of use-case-specific Conversational Question Answering (Conv QA) systems. Mainly the retriever component and the associated extraction of passages are drawbacks. In the end-to-end evaluation, 46.5% of errors were attributed to the retriever component and 18.5% to the extraction component directly.

Those cover the main contributions and insights of this thesis, for detailed insights, please refer to the respective chapters.

1.3. Thesis Structure

The construction of a conversational question-answering system is approached in two stages within this thesis. Firstly, we focus on question-answering based on a single natural language query, followed by the addition of the conversational component. Thus, Chapter 2 provides a comprehensive overview. Section 2.1 presents various approaches and concepts in the field of question-answering, while Section 2.2 extends this perspective to include conversational elements. Narrowing down the range of potential solutions, Chapter 3 introduces the Conversational Retrieval-Augmented Generation (ConRAG) framework, which lays out the spectrum of challenges, from handling document collections to enabling conversational question-answering. It's important to note that we streamline the choice of possible implementations to a system consisting of four distinct components: Extract, CQU, Retriever and Reader. The validation of this newly established system approach is conducted in a PoC manner in Chapter 4. Here, we implement and evaluate an exemplary system using the collection of examination regulations from Heidelberg University. This chapter also includes insightful reflections on the evaluation of conversational question-answering systems. This concludes the content outline of this thesis.

2. Background and Related Work

This chapter provides essential background information and reviews relevant prior research. It commences with an introduction to the sub-task of QA, as presented in Section 2.1. As mentioned in Chapter 1, this chapter maintains a clear distinction between QA and Conv QA. Consequently, Section 2.2 extends upon the foundational knowledge of QA and introduces the requisite concepts for the transformation of a QA-System into a Conv QA-System. Section 2.3 elaborates on different approaches towards the application of LLMs in resource-constrained settings. Section 2.4 delves into the related work, providing a comprehensive overview of the current state-of-the-art in the field of QA and Conv QA over textual knowledge sources. The goal of this chapter is to provide a holistic foundation of the research field of QA and Conv QA. Not all mentioned concepts and methods will further be utilized in Chapters 3 and 4 but are a necessary part of the research to understand the advantages and limitations of modern approaches and system design decisions. The introductions to certain sections will elaborate on the relevance of the respective topics for the research of this thesis.

2.1. Question Answering

In Section 2.1.1 we will lay the groundwork by introducing the fundamental aspects of QA-Systems and the techniques used to differentiate and categorize them. Following that, subsequent sections will delve deeper into the examination of specific system components.

The concepts and taxonomy of QA systems mentioned in Section 2.1.1 are essential background knowledge. While Section 2.1.2 outlines different conceptual approaches towards the holistic architecture of QA-systems, this thesis will later focus solely on RAGs. The following sections, 2.1.3, 2.1.4, and 2.1.5, will give a holistic image of possible approaches towards these components, which is necessary to understand on a higher level what possible system designs entail. However, the thesis will mainly focus on zero-shot retrievers, their optimization, and generative readers. The other mentioned methods can be viewed as possible alternatives and historical developments that should not be forgotten in the research process. The decision to include this information was

2. Background and Related Work

motivated by the dynamic nature of the research field of QA. While certain methods may have gained popularity in recent months, it's essential to recognize that the landscape of QA is continually evolving. Just because a particular method has gained traction does not render other approaches obsolete. By acknowledging the ongoing development and diversity within the field, we aim to provide a comprehensive understanding and context in this thesis.

2.1.1. Basics

Jurafsky and Martin define a QA-System as a system “designed to satisfy human information needs” [Jurafsky and Martin, 2023]. Hence, it primarily functions as an Information Retrieval System, with its primary objective being to provide users with the desired and accurate information in response to natural language requests.

The research community has yet to establish a universally accepted classification framework for QA systems. For instance, Hao et al. and Farea et al. [Hao et al., 2022, Farea et al., 2022] take a comprehensive approach to classify QA systems but differ in certain aspects, such as their treatment of question types and knowledge sources. On the other hand, other researchers [Zhu et al., 2021, Jurafsky and Martin, 2023, Etezadi and Shamsfard, 2023, Zhang et al., 2023b] employ a similar classification methodology but often focus solely on retrieval-based approaches, thereby lacking a holistic perspective.

The classification proposed by Farea et al. [Farea et al., 2022] goes a step further by distinguishing between the **QA-Framework** and **QA-Paradigms**, enhancing its versatility for comparing classical and modern QA systems. An adaptation of this classification will be utilized in this thesis. The originally proposed QA algorithms have been extended to include the Retrieval-based approach, and the Question Types have been revised based on the typology introduced by Mishra et al. in their 2016 survey [Mishra and Jain, 2016], which was further elaborated upon by Etezadi et al. [Etezadi and Shamsfard, 2023]. Also the answer types were adjusted to align with the classifications used in [McDonald et al., 2022, Dasigi et al., 2021]. In this context, a crucial distinction is made between a **QA** and **ConvQA** system, guided by the criteria outlined in [Zamani et al., 2023]: a **QA** system exclusively handles standalone questions, while any inquiry exceeding a single question and involving conversational context falls within the domain of a **ConvQA** system.

The **QA-Framework** encompasses external factors such as Question and Answer Types, while also considering system-related factors like the QA Algorithm and Knowledge Source [Farea et al., 2022, Hao et al., 2022]. Conversely, the **QA-Paradigm** defines the fundamental underlying concept of a system and can be seen as a subset of possi-

2. Background and Related Work

ble combinations within the **QA Framework**. Currently, three dominant paradigms prevail:

1. **Information Retrieval (IR)-Based QA**: This paradigm involves searching through extensive multi-modal data based on a user's question and using the retrieved passages to generate an answer.
2. **Knowledge Base (KB) QA**: In this approach, a semantic representation of the question is constructed, and a knowledge base is queried using this representation. The returned results are then used to generate an answer.
3. **Generative Question Answering**: Here, knowledge is fully implicit, and a Neural Network (NN) generates answers based on its trained parameters.

For visual clarity, a diagram illustrating the adjusted QA Framework Classification by Farea et al. is provided in Figure 2.1.

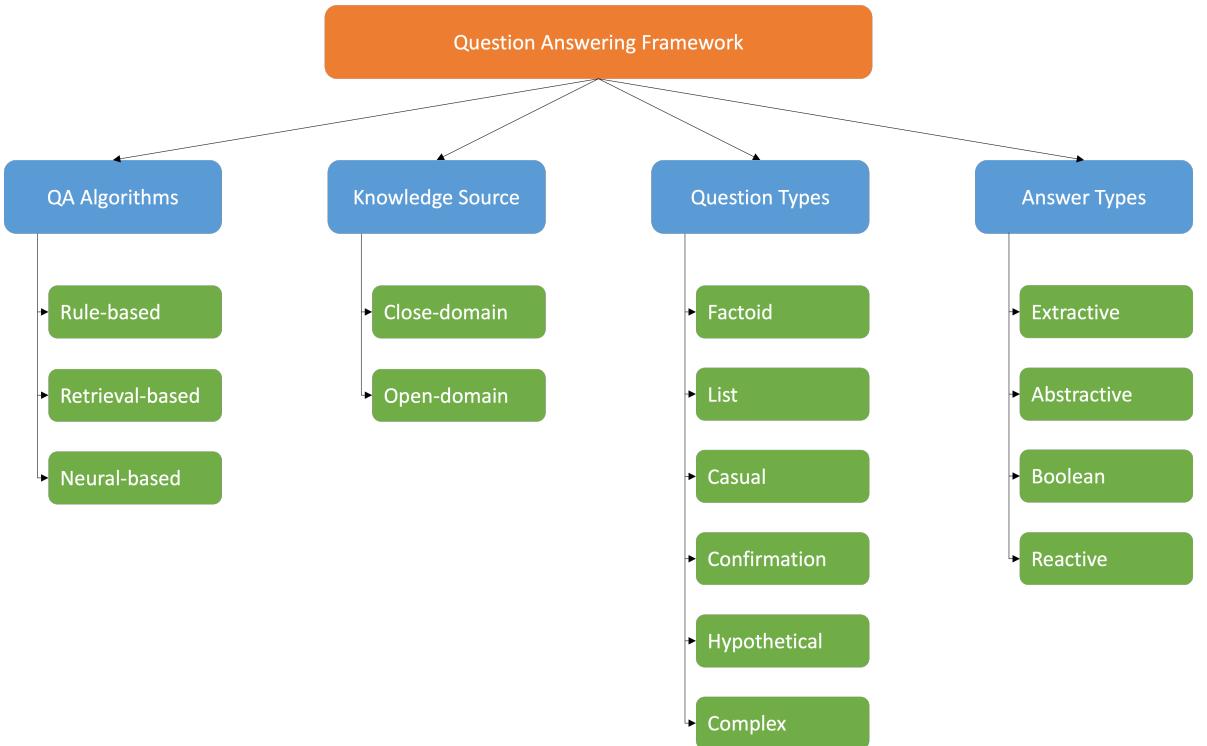


Figure 2.1.: Adjusted QA Framework Classification by Farea et al.
[Farea et al., 2022]

Figure 2.1 illustrates the aforementioned classification. The primary distinguishing factor is the employed **QA Algorithm**. Rule-based approaches involve the manual crafting of feature extractions from user questions, which are then compared to the

2. Background and Related Work

knowledge base. Rule-based approaches are typically employed in closed-domain QA systems exclusively [Etezadi and Shamsfard, 2023].

Retrieval-based approaches are the classic Information Retrieval (IR)-based QA systems, comprising two key components: an intent classifier and a retriever. The intent classifier’s objective is to discern the question’s intent and identify important entities. Subsequently, the retriever searches the knowledge source and identifies the most relevant passages [Farea et al., 2022, Zhu et al., 2021].

The Neural-based approach, often referred to as the generative approach, utilizes a Sequence-to-Sequence (S2S) model to generate accurate answers to given questions. In this paradigm, the information is stored directly in the neural network’s parameters, otherwise, the neural network is part of a Retrieval-based approach. Most datasets in these contexts consist of triples of question, context, and answer pairs [Jurafsky and Martin, 2023]. Notably, widely used datasets such as SQuAD and QASPER originally emerged from the field of machine reading comprehension, representing a foundational step in the evolution of QA systems [Rajpurkar et al., 2016, Dasigi et al., 2021, Zhu et al., 2021].

In addition to the **QA Algorithms**, the **Knowledge Source** plays a pivotal role in distinguishing various aspects of Question Answering (QA) systems. The nature of the knowledge source can range from structured to unstructured or semi-structured, and it may encompass diverse data modalities, including text, audio, and video. A common point of comparison in the QA landscape is between closed and open-domain systems.

In the broad sense, a **closed-domain** QA system operates within the confines of a specific knowledge domain, which means it has limited access to information. In contrast, **open-domain** QA systems grapple with an extensive array of knowledge sources, necessitating a more versatile approach [Farea et al., 2022].

Furthermore, a closed-domain setup often entails limitations on the types of questions it can handle, primarily focusing on factoid questions or predefined templates. Additionally, it frequently relies on structured knowledge bases like graphs or logically organized repositories [Hao et al., 2022].

Conversely, open-domain QA systems are designed to tackle a wide spectrum of user queries, ranging from factoids to more complex inquiries. They typically deal with unstructured knowledge sources, which can be substantial and diverse in content [Zhu et al., 2021, Farea et al., 2022, Jurafsky and Martin, 2023].

An alternative perspective for distinguishing QA-Systems lies in the **Question Types** that users can input into the system. Questions can fall into various categories, such as *factoid, list, casual, confirmation, hypothetical* [Mishra and Jain, 2016], or *complex* [Etezadi and Shamsfard, 2023].

2. Background and Related Work

- *Factoid questions*, the most common type, is typically signaled by question words (what, when, which, who, how) and yields a concise factual answer.
- *List questions* represent a specialized subset of factoid questions, where the answer comprises a list of facts.
- *Casual questions* encompass inquiries that deviate from the factoid format, often involving words like *how* or *why* and requiring more advanced reasoning.
- *Confirmation questions* seek simple yes or no responses, frequently employed in personal assistant applications.
- *Hypothetical questions* delve into hypothetical scenarios (e.g., "what would happen if"), aiming for plausible rather than definitive answers.
- *Complex questions* can be further categorized into *answer-retrieval-complex* and *question-understanding-complex*. In the case of question-understanding-complex questions, the complexity arises from nuances like multiple constraints, making the question itself intricate to comprehend. In contrast, answer-retrieval-complex questions involve complexities in finding the correct answer, often requiring the combination of information from multiple documents or similar sources. This is commonly referred to as long-form QA.

Lastly, a QA-System can be characterized by the **Answer Types** it offers, a concept closely intertwined with Question Types. Farea et al. [Farea et al., 2022] delineate four categories of answers: *extractive*, *abstractive*, *boolean* and *reactive*.

- *Extractive answers* represent the most common type, where the answer is a specific factual excerpt presented as a span of tokens.
- *Abstractive answers* typically correspond to complex questions that necessitate the system to consider multiple documents and information sources to formulate a response. In such cases, no predefined or annotated answer exists.
- *Boolean answers* are typically the result of confirmation questions, where the answer is either *yes* or *no*.
- *Reactive answers* often arise in response to confirmation questions and can be a system-generated reaction based on the user's provided answer.

2. Background and Related Work

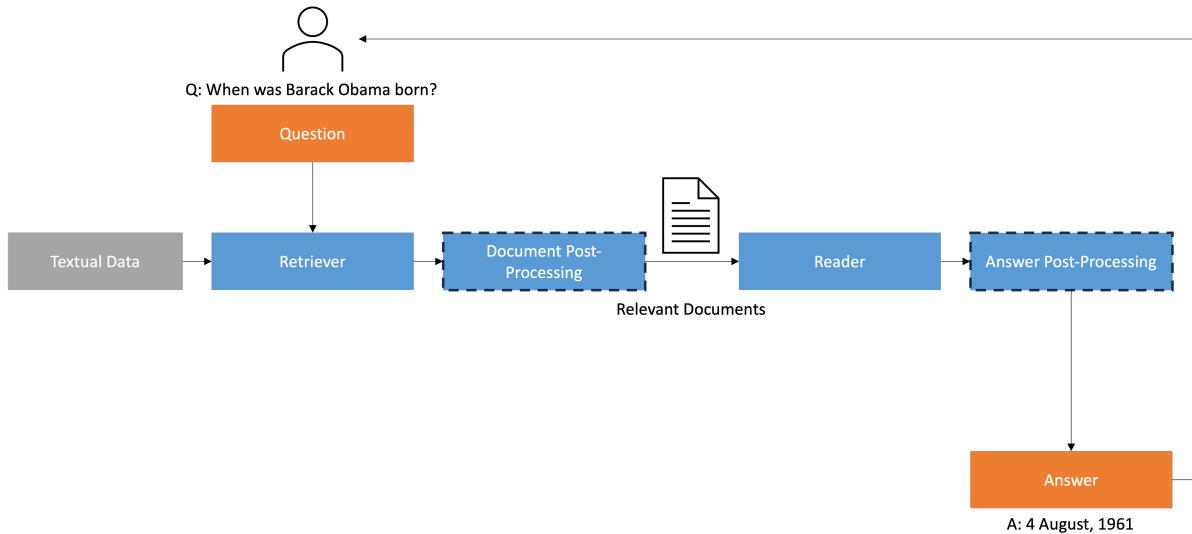


Figure 2.2.: Reader-Retriever-System Architecture for QA by Zhu et al. [Zhu et al., 2021].

2.1.2. Information Retrieval Architectures

As stated in the previous section, there are three major paradigms in QA: Information Retrieval (IR)-based QA, Knowledge Base (KB)-based QA, and Generative QA. This section will primarily concentrate on the first paradigm, IR-based QA, as it holds the most promise for addressing the objectives of this thesis topic.

This thesis will not focus on KB QA, as this approach requires the mapping of the query to a structured data representation. As the task of this thesis is to develop a general system, that is adaptable to different data inputs, KB QA will be excluded [Dimitrakis et al., 2020] (See Section 3.1).

Generative QA is often denoted as *Retriever-free* or *Neural-based* approaches. The central characteristic of this paradigm is that knowledge resides within the parameters of a neural network. Consequently, the knowledge is implicit, and the QA system will not furnish a specific document, passage, or other source from which it extracted the information. Instead, it offers a textual excerpt. While these systems can achieve competitive performance compared to IR-based QA systems, they are not under consideration for this thesis due to their lack of reference, which is a crucial requirement for the system (See Section 3.1) to be developed [Roberts et al., 2020].

Figure 2.2 depicts the general architecture of a **Retriever-Reader-System**, as defined by Zhu et al. [Zhu et al., 2021]. This architecture serves as the foundational framework for IR-Based QA systems and was initially introduced by Harabagiu et al.

2. Background and Related Work

[Harabagiu et al., 2003]. In this framework, all modules operate independently, can be trained separately, and are subject to independent evaluation.

The **Retriever** module’s primary role is to retrieve relevant documents, passages, or other pertinent information from a knowledge source and rank them based on their relevance to answering the user’s query. Subsequently, the **Reader** module extracts the answer from the retrieved documents and presents it to the user. This task bears a close resemblance to Machine Reading Comprehension (MRC), with the key distinction that in IR-Based QA, the system must handle multiple documents and comprehend them to formulate a response, unlike classical MRC tasks, which typically involve only one context document.

The **Document Post-Processor** module’s role is to curate and refine the set of documents that will be forwarded as *Relevant Documents* to the subsequent stage, the Reader. Concurrently, the **Answer Post-Processor** assists the Reader in addressing complex questions for which the answer may not be found in a single document alone [Zhu et al., 2021, Jurafsky and Martin, 2023].

It’s worth noting that some researchers include a **Question Analysis** module preceding the Retriever, which aims to preprocess the received question for more efficient query execution in the Retriever [Nassiri and Akhloufi, 2023]. However, for the purposes of this thesis, we adhere to Zhu et al.’s definition [Zhu et al., 2021], where this functionality is considered part of the Retriever.

Conceptually, there are three distinct approaches to the Retriever itself: *Sparse Retrieval*, *Dense Retrieval*, and *Iterative Retrieval*. The specifics of these approaches will be thoroughly explored in Section 2.1.4.

Document Post-Processors can be categorized into *Supervised Learning*, *Reinforcement Learning*, and *Transfer Learning*-based approaches. A detailed discussion of these approaches is also provided in Section 2.1.4.

In Section 2.1.5, we will delve into the finer details of Reader approaches and Answer Post-processing. Broadly speaking, there are two primary types of Readers: *Extractive* and *Generative* Readers. As for Answer Post-processing, it involves two key categories: *Rule-based* and *Learning-based* approaches.

There are also **End-to-End** approaches that employ a single module to execute the entire QA task. Excluding generative approaches, two common categories of such approaches are **Retriever-Reader** and **Retriever-only** models.

An End-to-End Retriever-Reader aims to train both the Retriever and Reader in a single backpropagation step, and in some cases, it introduces additional knowledge sources beyond the traditional IR framework. An illustrative example is RAG [Lewis et al., 2020b]. RAG consists of a pre-trained Generator with implicit knowledge encoded

2. Background and Related Work

in its parameters and a pre-trained Retriever. For each question, the Retriever identifies the most relevant documents and generates a latent vector based on them. This latent vector, along with the original question, is fed into the Generator. Section 2.1.5 will delve into details regarding the RAG architecture.

Another end-to-end approach, similar to RAG, is Retrieval-Augmented Language Model pre-training (REALM) [Guu et al., 2020]. While these previous two approaches extended the capabilities of pre-trained Sequence-to-Sequence (seq-2-seq) models, Nishida et al. pursued a different path by training a single NN to perform both tasks simultaneously: IR and MRC [Nishida et al., 2018].

It is noteworthy that all these end-to-end approaches have demonstrated competitive performance compared to state-of-the-art methods on specific QA datasets.

An essential yet often underestimated question is: What defines textual data, and how should one preprocess formats such as PDFs to extract this textual content? While many datasets already comprise small contextual snippets [Wang, 2022], it's crucial not to overlook the entire process of extracting snippets from unstructured PDFs, for example. Approaches to tackle this challenge will be explored in detail in the upcoming Section 2.1.3.

2.1.3. Extraction Approaches

As discussed in the previous Section 2.1.1, the knowledge source for a QA-System can take the form of textual or multimodal data. The specific type of data may necessitate certain requirements or specific adjustments to the Retriever used for IR.

In the context of this thesis, the primary knowledge source to be employed is PDF documents (See Section 3.1). In the research field, three major approaches exist for extracting textual information from unstructured data types like PDFs: *visual* [Tito et al., 2021], *direct* [Wang et al., 2019], and *alternative* [Dasigi et al., 2021] extraction methods.

It's important to note upfront that the chosen extraction method is intricately connected to the subsequent retrieval approach. The specifics, including metadata alongside pure textual data and quality requirements, may vary among different extraction and retrieval methods.

The visual approach is closely aligned with the research field of *Document Question Answering*. A well-known example dataset in this field is Document Visual Question Answering (DocVQA) [Tito et al., 2021]. The primary concept behind the visual approach to document question-answering is to capture not only the text of a PDF but also additional information such as the document's structure, various hierarchies on a

2. Background and Related Work

page (e.g., sections, subsections), and the ability to analyze tables and figures. These hierarchical structures can be leveraged to create two-stage retrieval approaches. In these approaches, initially, a collection of relevant files is identified based on higher-level attributes like the document’s title and abstract. Subsequently, a more granular retrieval process is executed over lower-level attributes such as passages within the relevant files. This *Iterative Retrievers* will be further discussed in Section 2.1.4 [Liu et al., 2021].

The challenge of *Visual Document Question Answering* typically involves taking images of PDF pages as inputs and mapping question-answer pairs to them. The answers are extracted from either a single paragraph or a combination of multiple paragraphs [Mathew et al., 2021]. Nonetheless, the extraction pipeline in this case usually resembles the *Retriever-Reader* architecture, where the extracted information from the visual processing is fed into such a system afterward. Researchers in this field often employ a pipeline that includes a *Document Layout Analysis* model, followed by the application of an Optical Character Recognition (OCR) tool to the detected regions [McDonald et al., 2022]. Examples of a *Document Layout Analysis* model include the Document Image Transformer by Li et al. [Li et al., 2022a].

The direct approach is the most prevalent method in the field of Question Answering (QA) and Information Retrieval (IR). The primary concept behind this approach is to extract textual information from PDFs and store it in a database. The extraction process can be accomplished using various tools such as *PDFMiner* or *Adobe Extract* [Meuschke et al., 2023]. However, a lingering question is how to effectively split the extracted textual data, especially considering that they are often not cleaned after extraction.

A common practice when employing a Language Model (LLM) is to optionally cleanse the text corpus and then divide it based on a predefined token size. This approach is evident in two notable open-source LLM projects: *Langchain* and the *Retrieval Plugin for ChatGPT* by OpenAI [Langchain, 2023, OpenAI, 2023]. In the original Dense Retrieval paper by Karpukhin et al., a sliding window of token size 5 was utilized [Karpukhin et al., 2020]. Therefore, it can be assumed that for contemporary LLM applications, the precise quality of the data, ensuring that a document contains syntactically correct sentences, may not be as critical.

Apart from modern approaches involving text clipping, previous methods aimed to identify paragraphs and similar structures within the extracted texts [Zhu et al., 2021].

An alternative approach involves the methodology employed in constructing the QASPER dataset. In this case, the authors conducted a pre-filtering of scientific papers’ PDFs, selecting only those with freely accessible LaTeX files. They then utilized the S2ORC tool to extract cleaned textual data from these LaTeX files [Dasigi et al., 2021]. It’s important to note that this approach is highly specific to the QASPER dataset and

2. Background and Related Work

cannot be universally applied. Nonetheless, it serves as an illustration of alternative methods for extracting textual data from PDFs.

2.1.4. Retrieval Approaches

The traditional state-of-the-art in IR relies on **Sparse Retrievers**, with one notable example being BM25. BM25 is renowned as "one of the most empirically successful retrieval models and is widely used in current search engines" [Zhu et al., 2021]. Nandan et al. even demonstrated that on modern Open-Domain Question Answering (ODQA) datasets, BM25 remains a viable baseline for zero-shot IR [Thakur et al., 2021].

BM25 was originally introduced by Robertson et al. [Robertson and Zaragoza, 2009]. It operates by utilizing the TF-IDF token weights between a question q containing tokens q_1, \dots, q_T and a set of passages P , where $p \in P$.

$$\mathbf{s}_{q,p}^{\text{BM25}} = \sum_{i=1}^T \log \left(\frac{|\mathcal{P}|}{N(q_i, \mathcal{P})} \right) \frac{n(q_i, p)(k_1 + 1)}{k_1 \left(1 - b + \frac{b|p|}{avpl} \right) + n(q_i, p)} \quad (2.1)$$

Equation 2.1 illustrates the BM25 score for a question q and a passage p . In this equation, $N(q_i, \mathcal{P})$ represents the count of passages in \mathcal{P} that contain the token q_i , while $n(q_i, p)$ indicates the frequency of token q_i within the passage p . The variable $|p|$ signifies the length of passage p , and $avpl$ stands for the average passage length in \mathcal{P} . The parameters k_1 and b are free parameters, typically set to $k_1 = 0.9$ and $b = 0.4$ [McDonald et al., 2022, Robertson and Zaragoza, 2009].

Traditionally, this lexical Information Retrieval (IR) approach has been capable of providing satisfactory retrieval results. However, in 2020, Karpukhin et al. demonstrated for the first time that a **Dense Retrieval** approach could outperform the Sparse Retrieval approach across multiple ODQA datasets [Karpukhin et al., 2020]. Consequently, the search for a general Dense Retrieval model has been ongoing, as these Dense Retrieval approaches offer advantages such as semantic matching and the ability to handle lengthy documents.

In general, there are three types of Dense Retrieval approaches [Zhu et al., 2021]: the **Representation-based Retriever**, often referred to as the *dual-encoder* [Karpukhin et al., 2020]; the **Interaction-based Retriever**, often referred to as the *cross-encoder*; and the **Representation-interaction Retriever**, often referred to as the *multi-stop retriever*. Figure 2.3 illustrates the general architecture of these three types of Dense Retrievers.

The **Dense Passage Retriever (DPR)** by Karpukhin et al. serves as a notable example to explain the **Representation-Based Retriever**. Given a collection M of

2. Background and Related Work

text passages p and a question q , the objective of DPR is to identify the k most similar passages to the question. To achieve this, DPR employs two distinct **BERT** [Devlin et al., 2019] Encoders. One Encoder, denoted as $E_Q(\cdot)$, encodes the question q into a d -dimensional vector, where $d = 768$. The other Encoder, labeled as $E_P(\cdot)$, encodes the passage p into a d -dimensional vector at the [CLS] token. The similarity between these two vectors is computed using the inner product:

$$\mathbf{s}_{q,p}^{DPR} = \mathbf{E}_Q(q)^\top \mathbf{E}_P(p) \quad (2.2)$$

The choice of the inner product as the similarity function is motivated by its computational efficiency and the demonstrated, comparable performance [Karpukhin et al., 2020]. The dot-product must yield a small value for pairs of questions and passages that are genuinely related. The training dataset D comprises m instances, where q_i represents the question, p_i^+ denotes the positive passage, and p_{i,n^-} represents the negative passage:

$$\mathbf{D} = \{(q_i, p_i^+, p_{i,1}^-, \dots, p_{i,n^-}^-)\}_{i=1}^m \quad (2.3)$$

The loss function is optimized using the negative log-likelihood of p_i^+ :

$$\mathcal{L}_{DPR} = -\log \frac{\exp(\mathbf{s}_{q_i, p_i^+}^{DPR})}{\exp(\mathbf{s}_{q_i, p_i^+}^{DPR}) + \sum_{j=1}^n \exp(\mathbf{s}_{q_i, p_{i,j}^-}^{DPR})} \quad (2.4)$$

It's important to note that in [Karpukhin et al., 2020], the selection of negative passages was not arbitrary. Instead, two additional approaches were employed: BM25 top passages that do not contain the answer and positive passages paired with other questions.

One significant advantage of the Representation-Based Retriever is that passages can be pre-indexed locally rather than at runtime. This reduction in latency between the question and the response may, however, come with trade-offs in the quality of the retrieved passages.

The **Interaction-Based Retriever** incorporates both the question q and the passage p within a single model, separated by a [SEP] indicator. These models offer various approaches for modeling the relationship between q and p . For instance, one common method is to utilize the [CLS] classifier as an indicator of whether the passage is relevant to the question. This approach was first introduced with Bidirectional Encoder Representations from Transformers (BERT) [Devlin et al., 2019]. While these models perform competitively with previous Representation-Based Retrievers, it's important

2. Background and Related Work

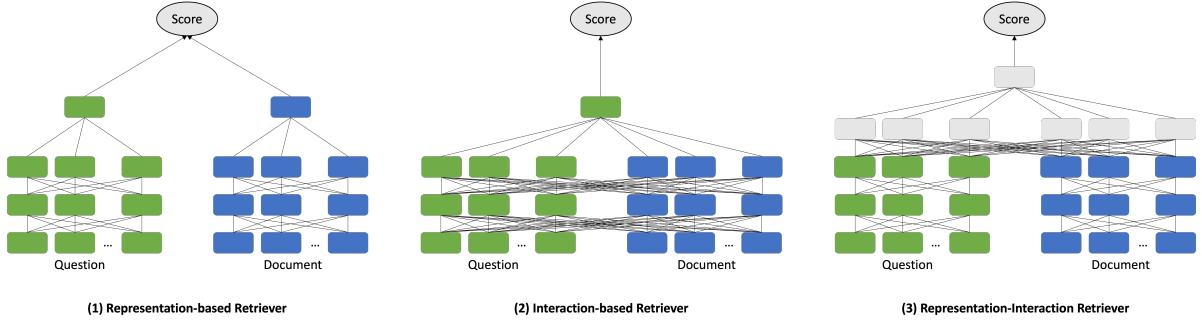


Figure 2.3.: Types of Dense Retriever by Zhu et al. [Zhu et al., 2021].

to note that they are 100-1000 times more computationally expensive [Khattab and Zaharia, 2020].

To address this latency issue, models like ColBERT introduced the concept of **contextualized late interaction** [Khattab and Zaharia, 2020]. In this thesis and subsequently in research, it is referred to as the **Representation-Interaction Retriever** [Zhu et al., 2021].

ColBERT, like Dense Passage Retrieval (DPR), employs two BERT Encoders, denoted as $E_Q(\cdot)$ and $E_P(\cdot)$. However, it introduces a late interaction mechanism. When provided with a query q , it is initially tokenized into BERT-based Wordpiece tokens, resulting in q_1, \dots, q_T . Following the [CLS] token, a [Q] token is appended to signify the question. If the length of the tokenized question is less than N_q , a predetermined token length, the remaining portion of the question is padded with BERT's [mask] token. Otherwise, it is truncated. This process, known as *query augmentation*, allows BERT to re-weight existing terms or expand the query, and it is pivotal to ColBERT's performance. The generated embeddings are then passed through a linear layer to reduce the output dimensions to a fixed size m , which is smaller than the original dimensions of BERT. The output is subsequently normalized to ensure that the L2 norm of each result equals one.

For each passage p , $E_P(\cdot)$ is employed for encoding. Similar to the question encoding process, p is segmented into its p_1, \dots, p_{T_d} Wordpiece tokens. The special token [D] indicates a passage. Short passages are not padded with a [mask] token. After the classical BERT output, a similar post-processing step is applied to the encoded passages, and all embeddings corresponding to punctuation are filtered out.

$$\mathbf{E}_q := \text{Normalize}(\text{CNN}(\text{BERT}([Q]q_0q_1 \dots q_T[mask] \dots [mask]))) \quad (2.5)$$

$$\mathbf{E}_p := \text{Filter}(\text{Normalize}(\text{CNN}(\text{BERT}([D]p_0p_1 \dots p_T)))) \quad (2.6)$$

2. Background and Related Work

The late interaction mechanism applied to the encodings involves computing the maximum similarity, which utilizes cosine similarity through dot products. This is made possible by the earlier normalization applied to the embeddings:

$$\mathbf{s}_{q,p}^{ColBERT} = \sum_{I \in [\|\mathbf{E}_q\|]} \max_{j \in [\|\mathbf{E}_d\|]} \mathbf{E}_{q,i} \cdot \mathbf{E}_{p,j}^\top \quad (2.7)$$

The interaction mechanism has no trainable parameters. ColBERT is differentiable end-to-end. During training, for example, with a triple (q, p^+, p^-) , ColBERT independently produces a score for each passage and is subsequently optimized pairwise using softmax cross-entropy loss over the scores of p^+ and p^- [Khattab and Zaharia, 2020].

Another type of Retriever is the **Iterative Retriever**. Iterative Retrievers are necessary when dealing with questions that are more complex than simple factoid questions, which can be answered by identifying the right passage in the knowledge source. An example is the HotpotQA dataset [Yang et al., 2018], designed specifically for multi-hop questions. The fundamental concept here is that such questions cannot be answered with just one precise piece of evidence. They require multiple passages from different documents at the very least. Iterative Retrievers encompass three stages in the pipeline: (1) document retrieval, (2) query reformulation, and (3) retrieval stopping.

An example is BEAM, currently holding the title of the highest-performing¹, QA-System across multi-hop QA datasets such as HotpotQA [Zhang et al., 2023a]. The document retrieval component can take the form of any retrieval model, including options like ColBERT, BM25, or DPR. In the case of BEAM, it leverages an Interaction-Based Retriever using DeBERTa. For each candidate passage p_c , BEAM calculates a relevance score concerning this passage within the context of all previously identified relevant passages p_r and the question q , using the embeddings of the [CLS] tokens [He et al., 2021]. The second step, query reformulation, can be executed explicitly or implicitly, meaning it can either be expressed in natural language or as a dense embedding. The advantage of using natural language lies in its interpretability; while employing dense embeddings operates within a semantic space and does not lack vocabulary interpretability [Zhu et al., 2021]. BEAM adopts a natural language-based approach. Specifically, after each hop, it appends the newly identified passage to the previously identified ones and feeds this information into DeBERTa:

$$\mathbf{s}_{q,p}^{BEAM} = \text{Classifier}(\text{DeBERTa}([CLS]q p_{r_1} \dots p_{r_i})) \quad | \quad p_c \in P \quad (2.8)$$

The nature of query reformulation depends on the type of retriever in use. Lastly

¹Status as of September 23, 2023, according to <https://paperswithcode.com> and the authors of [Zhang et al., 2023a]

2. Background and Related Work

retrieval stopping poses its own set of challenges. A common approach involves setting either a fixed number of hops or a maximum limit on the retrieved documents. Alternatively, some methods introduce a new token, such as [EOE] (End-of-Evidence), to signal the end of retrieval [Zhu et al., 2021]. BEAM, for example, employs a fixed number of hops, specifically 2, as determined through empirical evaluation.

The task of **Document Post-Processing** is to reduce the number of passages forwarded to the Reader, aiming to eliminate irrelevant ones. Traditional Retrievers, like Sparse Retrievers, often require a Document Post-Processor. However, Dense Retrievers often incorporate ranking and retrieval simultaneously, rendering this module unnecessary [Zhu et al., 2021]. Nevertheless, it remains possible to construct multi-stage Retrievers. This can be achieved by using a simpler Dense Retriever for pre-filtering passages and subsequently applying a more accurate one [Liu et al., 2021].

2.1.5. Reader Approaches

Readers originally emerged from the field of MRC, where the objective is to extract an answer from a given context. A well-known example is the SQuAD [Rajpurkar et al., 2016] dataset, which was mentioned in Section 2.1.1. However, unlike the original MRC task, a Reader in a Retrieval-Reader-System must process multiple passages to determine the relevant information needed to answer a given question [Zhu et al., 2021].

Modern readers rely on Transformer-based Pre-trained Language Model (PrLM)s since they establish new baselines on well-known datasets [Luo et al., 2022]. In general, there are two types of Readers that use PrLMs: **Extractive Readers** and **Generative Readers** [Jurafsky and Martin, 2023, Zhu et al., 2021, Luo et al., 2022].

In general, an **Extractive Reader** employs an encoder to identify the token sequence span that is relevant for answering a question. These encoders can be any autoencoder models, such as BERT [Devlin et al., 2019], DeBERTa [He et al., 2021], or RoBERTa [Liu et al., 2019]. Luo et al. [Luo et al., 2022] even utilized the encoder components of established encoder-decoder models like T5 [Raffel et al., 2020] and BART [Lewis et al., 2020a]. They demonstrated that, after fine-tuning, these models can outperform encoder-only models on certain tasks.

Figure 2.4 illustrates the span labeling process performed by the extractive reader. The question tokens q_1, \dots, q_n and the passage tokens p_1, \dots, p_m are input into the encoder, separated by a [SEP] token. The encoder learns two new embeddings, S and E , which represent span-start and -end tokens, respectively. To obtain the span start probability for an output token p'_i , the dot product between the output token and S is computed and then normalized by a softmax function over all output tokens. The

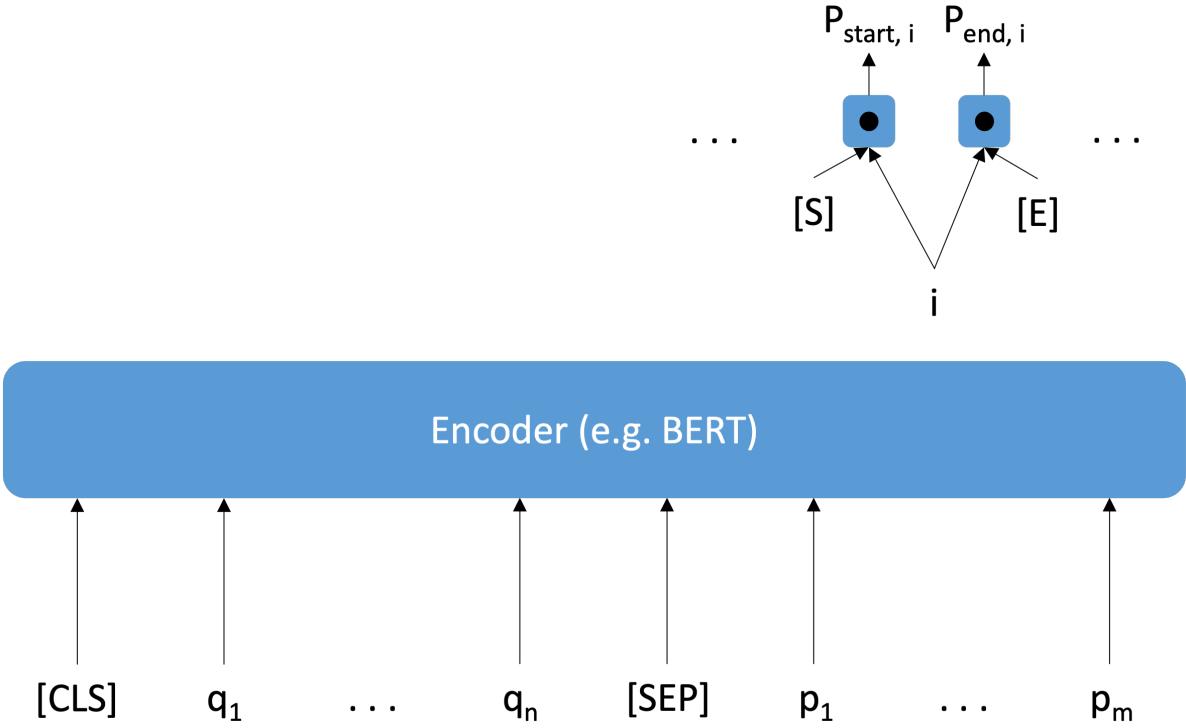


Figure 2.4.: Adjusted Graphic of the Extractive Reader by Jurafsky et al. [Jurafsky and Martin, 2023]

process is similar for the span-end token. The score of a span from position i to j is calculated as $S * p'_i + E * p'_j$. The span with the highest score, where $j \geq i$, is selected as the answer span. If the total length of tokens in q and p exceeds the maximum input length of the encoder, the passage is split into multiple segments, and the process is repeated for each segment [Jurafsky and Martin, 2023, Luo et al., 2022].

The **Generative Reader** operates straightforwardly when familiar with a seq-2-seq encoder-decoder model. Given a dataset containing (q, p, a) tuples, the encoder takes q and p as input and outputs the contextual representation h . Then, it is the decoder's task to generate a token sequence based on h and attention. The training objective can be described as minimizing the following loss function:

$$\mathcal{L}_{\text{Gen}} = \sum_{i=1}^K \log P(a_i | h, a_{:i}) \quad (2.9)$$

Here, K represents the length of tokens in a , a_i is the i^{th} token in a , and a_0 is a special beginning of sequence token. In cases where the answer is not contained within the passages, the [CLS] token indicates this situation [Luo et al., 2022, Zhu et al., 2021].

Latest research projects like Visconde [Pereira et al., 2023] even employ LLM as

2. Background and Related Work

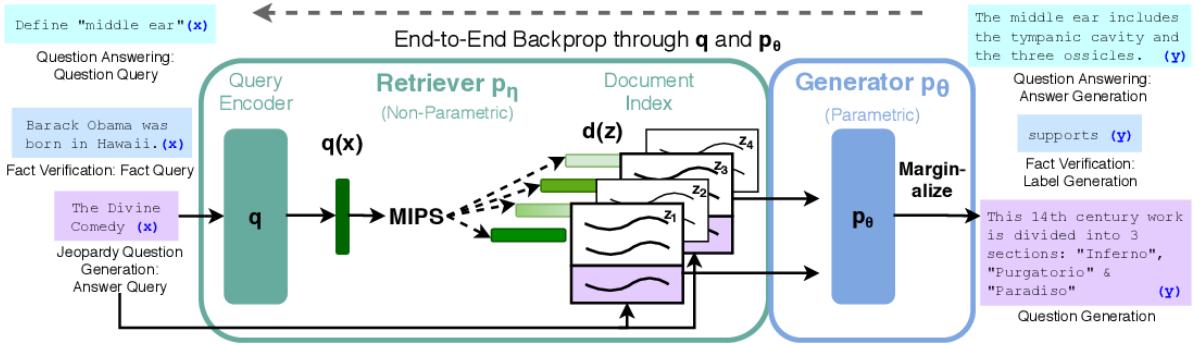


Figure 2.5.: Overview of RAG by Lewis et al. [Lewis et al., 2020b]

Generative Readers. The performance and usability of these models remain active topics of research.

Luo et al. conducted the first survey comparing state-of-the-art Extractive and Generative Readers [Luo et al., 2022]. They discovered that “on average, extractive readers perform better than generative ones” [Luo et al., 2022], except in cases involving long context passages, where generative approaches outperform the extractive ones.

RAG can be seen as a generative reader, but with a much more capable NN as the reader, specifically the idea is that the reader itself is a LLM with implicit knowledge encoded in its parameters, which it uses to generate an answer, the retrieved passages intentionally function as support in order to guide the reader and reduce risk of hallucination.

Figure 2.5 is taken from the original paper by Lewis et al. [Lewis et al., 2020b] and displays the general approach of RAG. The original idea of RAG is to have an end-to-end backpropagation in order to train the retriever and reader (generator) at once and on the same data, not separate as in most Retriever-Reader-Systems. The used retriever in the original RAG is a DPR. Other retrievers can be used, as this is just a decision to make as the generator does not directly depend on the type of retriever. More interestingly is the kind of implementation of the generator. RAG implements a *sequence-based* generator, while future work, such as Fusion-in-Decoder (FiD) [Izacard and Grave, 2021] use an *attention-based* generator. The sequence-based generator works the following way: Given an arbitrary encoder-decoder $p_\theta(y_i|q, p, y_{1:i-1})$, the query q , the k -relevant passages p , and the previously generated tokens $y_{1:i-1}$, the generator computes the probability distribution over the next token y_i . q and p are simply concatenated.

Further RAG generators are attention-based like FiD [Izacard and Grave, 2021]. Here the encoder and decoder of the generator are slightly decoupled as to the classic RAG. Given a question q , the retriever retrieves the top- k passages p . The encoder encodes every single passage in a question, title, passage triple (q, t, p) . The encodings of multiple

2. Background and Related Work

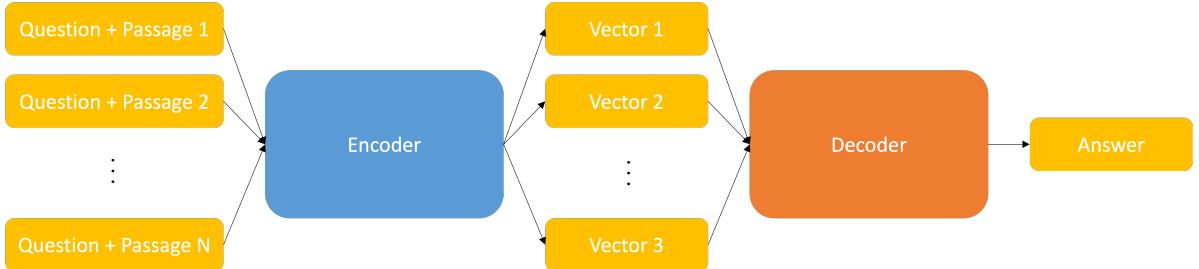


Figure 2.6.: Overview of FiD by Izacard et al. [Izacard and Grave, 2021]

passages are afterward concatenated and passed into the encoder-decoder attention of the decoder. An illustration can be found in Figure 2.6 This allows for the combination of multiple passages, so there is no input token limitation as for the classical RAG. Also, experiments by Izacard et al. showed, that the performance improves over multiple tasks, as multi-passage relations can easily be resolved by the decoder. The latest work of Izacard et al. is ATLAS, which set new state-of-the-art benchmarks on multiple evaluation tasks [Izacard et al., 2023]. ATLAS extends on the idea of FiD.

Still overall RAG approaches, the main idea is to have a fully end-to-end backpropagation during training or fine-tuning of the systems and a LLM whose generation is supported by passages.

The **Answer Post-Processor** is similar to the Document Post-Processor, serving as an optional component. Its primary task is to provide support for multi-hop complex questions, helping determine the final answer from a set of answers extracted by the reader component [Zhu et al., 2021]. Depending on the implementation of the Reader, this component may become obsolete.

2.1.6. Limitations

The evaluation metrics for IR systems will be discussed in detail in Chapter 4. In general, selecting the components and models for an IR system always involves a trade-off between accuracy, memory consumption, and inference speed [Zhang et al., 2023b].

Accuracy is primarily determined by the chosen Retriever-Reader-System. Sparse retrievers often lack a certain degree of semantic understanding, resulting in less accurate retrieved passages. In contrast, Dense Retrievers can achieve higher levels of accuracy but require thorough evaluation and training for the desired use case. Thakur et al. demonstrated that high-accuracy Dense Retrievers like DPR can underperform in zero-shot scenarios compared to BM25 by -47.7% [Thakur et al., 2021]. This highlights another crucial limitation of all NN-based retrievers and readers: training. BM25

2. Background and Related Work

is, by nature, an unsupervised model for IR, while common approaches for Dense Retrieval usually belong to the group of supervised models. These models heavily depend on their training data, whereas a Sparse Retriever like BM25 can be used without any training. According to experiments conducted by Thakur et al. [Thakur et al., 2021], the best-performing out-of-distribution standalone Retrievers are Representation-Interaction Retrievers like ColBERT. Nevertheless, there exist enhanced approaches like *Mixture-of-Experts/Hybrid-Search* or the multi-stage retrievers like *BM25+CE* (will be discussed in Section 3.3.3), which are an implementation of the mentioned *Document Post-Processing*. These approaches can outperform standalone Retrievers.

Constructing a training dataset for a QA task can be a tedious process, as these datasets must consist of tuples in the form of (question, context, answer), which is not always feasible. One established research direction to address this issue is Automatic Question Generation (QG) [Serban et al., 2016]. In QG, a seq-2-seq model is employed to generate questions and answers based on a given passage.

Zhang et al. provide an example of DPR applied to the Natural Questions dataset in their survey on efficient ODQA [Zhang et al., 2023b]. The total processing time for a query is 0.91 seconds². This time is divided into 74.79% for evidence search and 23.95% for reading. The total memory cost is 79.32GB, with the index occupying 81.95%, the raw corpus 16.39%, and the model 1.66%. Approaches to optimize this may include:

1. Reducing Processing Time: (1) Accelerating Evidence Search, (2) Accelerating Reading
2. Reducing Memory Cost: (1) Reducing Index Size, (2) Reducing Corpus Size, (3) Reducing Model Size
3. One-stage Frameworks: (1) Directly Generating Answers, (2) Directly Retrieving Answers

Techniques used in this context may include:

1. Data-based: (1) Passage Filtering, (2) Dimension Reduction, (3) Product Quantization
2. Model-based: (1) Model Pruning, (2) Knowledge Distillation, (3) Knowledge Source

²It's important to mention that DPR is a Representation-based Retriever, which allows offline storage of passage embeddings. The result was obtained using an Nvidia GeForce Rtx 2080 Ti GPU, averaged over 1000 examples

2. Background and Related Work

A common technique, which is used in nearly every experimental setup for QA-Systems, is FAISS [Johnson et al., 2021], a GPU-optimized implementation of the exact k -means clustering algorithm.

For a detailed overview of approaches towards more efficient ODQA systems, please refer to the comprehensive survey by Zhang et al. [Zhang et al., 2023b].

2.2. Conversational Question Answering

The differentiation of Conv QA towards QA will be discussed in Section 2.2.1. This section also introduces the fundamental concepts of Conv QA which are necessary to understand challenges and necessary components compared to a regular QA-System. Section 2.2.2 will cover approaches toward the concept of query expansion. Section 2.2.3 will clampse on the concept of initiative and further approaches towards a Conversation Manager.

Most crucial for the research of this thesis are the concepts introduced in Section 2.2.1. Especially how a Conv QA distinguishes from a regular QA-System and the core concepts of a conversation. The approaches introduced in Section 2.2.2 are good to know, but the thesis will later focus mostly on query rewriting. The other approaches are necessary for alternative system designs which will not be discussed. Section 2.2.3 only scratches the surface of the research field of *Initiative*. This is a highly complex field, with ongoing research and won't be covered in this thesis.

2.2.1. Basics

Core concepts in the field of Open-Domain Conversational Question Answering (ODCQA) towards a conversation in terms of Conv QA are: *Turns*, *Hisotry*, *Memory*, *Session* and *Dialog Features* and *Dialog State* [Zamani et al., 2023]. It's important to mention, that in other subdomains/-tasks of CIS more concepts are introduced, such as *State*, those are not necessary or applicable for ODCQA [Zaib et al., 2022].

Figure 2.7 shows the core concepts based on a chat. Firstly, a **Turn** is a question-response pair. Whereas a conversation usually consists of multiple turns (multi-turn). Conversational Question Answering (CoQA) is a dataset published in 2019 by researchers at Stanford in order to extend the known QA dataset SQuAD towards a conversational dataset, whereas on average one conversation session consists of 15 turns [Reddy et al., 2019]. Multi-turns are the main distinguisher between the in Section 2.1 introduced task, to a Conv QA task. In a multi-turn scenario natural language phenomena like *coreference* (multiple expressions referring to the same thing) or *ellipsis* (omitting words or

2. Background and Related Work

topics implied by the context) can occur. While in regular QA the System will only be challenged with single-turn scenarios, so only one question, which needs an answer, in Conv QA the systems have to face multi-turn scenarios, where a user might also ask followup question or in general multiple questions after each other. A **Hisotry** is consequently a set of turns that belong to one conversation session. A **Session** is a in it completed conversation. Lastly, the **Memory** is an abstract entity in which the Conv QA-System stores knowledge related to a history, session or even user in general [Zamani et al., 2023, Gao et al., 2023]. The extent of this depends on the implementation of memory in the Conv QA pipeline, which will be discussed in Section 2.2.2.

Dialog Features need to be assessed extra to the other mentioned concepts. While the other concepts tackle the frame of the conversation, the dialog feature evaluates the user questions themself. Possible dialog features may include *drilling-down* questions, *topic-shift*, *clarification* or *definition*. Different dialog features call for different responses by the system [Gupta et al., 2020]. The **Dialog State** has to be assessed similarly. The dialog state represents the relation between turns. In cases of pre-defined domains methods like state slots are used, e.g. **Date _**, **Location _**, **Artist _** have to be filled during the conversation in order to retrieve the correct information from the KB [Rastogi et al., 2020]. Open-Domain Conv QA usually don't track the state *explicitly*, but rather track it *implicitly* via the type of implementation of the *Contextual Query Understanding* unit.

Regarding the System architecture of a Conv QA there is no one fits them all solution at the moment, but Gao et al. [Gao et al., 2023] presented a modern system architecture, which represents commonly used approaches and their corresponding components in a general fashion. This general architecture can be observed in Figure 2.8. Modern Conv QA systems are closely related to QA systems, but lag certain generalizing components in order to be full CIS systems [Zamani et al., 2023].

Similar to the retriever-reader architecture introduced in Section 2.1.2, a Conv QA is made up of those two components as well. The retriever has to understand the context, so the history of a conversation and retrieve based on that the most relevant documents. The reader on the other hand is closely related to the reader of a classic retriever-reader architecture [Zamani et al., 2023, Gao et al., 2023]. Some implementations even feed into the reader component the context in order to rank the retrieved passages better and generate a more accurate answer [Owoicho et al., 2022].

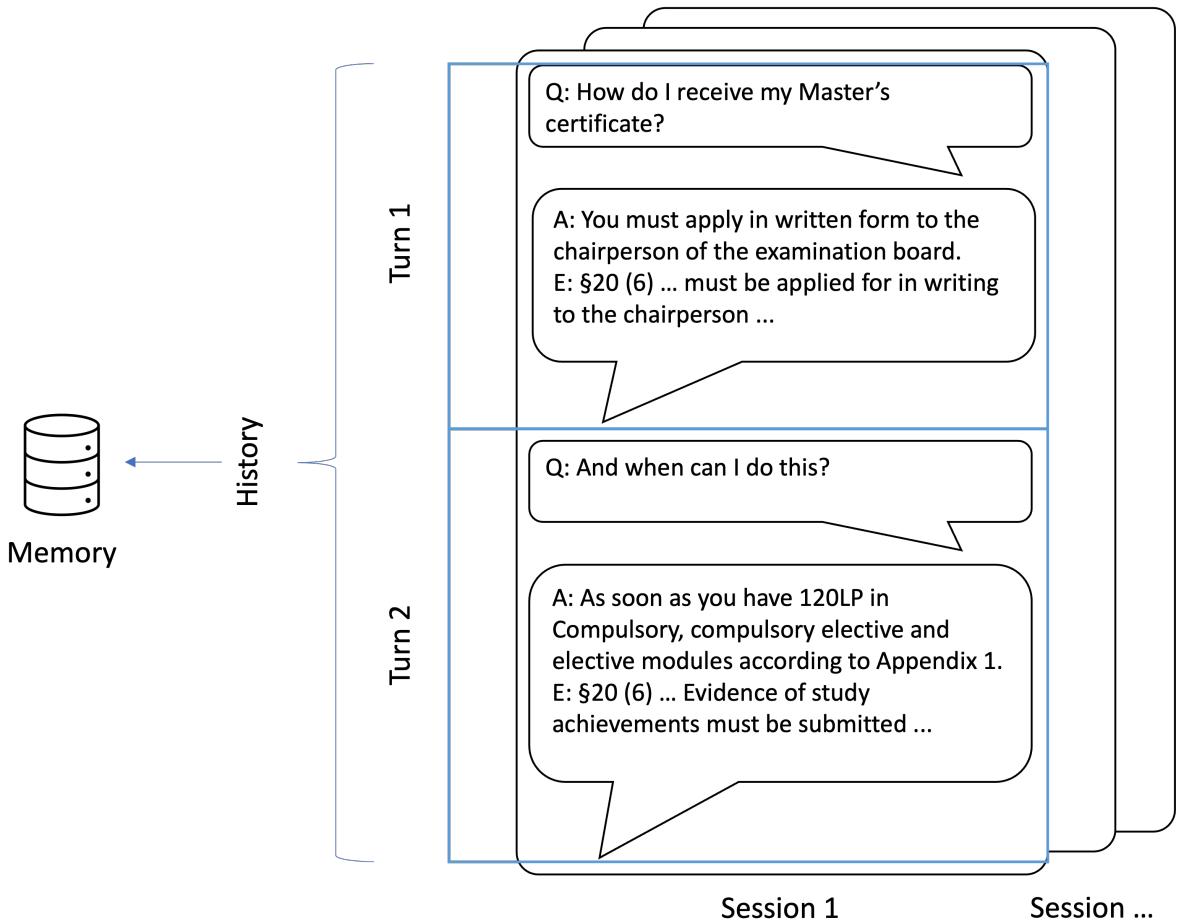


Figure 2.7.: Concepts of a Conversation in regards to a CIS

2.2.2. Contextual Query Understanding

There are two main distinguishing approaches towards history implementation. The first is a simple heuristic of using the last- k turns for **Query Expansion**, **Query Rewriting** or **Conversational Retrievers**. The second is to extract the important parts of the history in regards to a question and use them for Query Expansion or Rewriting [Gao et al., 2023].

A good example to explain the second approach towards extracting important parts of the history $H = (q_1, a_1), \dots, (q_i, a_i)$ given a new question q_{i+1} is **Query Resolution by Term Classification** (QuReTeC)[Voskarides et al., 2020]. QuReTeC consists of two components essentially: one BERT-based model and a trainable classification layer. The H is being passed through the BERT model, whereas the following structure of concatenation is being used:

$$\text{BERT}([CLS], H, [SEP], q_{i+1}) \quad (2.10)$$

2. Background and Related Work

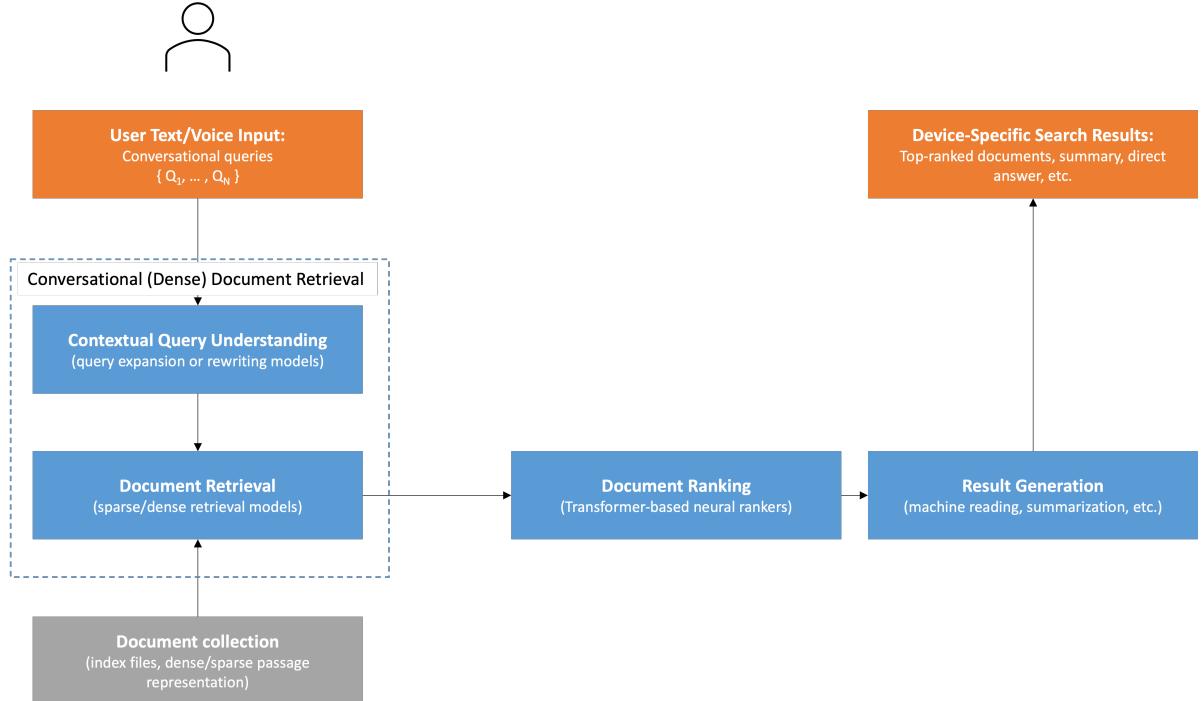


Figure 2.8.: General System Architecture of a Conv QA System
by Gao et al. [Gao et al., 2023]

On every first sub-token of a term of the H the term classification layer is applied, which is a network consisting of a dropout layer, a linear layer and a sigmoid function. The term classification layer predicts a label between 0 – 1 indicating its importance for answering the new question q_i . This leads to a set of terms I which need to be incorporated into the retrieval [Voskarides et al., 2020]. This is generally also known as **Query Expansion**, whereas we add terms to a given query for retrieval. Next to this supervised, trained approach, there are also implementations that work unsupervised like Historical Query Expansion (HQExp) [Yang et al., 2019], which was one of the best-performing models in the TREC CAsT 2019 [Dalton et al., 2020].

Modern neural approaches more often implement a **Query Rewriting** module which is built on top of seq-2-seq-models to rewrite a query q_{i+1} given a history H in order to use the generated new query for retrieval using an established QA retriever [Owoicho et al., 2022]. The main advantage of this approach is the absence of the need for large supervised datasets as for Conversational Retrievers [Dai et al., 2022a]. One of the top performing models in the TREC CAsT 2022 was HEATWAVE by a Team of the University of Cambridge England [Liusie et al., 2022]. HEATWAVE utilized a query rewriter and a classical lexical BM25 retriever in combination with a BERT-based reranker. The rewriter uses a T5-based Transformer model and gives as input $ctx - n - m$,

2. Background and Related Work

where n refers to the last n -many user utterances and m to the m -many system responses. In general, the task can be simply broken down into the following:

$$q_{\text{rewritten}} = \text{Rewriter}(ctx - n - m) \quad (2.11)$$

For training of this model, they used among others the canard dataset [Elgohary et al., 2019] a manually annotated version of the QuAC dataset, specifically for the task of query rewriting given a conversation history H .

State-of-the-art research utilizes more and more LLM for the task of Query Rewriting, as they can handle long context histories H and are in general strong zero- or few-shot models [Mao et al., 2023]. This is also the main approach frameworks like *Langchain* [Langchain, 2023] or *ChatGPT by OpenAI* [OpenAI, 2023] use.

Lastly the approach of **Conversational Retrievers** exists. Those use compared to classical QA retrievers not a pair of (q, p) in order to calculate a similarity $\text{sim}(q, p)$ between the question q and the passage p , but use conversational interactions like $(q_1, a_1, \dots, q_i, a_i, q_{i+1}, p)$, in short (H, q_{i+1}, p) , so combining a conversation history H with a new question q_{i+1} and the relevant passage p to answer this question given the history H [Gao et al., 2023, Dai et al., 2022a]. High-performing zero-shot or subdomain-adapted Conversational Retrievers do not exist, as it is extremely time-consuming to create a dataset for this type of Retriever. To close this gap, researchers proposed sufficient data augmentation techniques to generate those datasets, given a document. One example is the work of Dai et al. [Dai et al., 2022a] which introduced the technique of “Dialog Inpainting” [Dai et al., 2022a].

2.2.3. Initiative

Most modern human-computer interactions follow a one-sided initiative model, where either the user- or the system-initiative is given. In mixed-initiative scenarios of ConvQA the system can take initiative without explicit commands of the user. Examples of initiative are: *Topic Shifts*, *Clarification Questions* or *Question Recommendations* [Zamani et al., 2023]. In this thesis, we will focus on *Clarification Questions* only.

Asking *Clarification Questions* is the task of identifying ambiguity in a user’s search request and resolving it by posing a question with the intent to eliminate ambiguity. A common taxonomy to use for the types of ambiguous questions includes: 1) questions where the focus is ambiguous, and 2) questions with several distinct possible answers [Larsson, 2002]. Several studies have been conducted to understand user behavior in relation to clarification in search. Tavakoli et al. [Tavakoli et al., 2021] found that users are more likely to engage when a *Clarification Question* aims at clarifying ambiguous

2. Background and Related Work

information instead of seeking confirmation or similar. Zamani et al. [Zamani et al., 2020], based on search-engine log analysis, found that users are more likely to engage with a *Clarification Question* if their own question has high ambiguity and, therefore, multiple possible resolutions or when there is a dominant assumed search intent.

To resolve ambiguity in the context of Conv QA, the most modern solutions follow a two-step approach, where, in the first step, the system identifies the ambiguity, and in the second step, the system generates a *Clarification Question* to resolve the ambiguity [Kuhn et al., 2023, Guo et al., 2021]. The exact implementation may differ, whereas Guo et al. [Guo et al., 2021] developed a seq-2-seq-model to predict the ambiguity of a question given context, Kuhn et al. [Kuhn et al., 2023] used a chain-of-thought-like approach, where a LLM performs both steps sequentially. In general, it has to be said that, in the context of ODQA, the niche of *Clarification Questions* is not well-researched yet, and there is no established benchmark dataset.

2.3. Efficient Large Language Models

With the increasing size of large language models (LLM), **Large Language Models with Adapters (Llama) 2** offers models ranging from 7 billion to 70 billion parameters [Touvron et al., 2023]. Even these models are considered relatively small compared to the largest models like PaLM 2 [Anil et al., 2023] with 340 billion parameters. The challenge arises when running such models in scenarios with limited computational resources, especially on smaller domains or tasks. This challenge is particularly relevant to the task presented in this thesis, which involves building a Conv QA system for a custom set of documents.

While several surveys [Ling et al., 2023, Zhao et al., 2023] have explored the topic of efficient LLM usage in resource-constrained systems, Treviso et al. [Treviso et al., 2023] present the most comprehensive taxonomy of methods and approaches in this context. Figure 2.9 provides a high-level overview of the stages at which efficiency-improving methods can be implemented in LLMs. Given the specific focus of this thesis, not all stages will be discussed in detail. For more comprehensive insights, please refer to the original survey by Treviso et al. [Treviso et al., 2023].

Section 2.3.1 will explore possibilities to enhance efficiency during the fine-tuning process, while Section 2.3.2 will delve into the topic of model compression, which is applicable to the *Inference* step in Figure 2.9.

This thesis itself won't delve into fine-tuning (the content of Section 2.3.1) or model compression (the content of Section 2.3.2) itself. Still, concepts like multi-task learning are necessary to understand the capabilities and where those are coming from for modern

2. Background and Related Work



Figure 2.9.: Adapted Stages of Efficiency Improvement for LLMs
by Treviso et al. [Treviso et al., 2023]

LLMs, and the technique of *Prompting* will be used in Chapter 4. Also, in Chapter 4, quantized models will be used to evaluate the performance of the Conv QA-System, out of the necessity to run LLMs on the available hardware resources. Nevertheless, the following section can be seen as an excursion to the main topic of this thesis for real-world implementations facing resource-constrained systems.

2.3.1. Fine-Tuning

Hu et al. [Hu et al., 2022] demonstrated the significant benefits of fine-tuning GPT-3 for few-shot applications, highlighting the remarkable improvements fine-tuning can achieve. This is further supported by the experiments conducted by Chung et al. [Shen et al., 2023].

Efficient fine-tuning of LLMs can be categorized into three distinct approaches: *Parameter Efficiency*, *Multi-task Learning*, and *Prompting*. Figure 2.10 provides an overview of these approaches along with their corresponding methods.

Parameter Efficiency is commonly referred to as **Parameter-efficient Fine-tuning (PEFT)** [Huggingface, 2023a]. A notable PEFT approach is Low-Rank Adaptation (LoRa), developed by Hu et al. [Hu et al., 2022]. LoRa falls under the category of adapters, a term coined because it revolves around the concept of freezing the parameters of the LLM and fine-tuning only a small set of task-specific parameters, which can be swapped depending on the desired downstream task. Unlike some other adapter-based methods [Houlsby et al., 2019], LoRa does not introduce additional inference latency due to the merging of trainable matrices with frozen weights. Moreover, LoRa can be seamlessly combined with many other PEFT methods.

In practice, given a pre-trained weight matrix $W_0 \in \mathbb{R}^{d \times k}$, which is typically full-rank between layers, the update can be constrained to be a low-rank composition: $W_0 + \Delta W = W_0 + BA$, where $B \in \mathbb{R}^{d \times r}$, $A \in \mathbb{R}^{r \times k}$, and the rank $r \ll \min(d, k)$. While W_0 remains frozen during training, A and B become trainable parameters. The forward pass $h = W_0x$ can be represented as the following sum:

2. Background and Related Work

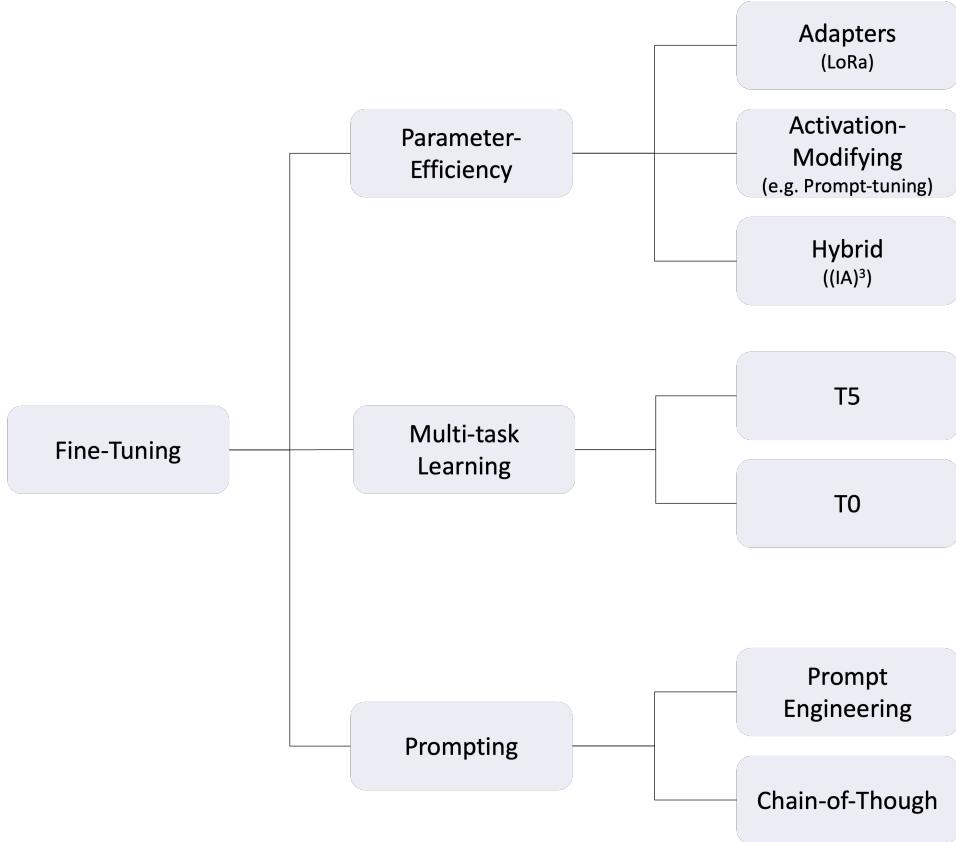


Figure 2.10.: Adapted Fine-Tuning Approaches for LLM by Treviso et al. [Treviso et al., 2023]

$$h = W_0x + \Delta Wx = W_0x + BAx \quad (2.12)$$

Figure 2.11 illustrates the architecture and initialization during training. The parameters of A are randomly sampled using Gaussian initialization, while B is initialized to 0.

Other PEFT techniques include **prompt-tuning** [Lester et al., 2021] and **prefix-tuning** [Li and Liang, 2021]. Both approaches are similar in the way they leverage task-specific modifications to the input to guide the model’s behavior. They involve concatenating learned vectors to activations or embedding sequences, making them activation-modifying PEFT methods.

A PEFT approach that can be considered a hybrid between LoRa and activation-modifying techniques is **(IA)³** [Liu et al., 2022]. What sets **(IA)³** apart is its focus on LLMs designed explicitly for multi-task learning, as all existing PEFT techniques significantly underperformed in experiments conducted by Liu et al. [Liu et al., 2022]. In **(IA)³**, the model’s activations are rescaled using element-wise multiplication with

2. Background and Related Work

learned vectors, known as adaptors. Specifically, (IA)³ employs three learned vectors: $l_k \in \mathbb{R}^{d_k}$ for keys and $l_v \in \mathbb{R}^{d_v}$ for values in self-attention and encoder-decoder attention mechanisms, as well as $l_{ff} \in \mathbb{R}^{d_{ff}}$ for the feed-forward network. The rescaling is incorporated into the attention mechanism as follows:

$$\text{softmax} \left(\frac{Q(l_k \odot K^T)}{\sqrt{d_k}} \right) (l_v \odot V) \quad (2.13)$$

For the feed-forward network, the rescaling is implemented as follows, where γ represents the feed-forward activation:

$$(l_{ff} \odot \gamma(W_1 x)) W_2 \quad (2.14)$$

In summary, (IA)³ is a PEFT approach specifically designed for multi-task learning, and it appears to outperform LoRa in terms of the number of parameters added and training computation costs.

Multi-task Learning refers to the concept of fine-tuning a PrLM on various tasks to achieve better zero-shot and fine-tuning performance. One of the most prominent PrLMs trained on multiple NLP tasks is T5 [Raffel et al., 2020]. Liu et al. [Liu et al., 2022] demonstrated that a generically trained T5 model (T0) can be few-shot fine-tuned with approximately 10% of the parameters of LoRa, at a lower computational cost, while achieving higher accuracy in classification tasks³. This is made possible, due to using (IA)³ as PEFT and the multi-task pre-training of the used model. Still, this also includes the drawback of having in general a larger model.

Prompting refers to the general concept of presenting a task as a textual instruction to a LLM [Brown et al., 2020]. Recent advances have even led to the development of a new sub-task and job role called *Prompt Engineering* [White et al., 2023]. It's worth noting that different prompts with the same intent can yield different results, making the selection of the right prompt a challenge in itself [Liu et al., 2023a]. Furthermore, concepts like chain-of-thought (CoT) prompts have been developed. In CoT prompts, the given example in a few-shot prompt is redesigned to mimic step-by-step reasoning and conclusions known from the way humans think, aiming to achieve higher performance in zero- and few-shot scenarios simply by adjusting the explicit natural language prompt [Wei et al., 2022]. CoT requires retraining of the LLM.

³Currently, there is no experiment comparing LoRa and (IA)³ on QA tasks.

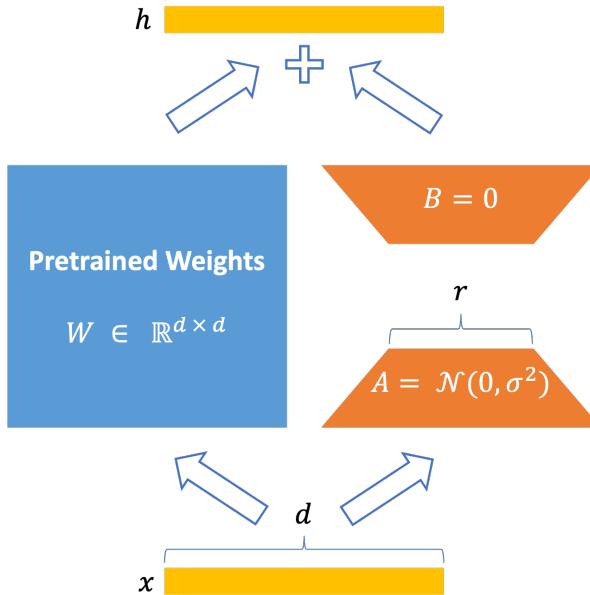


Figure 2.11.: General Idea of LoRa by Hu et al. [Hu et al., 2022]

2.3.2. Compression

Compression aims to reduce the size of a model, whether it's the number of parameters while maintaining the same level of accuracy on the downstream task, or the actual storage required for the model. There are three primary approaches to this task: *Pruning*, *Knowledge Distillation*, and *Quantization* [Treviso et al., 2023, Zhu et al., 2023]. Figure 2.12 provides an overview of these approaches and their corresponding methods.

Pruning can be further categorized into *structured* and *unstructured* pruning. Structured pruning involves removing specific patterns of weights or activations from a model, with the goal of maintaining a dense matrix representation to ensure compatibility with existing implementations and hardware. One notable example of a structured pruner for LLMs is LLM-Pruner [Ma et al., 2023]. On the other hand, unstructured pruning entails removing individual weights or activations from a model, resulting in a sparse matrix representation. This approach may require specialized hardware or software implementations to efficiently compute and achieve speed improvements of $1.5\times$ to $2.16\times$, while reducing up to 60% of the parameters [Frantar and Alistarh, 2023]. Examples of engines designed specifically for unstructured pruning include NVIDIA's CUTLASS library for GPUs [Frantar and Alistarh, 2023] and DeepSparse [Magic, 2023] for CPUs.

Knowledge Distillation is an approach that involves using a generally well-performing LLM as a teacher to instruct a significantly smaller student model [Hinton et al., 2015]. Zhu et al. distinguish between *White-box* and *Black-box* knowledge distillation. In the former, the student has full access to the teacher's parameters, while in the latter, only

2. Background and Related Work

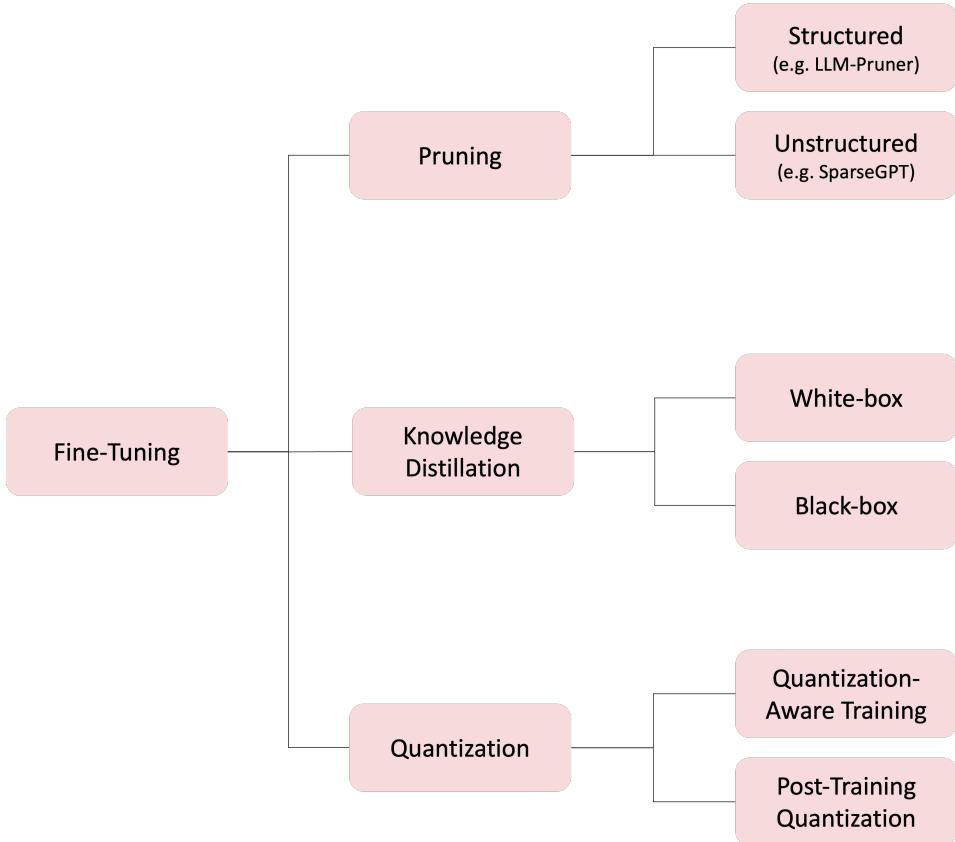


Figure 2.12.: Adapted Compression Approaches for LLM by Treviso et al. [Treviso et al., 2023]

the teacher’s predictions are accessible to the student [Zhu et al., 2023]. An example of white-box knowledge distillation is MiniLLM [Gu et al., 2023], where the distribution of the final layer’s outputs for both the teacher and the student, given a prompt, is compared using the Kullback-Leibler divergence. This comparison is used in a loss function for backpropagation in the student model.

Black-box approaches are more commonly used in knowledge distillation. In these cases, a LLM is employed to either directly provide its predictions based on a prompt [Huang et al., 2022], offer assisting explanations [Li et al., 2022b], or sort the training data by difficulty and artificially generate more data points [Jiang et al., 2023], among other techniques. For a comprehensive overview, please refer to Section 2.2 *Knowledge Distillation* in Zhu et al.’s survey [Zhu et al., 2023]. Experiments with different distillation approaches have shown that distillation has its limitations, and for specific downstream tasks, fine-tuning can outperform knowledge distillation [Zhu et al., 2022].

Quantization is an approach that involves reducing the datatype representation of weights or activations, which are typically floating-point numbers, to smaller representations in terms of bits, such as 8-bit integers or even smaller discrete formats [Gholami

2. Background and Related Work

et al., 2022]. Generally, there is a distinction between *Quantization-aware Training* and *Post-Training Quantization*. The names are self-explanatory; the former involves applying and adjusting quantization during the training process (either pre-training or fine-tuning) [Liu et al., 2023b], while the latter pertains to quantization after the training is completed [Frantar et al., 2023]. In both cases, numerous approaches and methods exist, applying different paradigms and quantization techniques, including decisions regarding which parameters to quantize, structured vs. unstructured quantization, quantization strength, and many others. It’s not possible and necessary to discuss all of these here, but for a comprehensive overview, please refer to Section 2.3 *Quantization* in Zhu et al.’s survey [Zhu et al., 2023].

The most prominent example of Post-Training Quantization is GPTQ, which is the only ready-to-use implementation available in the Huggingface Transformer Library [Huggingface, 2023b]. GPTQ was the first method to achieve high compression for LLMs with over 175 billion parameters while maintaining high accuracy compared to prior state-of-the-art algorithms. Specifically, with a 4-bit quantization of the weights, GPTQ achieved approximately $5\times$ compression for BLOOM-176B and OPT-175B, two openly available LLMs, while experiencing only a ≤ 0.25 decrease in perplexity compared to the original model. Therefore, the following section will explain the Post-Training Quantization approach of GPTQ in detail.

GPTQ builds upon Optimal Brain Quantization (OBQ), the previous work by Frantar et al. on Post-Training Quantization [Frantar and Alistarh, 2022]. With GPTQ, their objective was to reduce the runtime complexity of OBQ, which is $O(d_{row} * d_{col}^3)$, making it compatible with LLMs containing billions of parameters. The central idea behind GPTQ is a layer-wise optimization approach. The aim is to discover quantized weights $\widehat{\mathbf{W}}$ that minimize the squared error compared to the full-precision layer \mathbf{W} output, using a given set of input data points \mathbf{X} :

$$\operatorname{argmin}_{\widehat{\mathbf{W}}} \|\mathbf{W}\mathbf{X} - \widehat{\mathbf{W}}\mathbf{X}\|_2^2 \quad (2.15)$$

In OBQ, we denote the next weight to be quantized as w_q . We define the function $\text{quant}(w)$, which rounds a weight w to the nearest value on the quantization grid.

$$w_q = \operatorname{argmin}_{w_q} \frac{(\text{quant}(w_q) - w_q)^2}{[\mathbf{H}_F^{-1}]_{qq}} \quad (2.16)$$

In the context of **GPTQ**, a column of weights is always updated simultaneously. Therefore, $\text{quant}(W_{:,j})$ refers to the following:

$$\text{quant}(W_{:,j}) := \forall w_q \in W_{:,j} \quad (2.17)$$

The Hessian matrix $\mathbf{H}_F = 2X_F X_F^T$ is utilized for both weight updates and quantization error calculations. Once all columns within a block B are quantized, the weight update is computed as follows, where Q represents the set of indices corresponding to quantized weights:

$$\boldsymbol{\delta}_F = -(\mathbf{w}_Q - \text{quant}(\mathbf{w}_Q)) \left([\mathbf{H}_F^{-1}]_{QQ} \right)^{-1} (\mathbf{H}_F^{-1})_{:,Q} \quad (2.18)$$

Furthermore, the Hessian is updated in the following manner, avoiding the need for recomputation; instead, columns corresponding to quantized weights are simply dropped from the Hessian.

$$\mathbf{H}_{-Q}^{-1} = \left(\mathbf{H}^{-1} - \mathbf{H}_{:,Q}^{-1} \left[(\mathbf{H}^{-1})_{QQ} \right]^{-1} \mathbf{H}_{Q,:}^{-1} \right)_{-Q}. \quad (2.19)$$

This leads to Algorithm 1.

Algorithm 1 Quantize W given inverse Hessian $H^{-1} = (2XX^T + \lambda I)^{-1}$ and blocksize B by Frantar et al. [Frantar et al., 2023]

```

1:  $Q \leftarrow \mathbf{0}_{d_{\text{row}} \times d_{\text{col}}}$                                 ▷ Quantized output
2:  $E \leftarrow \mathbf{0}_{d_{\text{row}} \times B}$                                 ▷ Block quantization errors
3:  $H^{-1} \leftarrow \text{Cholesky}(H^{-1})^T$                                 ▷ Hessian inverse information
4: for  $i \leftarrow 0, B, 2B, \dots$  do
5:   for  $j \leftarrow i, \dots, i+B-1$  do
6:      $Q_{:,j} \leftarrow \text{quant}(W_{:,j})$                                 ▷ Quantize column
7:      $E_{:,j-i} \leftarrow \frac{W_{:,j}-Q_{:,j}}{[H^{-1}]_{j,j}}$                                 ▷ Quantization error
8:      $W_{:,j:(i+B)} \leftarrow W_{:,j:(i+B)} - E_{:,j-i} \cdot H_{j,j:(i+B)}^{-1}$       ▷ Update weights in block
9:   end for
10:   $W_{:,i:(i+B)} \leftarrow W_{:,i:(i+B)} - E \cdot H_{i:(i+B),i:(i+B)}^{-1}$       ▷ Update all remaining weights
11: end for

```

To enable GPTQ to be applicable to LLMs with billions of parameters, the authors have introduced three key optimizations:

1. *Arbitrary Order*: In the case of large models, the order in which weights are quantized becomes irrelevant. Therefore, GPTQ updates all weights in the same order for all rows. This means that the set of unquantized weights, denoted as F , and H_F^{-1} , the Cholesky Form - Inverse Layer Hessian, remain constant across all rows. This is because H_F depends solely on X_F and is independent of the weights. This reduction in the number of times H needs to be updated simplifies the process from $d_{\text{col}} \times d_{\text{row}}$ updates to just d_{col} updates.

2. Background and Related Work

2. *Lazy Batch-Updates*: Quantization of a column depends solely on updates to that particular column. Therefore, GPTQ employs batches of columns (with a batch size of $B = 128$). Equations 2.18 and 2.19 can be executed after the computation of a full batch B . The set of indices Q corresponds to the indices of quantized weights in the batch.
3. *Cholesky Reformulation*: To address numerical errors that arise from repeated application of equation 2.19, a Cholesky reformulation is applied to calculate all the necessary information about H^{-1} in advance. As the complete Cholesky decomposition cannot be applied, a mild damping factor is applied to the diagonal.

Additionally, an accessible quantization package called **AutoGPTQ** has been developed, which implements the GPTQ algorithm in PyTorch [AutoGPTQ, 2023]. This package has been adopted by Hugging Face and is currently the only ready-to-use quantization technique available in the Transformers library [Huggingface, 2023b].

2.4. Related Work

2.4.1. Question Answering based on PDFs

PDF Question Answering is the task of providing answers to questions related to the content of one or multiple documents [Mathew et al., 2021]. The field of research which actively explores this the closest is Visual Document Question Answering. It works on the development of an IR-QA system that operates on images of documents. An exemplary architecture and a general pipeline for transforming PDFs into an IR-QA system is presented by McDonald et al. [McDonald et al., 2022]. They developed their zero-shot framework around the QASPER dataset but used the original PDFs instead of extracted text via LaTeX. Moreover, readily available open-source tools like V-Doc [Ding et al., 2022] simplify the deployment and testing of datasets, models, and IR-QA systems of the Visual Document Question Answering domain.

More recently, the open-source framework *Langchain* has gained tremendous attention⁴. Langchain focuses on harnessing LLMs using chains, which are essentially prompts for an LLM that can be chained together [Langchain, 2023]. They also provide documentation on building a QA system based on PDFs [Langchain, 2023]. Similarly, *OpenAI* offers a Retrieval Plugin for *ChatGPT* [OpenAI, 2023], also an open-source repository. These QA systems adhere to the paradigms established in previous works such as

⁴As of September 24, 2023, Langchain has received 63k stars on GitHub

2. Background and Related Work

[Karpukhin et al., 2020, Ni et al., 2022, Neelakantan et al., 2022, Lewis et al., 2020b]. Specifically, this entails:

- Given a text corpus, documents can be retrieved by extracting relevant passages. Data cleaning of the corpus is optional but not necessary. Therefore, these systems employ a *direct extraction* approach, especially when dealing with PDFs.
- Utilizing large-scale, diversely trained encoders. Representation-based Retrievers, when equipped with sufficient trainable parameters and diverse training datasets, often yield comparable results to fine-tuned, more complex retrieval models [Ni et al., 2022, Neelakantan et al., 2022].
- Using the LLM as a generative reader for QA, as demonstrated in the work of Izacard et al. [Izacard and Grave, 2021].

Non-LLM research for QA based on PDFs is notably scarce. In the field of ODQA, discussions regarding applicable frameworks that encompass the entire pipeline from PDFs to QA are infrequent. Instead, the focus often revolves around constructing QA systems using predefined and well-supervised datasets. However, there is some research that explores the feasibility of deploying high-performing QA systems in out-of-domain scenarios, bypassing the initial stage of data preprocessing (from PDFs to passages). This research strives to outline possibilities for using a QA system in real-world passage collections.

Applying Dense Retrievers Out-of-Domain: As emphasized by Thakur et al. in their “Heterogeneous Benchmark for zero-shot Evaluation of Information Retrieval Models” (BEIR) [Thakur et al., 2021], dense retrievers exhibit weak out-of-domain performance. Lyu et al. [Farea et al., 2022] also demonstrate the limited generality of dense retrievers when trained in one subdomain and subsequently applied in a different one. This underscores the conclusion that there are two approaches to employing retrievers in out-of-domain scenarios: (1) fine-tuning or (2) zero-shot, but with large encoders that have been trained on diverse datasets [Ni et al., 2022].

The challenge with fine-tuning lies in the unavailability of labeled data, which is typically required for supervised models in the form of tuples such as (*question*, *answer*, *context*). Several diverse approaches have been developed to address this issue. One approach employs QG techniques, as exemplified by PROMPTAGATOR [Dai et al., 2022b], which utilizes LLMs. Another strategy involves the use of Mixture-of-Experts and meta-learning algorithms [Chen and Ma, 2021]. Some researchers have explored semi-supervised training datasets, as demonstrated by Sachan et al. [Sachan et al., 2023], who developed ART, a training framework for dense retrievers that only requires

2. Background and Related Work

questions and surpasses the standard DPR training implementation. At the current point in time, there is no state-of-the-art approach to fine-tune a dense retriever on a small subdomain dataset.

In their study, Reddy et al. [Gangi Reddy et al., 2021] addressed the challenge of creating a QA-System for Covid-19-related documents, where no supervised QA dataset was available. Consequently, they conducted a comparison between the performance of zero-shot BM25 and DPR. Their findings revealed that BM25 outperformed DPR on the BiosQA QA dataset, closely related to the Covid-19 domain. Throughout their experiments, they evaluated various approaches, including simple zero-shot techniques, and fine-tuning of DPR using QG via BART, which yielded superior results. Notably, the most effective retriever for unsupervised domain adaptation was a combination of BM25 and unsupervised fine-tuned DPR.

Furthermore, Gururangan et al. [Gururangan et al., 2020] demonstrated in their experiments that fine-tuning PrLMs on domain-specific language or, even better, task-specific data led to a significant performance boost.

Gholami et al. [Gholami and Noori, 2021] experimented with non-fine-tuned dense retrievers on a non-QA dataset, specifically a collection of AWS documentations. Their results, particularly for the retrieval component, were sobering, aligning with the findings of benchmark studies by Thakur et al. [Thakur et al., 2021] and Lyu et al. [Farea et al., 2022].

On the other hand, there exist reader components with a high degree of generalizability, as demonstrated by UnifiedQA-v2 [Khashabi et al., 2022], an extractive reader and T5 [Raffel et al., 2020], a generative reader. So the main challenge, when building a IR-QA-System, lies within the implementation and adaptation of the retriever component.

2.4.2. Open-domain Conversational Question Answering

Datasets: Notable datasets for Conv QA include CoQA, TREC 2019, and QReCC, which primarily feature extractive questions [Reddy et al., 2019, Dalton et al., 2020, Dai et al., 2022a]. While these datasets address conversation-specific challenges like coreference resolution, the absence of an adequate benchmark for Conv QA is evident. Naveed et al.’s survey on LLMs [Naveed et al., 2023] highlights various alternative datasets in this domain, largely focused on specific challenges. Recent works like Llama2 underscore the lack of a state-of-the-art benchmark, necessitating human-based evaluation as a primary approach [Touvron et al., 2023].

2. Background and Related Work

Chat Fine-Tuned LLMs: Fine-tuning LLMs on human chat-like tasks, introduced by LaMDA and later widely spread by ChatGPT, has emerged as a promising approach for conversational information retrieval [Thoppilan et al., 2022, OpenAI, 2023]. These models are designed for conversational interactions and can generate human-like responses, although they may lack truthfulness, a crucial aspect for high-quality Conv QA systems.

Truthfulness: Addressing the issue of truthfulness, several approaches such as REALM, RAG, FiD, and WebGPT leverage external knowledge [Guu et al., 2020, Lewis et al., 2020b, Izacard and Grave, 2021, Nakano et al., 2022]. Among these, RAG has garnered significant research attention due to its ease of implementation and explicit knowledge, enhancing model understandability [Gao et al., 2024].

Use-Case Implementation: Research on implementing a real-world Conv QA system using a specific dataset is limited. The DialDoc 2021 Shared Task presents various solutions, mainly focusing on extractive readers [Feng, 2021]. However, there is a lack of research addressing the challenges, approaches, and considerations for designing such a system based on the latest advances in LLMs. The documentation of the Langchain Framework for Question Answering provides some guidens, although it lacks certain considerations [Langchain, 2023], such as complex retrievers or other paradigms then RAG for example FiD. Langchain mainly focuses on DPR retrievers and explicit RAG implementations.

3. Open-domain QA Chatbot

This chapter outlines the methods and techniques employed in the development of a conversational question-answering system designed for the use with a collection of documents as knowledge source. The chapter is structured as follows: Section 3.1 provides an overview of the desired use case, its objectives, and constraints concerning a Conversational Question Answering System. Section 3.2 will introduce the fundamental document inquiry model, and Section 3.3 presents a general framework that can be utilized for the implementation of a Conversational Question Answering System. Its subsections will highlight and discuss the components introduced within the framework. Section 3.4 will summarize the contribution of this thesis to the field of Conversational Question Answering Systems.

3.1. Overview and Objective

The primary use case addressed in this thesis can be summarized as follows: Imagine having a collection of documents, and our goal is to create a chatbot capable of engaging in conversations about the knowledge within these documents. This chatbot provides accurate answers to questions based on the content of the documents and furnishes supporting evidence from these documents. Furthermore, it enables users to have a conversational query experience, allowing them to ask follow-up questions and engage in dialogue with the chatbot based on its previous responses. Figure 3.1 illustrates an example of this use case.

Currently, to the best of our knowledge, there is no scientific paper or similar resource offering a comprehensive framework or pipeline to address this use case. This thesis aims to bridge this gap by presenting a framework and pipeline designed to tackle this specific scenario. Figure 3.2 provides an overview of the system architecture. The system follows the RAG architecture, as detailed in Section 2.1.4, which extends the classical Retriever-Reader with a LLM as a Reader, capable of incorporating parametric knowledge. To extend RAG to a Conv QA, a Contextual Query Understanding (CQU) unit, as introduced in Section 2.2.2, is essential. This novel architecture will be termed **Conversational Retrieval-Augmented Generation (ConRAG)**. The extraction

3. Open-domain QA Chatbot

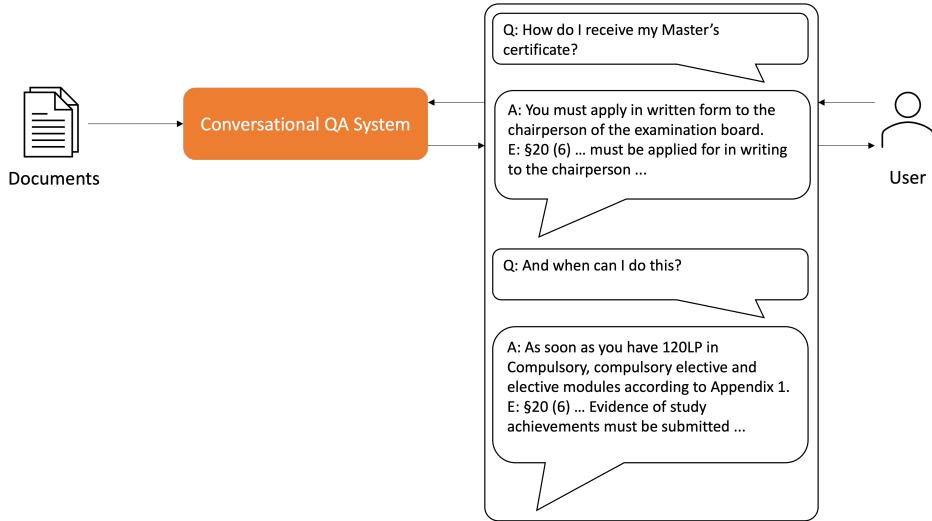


Figure 3.1.: Overview of the Example Use-Case

pipeline will be discussed in Section 3.3.1, with its primary tasks being the extraction of passages from the provided set of documents, the creation of an index, and the optional generation of synthetic training data. The three major modules comprising the architecture, namely the *Retriever*, *Reader*, and *CQU*, will be elaborated in their respective sections: 3.3.3, 3.3.4, and 3.3.2.

To summarize, the objectives of the QA capabilities of the system are as follows:

1. Utilize **documents** as the primary **knowledge source**.
2. Enable the QA-System to handle a **variety of answer types**, including: **extractive**, **abstractive**, and **boolean**.
3. **Provide references** to document snippets as **evidence to answers**.
4. Ensure the pipeline's generalizability by using **open-domain methods**, allowing it to adapt to new domains or knowledge sources with **minimal or no supervision** and **small datasets**.
5. Design the pipeline to be **feasible without the need for datacenter-grade hardware resources**, making it accessible for development on standard research hardware.
6. **Prioritize accuracy as the primary objective**, as constraining memory consumption is indirectly covered in point (5). **Latency is not a primary concern**, as the system is not intended for real-time use and will not be optimized for that in this thesis.

3. Open-domain QA Chatbot

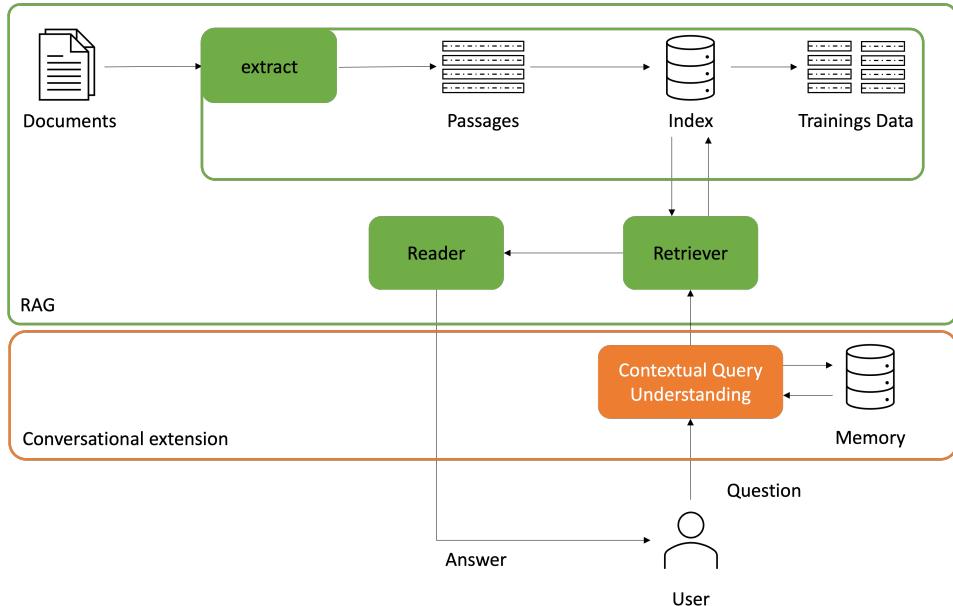


Figure 3.2.: Overview of the System Architecture

Regarding the ConvQA-System, the objectives are as follows:

1. Enable the ConvQA-System to **handle** the following follow-up **question types**: **drilling-down**, **clarification**, **topic shift** and **comparison**.
2. Be able to take Initiative in the form of **clarifying questions**.
3. The **memory** will be **limited to a session** and not multi-session.

3.2. Document Enquiry Model

This Section will lay out the problem of document-based Conv QA.

Substantially the fundamental source of knowledge is *documents*. A *document* can be any type of structured or unstructured file, which is being used for storing and displaying information. Examples are HTML (structured) or PDF (unstructured) files. A *document* consists of content C_d , a collection of information units c_d , whereas the content C_d has to have at least one c_d , but can contain also multiple. On top a *document* contains metadata M_d , which are a collection of key-value-pairs, an example is $M_d = \{(title, Examination Regulation Master Data and Computer Science), (publish_date, 14.06.2022)\}$. Next to the mentioned components of a *document*, it has also a unique identifier UID_d , therefore a *document* $d = (C_d, M_d, UID_d)$. For this thesis we will only consider textual information units in C_d and not figures or images. Out of the necessity, for knowledge granularity and precise information context, we will define

3. Open-domain QA Chatbot

passages next. A *passage* p is a subsequence of a string of textual content $c_d \in C_d$. The granularity of p can be defined as use-case specific. If p is a sentence, 100 tokens or the full textual content c_d . Nevertheless, every p contains a reference to the original document d it was taken from and has its own unique identifier UID_p . This leads to the following *Passage Model*:

Definition 3.1 (Passage Model) *A passage p is a subsequence of a textual content string $c_d \in C_d$ of a document $d = (C_d, M_d, UID_d)$, whereas $p = (\text{content}, UID_p, UID_d)$.*

Definition 3.1 indicates that the hierarchical or sequential order between passages inside a document won't be captured. This is true for this thesis work only, other approaches may use an incremental index instead of UID_p to indicate the order of passages inside a document d .

For ease of notation, we will refer to the *content* of a passage p as p itself in the following. The collection of all *passages* P will be referred to as the *Knowledge Source*.

For the following definitions, it's important to clarify the concept of *Intent* first. A simple example to illustrate intent is the following:

Question: When was Barack Obama born?

Answer 1: 4. August 1961

Answer 2: Barack Obama was the 44th president of the USA.

Answer 3: Either 04.08.1961 or 05.09.1962 I'm not sure.

The intent of the question is fulfilled given answer 1, so the answer contains the information the user was looking for when starting the search, but answer 2 misses the search intent. Answer 3 is somewhere in between, as it understands the search intent, but doesn't fulfill it correctly. Therefore intent can be summarized as the user's true information need he had when generating a search query. Mathematically we define intent in the following way:

Definition 3.2 (Intent) *Given two elements, which either can be questions, answers or passages, there exists an operation Intent $\mathcal{I}(x, y)$. This operation returns a value between 0 and 1, which indicates the overlap of the two intents of x and y , $\mathcal{I}(x, y) = [0, 1]$.*

Next, we need to define what a *Question* is. For this problem, a question is fundamentally a string. Generally, a question also has an intent, which can be measured against the golden answer (the correct answer to a question). The *Question Model* is therefore defined as follows:

Definition 3.3 (Question Model) A question q is a string. Given the golden answer a_q , $\mathcal{I}(q, a_q) = 1$.

Naturally, where there is a *Question*, there has to be an *Answer*. An *Answer* a is a string, which is the answer to a *Question* q . It can be considered a gold answer when $\mathcal{I}(q, a) = 1$. Formally, we define an *Answer* as:

Definition 3.4 (Answer Model) An answer a is a string, that answers a given question q . To which extent answer a answers question q is measured by the Intent. If $\mathcal{I}(q, a) = 1$, a is considered a gold answer to q .

In terms of conversations, we split an exchange between two agents into *Turns* as described in Section 2.2.1. Generally speaking, a *Turn* h consists of a tuple $\langle q, a \rangle$, whereas the a is the response to q . *Turns* happen in order and therefore have a logical relation. We refer to the sequence of multiple *Turns* within one conversation as *History* H .

Definition 3.5 (History Model) A history H is a sequence of turns h , whereas $h = \langle q, a \rangle$, $H = \langle h_1, h_2, \dots, h_i \rangle$.

As we now have elaborated, what *Questions*, *Knowledge Source*, *Answers* and *History* are, we're ready to define the problem of Conv QA:

Definition 3.6 (Conversational Question Answering Task) Given a new question q_{i+1} and a history $H = h_1, h_2, \dots, h_i$, a model (\mathbf{M}) generates an answer a_{i+1} , based on the provided knowledge in the knowledge source P , which satisfies the search intent of q_{i+1} . Next to the answer a_{i+1} , \mathbf{M} returns p as evidence from P . Formally:

$$\mathbf{M} : (q_{i+1}, H, P) \rightarrow (a_{i+1}, p)$$

3.3. Conversational Retrieval-Augmented Generation

In order to provide a solution to the Task of Conv QA as defined in Definition 3.6, the system must be able to perform evidence selection based on a *Knowledge Source* which is an important criterion also laid out in Section 3.1.

In order to now create a system architecture, that fulfills the task of model \mathbf{M} (see Definition 3.6), we will split the main task of Conv QA into multiple subtasks:

1. **Information Extraction:** Given a set of documents D , extract the textual content C_d of each document $d \in D$ and create a knowledge source P based on C_d of every document $d \in D$.
2. **Contextual Query Understanding:** Given a history H and a new question q_{i+1} , generate a contextualized question q_c based on H , such that $\mathcal{I}(q_c, q_{i+1}) = 1$.
3. **Passage Retrieval:** Given a contextualized question q_c and a knowledge source P , retrieve the most relevant passages p from P and combine them in an evidence set E .
4. **Response Generation:** Given a history H , a new question q_{i+1} and a set of passages E , generate an answer a to q_{i+1} based on E , so that $\mathcal{I}(q_{i+1}, a) = 1$.

There may exist other approaches to break down the task of **M** into sub-tasks, but for this thesis, we will focus on a solution based on the four sub-tasks outlined above. It is to be highlighted, that breaking the task of **M** implies also an order in which the sub-tasks have to be performed. Sub-task (1) will be performed once, while (2-4) will be repeated on every new question q_{i+1} .

In order to develop a system that can be applied to this abstract task, we match every task to a component. The *Information Extraction* sub-task will be solved by the *Extract* component, further detailed in Section 3.3.1. *Passage Retrieval* will be covered by the *Retriever* component, further described in Section 3.3.3. The *Response Generation* will be handled by the *Reader* component, more precisely in this thesis we will focus on LLMs with intrinsic parametric knowledge as *Reader*. This will lead to a RAG system consisting of the *Retriever* and *Reader*. This choice has been made due to the fact, that the latest research breakthroughs sparked the interest in RAG systems in comparison to classical Retriever-Reader systems (check therefore the related work Section 2.4). Details on the *Reader* component will be laid out in Section 3.3.4. In order to now handle conversations, a CQU unit as described in Section 2.2.2 is necessary to handle the sub-task of *Contextual Query Understanding*. Section 3.3.2 will dive into the details.

Figure 3.3 illustrates the combination of the four components, which make up the model **M** in their corresponding sub-tasks. This is the abstract Model **M** which leads, when applied to a real use-case, to the ConRAG system architecture.

3.3.1. Extract

The sub-task of *Information Extraction* was defined in the previous section as sub-task (1) of the model **M**. Given a set of documents D , the knowledge source P needs to be

3. Open-domain QA Chatbot

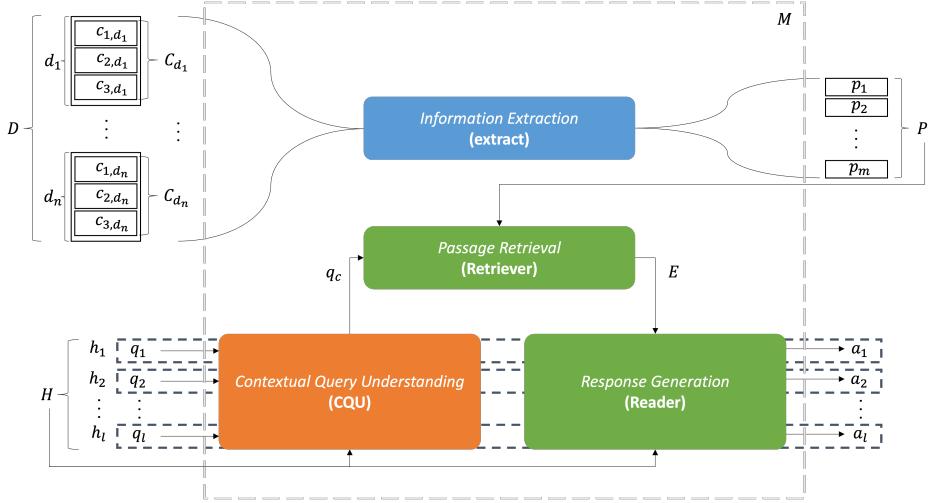


Figure 3.3.: Overview of the System Architecture in context of the sub-tasks of \mathbf{M}

extracted using *Information Extraction* techniques in combination with *Passage Extraction* operations. As synthetic data generation is also part of the extraction component according to Section 3.1, it will also be discussed in this section. This is originally not part of the model \mathbf{M} , but makes from a system architecture sense, to place these operations in this system component. Synthetic Data can be seen as a separate task, which has nothing to do with the original task of Conv QA, but is a necessary step in order to train components of the model \mathbf{M} or evaluate those.

Information Extraction: When it comes to extracting text from any document d , there are many approaches to choose from. Some extract structures, metadata, or similar, which can be further utilized, while others extract unstructured text only. In any case, this extraction process highly depends on the source document type. An HTML website requires different approaches and tools compared to a PDF, for instance. An example tool for direct extraction of PDFs is Py2PDF [PyPDF2, 2023]. Regardless of the source document and tool used to extract textual information C_d , there are two major possible outcomes given a set of documents D :

1. *Structured Extraction:* Denoted as $f_{StrucExt}(\cdot)$, in this extraction operation, C_d can be extracted into logical segments directly: $f_{StrucExt}(d) = \{c_{d_i} \subset C_d : i \in \{1, 2, \dots, n\}, C_d = \bigcup_{i=1}^n c_{d_i}\}$. If $f_{StrucExt}$ is applied to all documents d in D , the resulting set C contains all the outputs of $f_{StrucExt}$ for each document in D . Formally, $f_{StrucExt}(D) := C = \{C_{d_1}, C_{d_2}, \dots, C_{d_m}\}$, where each document d is transformed into a set of logical segments C_{d_i} representing its content.

3. Open-domain QA Chatbot

2. *Unstructured Extraction:* Denoted as $f_{UnstExt}(\cdot)$, when applied to D , it results in a concatenated text corpus $C_d = \{c_d\}$ containing the textual content of a document d : $f_{UnstExt}(D) := C = \{C_{d_1}, C_{d_2}, \dots, C_{d_n}\}$, where each document d is transformed into a single text snippet c_d representing its content.

In order to generate snippets p from C , *Passage Extraction* has to be applied on C .

Passage Extraction: The implementation of passage splitting depends on the nature of C and the desired granularity of the output. In general, there are three operations for constructing passages p based on a text snippet c_d :

1. *Paragraphs:* $f_p(\cdot)$ is an operation that transforms a collection C_d of texts c_d into a set of passages $P = \{p_1, p_2, \dots, p_n\}$. It operates similarly to $f_{StrucExt}(\cdot)$, but instead of D , it operates on C_d . The output of $f_p(\cdot)$ consists of passages p that represent logical segments of the text corpus c_d . Usually, it operates in a *rule-based* manner, meaning that a paragraph is defined by a token indicating a paragraph (e.g., $< p />$ in HTML). The length l of each p_i is variable, and the number of paragraphs $|P|$ can also vary.
2. *Snippets:* $f_s(\cdot)$, when you have a fixed passage length l , divides the concatenated text c_d into $|c_d| \bmod l + 1$ passages p per c_p . Alternative approaches may involve specifying minimum and maximum lengths, denoted as l_{\min} and l_{\max} . The exact point of division depends on whether a sentence ends within the specified window or not. If a sentence ending is found within the window, the snippet concludes at that point. Otherwise, it concludes at the end of the window. The individual length of the extracted passages $P = \{p_1, p_2, \dots, p_n\}$ is not fixed in the case of syntactic snippets, as well as the number of paragraphs $|P|$.
3. *Sliding Windows:* $f_w(\cdot)$ utilizes a window size l , a concatenated text c_d , and a step size s . The window slides over the text c_d , and the text within the window is used as a passage p . This results in $\frac{|c_d|-l}{s}$ passages, denoted as $P = \{p_1, p_2, \dots, p_n\}$.

These operations can be combined in a pipeline fashion. For example, first $f_p(C_d) := P_{paragraphs} = \{p_{paragraphs,1}, p_{paragraphs,2}, \dots, p_{paragraphs,n}\}$ is constructed, and afterward, on the logical paragraphs, $f_w(f_p(C_d)) := f_w(P_{paragraph}) = P_{window} = \{p_{window,1}, p_{window,2}, \dots, p_{window,i}\}$ is applied. The way these operations are combined is highly use-case specific and needs to be evaluated for each use-case individually. Another important factor influencing the decision regarding the parameters l , in general, the operation used, is the desired model for the *Retriever* and *Reader* as they may be trained on a

3. Open-domain QA Chatbot

specific token length as input or have a maximum length of tokens they accept as input or require a certain clarity of data points.

Synthetic Training Data Generation: This can be interpreted as a knowledge distillation task. The main idea is to use a model $ED(\cdot)$ to generate a synthetic dataset for the sub-task and problem field of *Passage Retrieval* and *Response Generation*:

$$ED(P) := (P, Q) \quad (3.1)$$

Where $P = \{p_1, p_2, \dots, p_n\}$ corresponds to a corpus of passages to retrieve from and Q is a set of questions. If $\mathcal{I}(q, p) = 1$, this means that the passage reflects the question's search intent. Given this task, there are several synthetic dataset types:

1. *Questions given Context* $ED_{qp}(p) := q_s$: Given a passage p , generate a synthetic question q_s that satisfies the desired search intent. Applying $ED_{qp}(P)$ to a set of passages will generate a set of questions $Q_s = \{q_1, q_2, \dots, q_n\}$, with one question for every passage, i.e., $|Q_s| = |P|$. ED_{qp} can also be applied multiple times with different intents to generate multiple different questions $q_{s,j,1}, q_{s,j,2}, \dots, q_{s,j,i}$ for a single passage p_j .
2. *Question-Context-Answer Triples* $ED_{qpa}(p) := (q_s, p, a_s)$: Given a passage p , generate a synthetic question q_s that satisfies the search intent and provides an answer a_s , where $\mathcal{I}(q_s, a_s) = 1 \wedge \mathcal{I}(q_s, p) = 1 \wedge \mathcal{I}(a_s, p) = 1$. The result of applying $ED_{qpa}(P)$ to a set of passages is a set of question-passage-answer triples $QPA_s = \{(q_{1,s}, p_1, a_{1,s}), (q_{2,s}, p_2, a_{2,s}), \dots, (q_{n,s}, p_n, a_{n,s})\}$, where $|QPA_s| = |P|$. ED_{qpa} can also be applied multiple times with different intents to generate multiple different questions $q_{s,j,1}, q_{s,j,2}, \dots, q_{s,j,i}$ and corresponding answers $a_{s,j,1}, a_{s,j,2}, \dots, a_{s,j,i}$ for a single passage p_j .

While the previous two approaches focus on single tuples of either (q, p) or triples of (q, p, a) , there is also a problem field of generating conversations based on passages $P = \{p_1, p_2, \dots, p_n\}$ from the same underlying document d . Therefore, the task given in Equation 3.1 is changed to:

$$ED(P) := (H, P) \quad (3.2)$$

Where H corresponds to a set of conversation histories h containing multiple turns. The first turn is always a question q_1 , followed by an answer a_1 based on a passage $p \in P$, also given a search intent I over the whole history h . Synthetic datasets for this task can be generated in the following ways:

3. *Conversational Question-Context-Answer Histories* $ED_{cqpa}(E) := H_s$: Given a subset of passages $E \subset P$, the task of the model ED_{cqpa} is it to construct a realistic conversation with an intent over the corpus E . The result of applying $ED_{cqpa}(E)$ on a set of passages is a set of conversation histories $H_s = \{h_1, h_2, \dots, h_n\}$. ED_{cqpa} can also be applied multiple times with different intents in order to generate multiple different conversations $h_{s,j,1}, h_{s,j,2} \dots, h_{s,j,i}$ for a single subset of passages $E \subset P$.

There are several ways to implement the models ED_{qp} , ED_{qpa} or ED_{cqpa} . Implementations of those approaches will be discussed in Chapter 4.

3.3.2. Contextual Query Understanding

Previously the the model **M** for Conv QA was introduced. The CQU component is responsible for the *Contextual Query Understanding* task. Essentially, the goal of this task is to identify the necessary information in the history H of a conversation to adapt the expression of the question q_{i+1} into a contextualized question q_c , such that $\mathcal{I}(q_c, q_{i+1}) = 1$. This adaptation is crucial, as the *Passage Retrieval* task is executed using a single question q rather than the entire history H . Therefore, language challenges may arise between turns, such as pronoun resolution (e.g., it, he, she) or turn references.

$$CQU(h_{i-k:i}, q_{i+1}) := q_c \mid \mathcal{I}(q_c, q_{i+1}) = 1 \quad (3.3)$$

The depth of the history $h_{i-k:i}$ that *CQU* operates on is a hyperparameter to be chosen. Generally, there are three levels of depth:

1. **K-many:** Consideration of the last k turns.
2. **All:** Utilization of the entire history, $k = i - 1$.
3. **Memory:** Inclusion of not only the current history but also previous chat histories by this user. The concept of memory will not be further discussed in this thesis, as it opens a whole new field of research and was excluded in the problem statement in Section 3.1.

The decision on depth already implies limitations, as the *K-many* approach may lead to reference errors in a conversation, while too large k increases the computational complexity. This leads to $h_{i-k:i} \subset H$, which is used by the *CQU*.

The CQU unit itself can be either an implicit or explicit method. Implicit methods use a one-step approach and are based on the *Transformer* architecture, while explicit methods apply a two-step approach and can vary in implementation:

$$CQU_{\theta}(h_{i-k:i}, q_{i+1}) := q_c \mid \mathcal{I}(q_c, q_{i+1}) = 1 \quad (3.4)$$

Equation 3.4 shows the transformer-based CQU unit with trainable parameters. During the training of the parameters, it is necessary to have an annotated dataset of turns and contextualized questions $q_{c,i}$ for every turn question q_i . The loss of the CQU will be defined between the human-annotated contextualized question q_c^* and the generated contextualized question q_c .

Explicit approaches consist of two operations:

$$rel_{CQU}(x, q_{i+1}) := \begin{cases} x \in C_{q_{i+1}}, & \text{if } x \text{ is relevant for } q_{i+1} \\ x \notin C_{q_{i+1}}, & \text{if } x \text{ is not relevant for } q_{i+1} \end{cases} \quad (3.5)$$

$$rep_{CQU}(C, q_{i+1}) := q_c \quad (3.6)$$

Where rel_{CQU} is an operation to determine, for a string $x \subset h_{i-k:i}$ (which can be a token, word, or snippet of a turn $h_{i-k,i}$), if it is relevant context to enrich the question q_{i+1} . If so, it is added to the context set $C_{q_{i+1}}$. The operation rep_{CQU} works over $C_{q_{i+1}}$ and q_{i+1} to rephrase and generate a contextualized question q_c .

The relevance operation can be implemented in multiple ways. Some include seq-2-seq models or entity matching. The rephrasing operation is a typical MRC summary task. Its inputs are the set of context-relevant information $C_{q_{i+1}}$ and the question q_{i+1} .

In general CQU is an ongoing field of research, especially for follow-up questions and resolving ambiguity. The latest research utilizes more often LLMs for this task. This thesis won't focus on CQU. Please refer to [Mao et al., 2023] for a state-of-the-art approach.

3.3.3. Retriever

Previously we have divided the main task of a model \mathbf{M} for Conv QA into multiple sub-tasks. The Retriever component handles the *Passage Retrieval* task. The goal of this component is to identify relevant passages p from the knowledge source P given a contextualized question q_c . The contextualized question q_c is generated by the CQU unit (see Section 3.3.2). Approaches, that are using instead of q_c the whole history H for retrieval, will not be discussed in this thesis. The identification of passages includes a scoring of relevance for each passage p given a question q_c . The scoring function $p_{Ret}(p|q_c)$ is defined as:

$$p_{\text{Ret}}(p|q_c) = \text{Score}(q_c, p) \quad (3.7)$$

Here, $\text{Score}(q_c, p)$ is a value determining the relevance of a passage p in relation to q_c , enabling the ranking of passages p given a question q_c . Based on the score, the passages will be ordered in descending order, and the top- k passages will be combined into an evidence set E_{q_c} . A different approach is to set a threshold for the $\text{Score}(q, p)$ and add all passages p to E_{q_c} for which $\text{Score}(q, p)$ surpasses the threshold. Within the evidence set E_{q_c} , a passage p will be represented as a tuple $(p, \text{Score}(q_c, p))$. Concerning the evidence set E_{q_c} in relation to the question q and the underlying search intent, the following cases exist:

$$\forall q \in Q, \exists! E_{q_c} \subset P := \begin{cases} 1. \mathcal{I}(q, E_{q_c} = \emptyset) = 1 \\ 2. \mathcal{I}(q, E_{q_c} = \{p\}) = 1 \\ 3. \mathcal{I}(q, E_{q_c} = \{p_1, p_2, \dots, p_n\}) = 1 \end{cases} \quad (3.8)$$

For every question, there exists an evidence set. In order for intent $\mathcal{I}(q, E_{q_c})$ to be 1, there are three cases: either the evidence set E_{q_c} has to be empty, E_{q_c} has to have exactly one element p , or E_{q_c} has to have multiple passages p . To put it more simply, for a given question, the answer is determined either by the evidence that there is no supporting passage p , or that there is exactly one passage p containing the requested information, or that multiple passages p contain the necessary information to answer the question q .

Retriever: Depending on the chosen method for a Retriever, the Retriever component can either be trainable or have adjustable parameters:

$$r_{\text{Ret}, \theta} = \text{Retriever}_{\theta}(q_c, P) \quad (3.9)$$

$$r_{\text{Ret}, \Theta} = \text{Retriever}_{\Theta}(q_c, P) \quad (3.10)$$

Here, θ is the set of trainable parameters, and Θ is a set of adjustable parameters. Examples of the first type of retriever include dense retrievers like *DPR*, as introduced in Section 2.1.4. Examples of the second type of retriever include sparse retrievers like *BM25*, also introduced in the same section. Regardless of the type of retriever, the operation from Equation 3.7 is applied to each passage $p \in P$, and the top- k passages are selected to form the evidence set E_{q_c} . As outlined in Equation 3.8, some questions may require no evidence to fulfill the search intent. However, with the proposed Retrievers,

3. Open-domain QA Chatbot

this is not possible. A retriever r_{Ret} will always assign a $\text{Score}(q, p)$ to every passage given a question and its search intent. Filters are applied to filter out passages with a low score, but this is not part of the retriever itself. The evaluation of the evidence with respect to the search intent and correctly identifying the case at hand will be handled by the Reader component (see Section 3.3.4).

Mixture-of-Experts: To enhance the quality of the evidence set E_{q_c} , state-of-the-art research and related work have emphasized the effectiveness of combining multiple retrievers (see Section 2.4). This combination can be achieved in various ways, which can also be combined further:

- **Re-Ranking:** This involves the combination of multiple retrievers in a pipeline fashion. In most cases, a fast but imprecise Retriever is used to retrieve a large evidence set, e.g., $r_{BM25}(q_c, P) = E_{q_c,1}$, where $|E_{q_c,1}| = 100$. Subsequently, a second, more accurate Retriever is employed to re-rank the evidence set E_{q_c} , e.g., $r_{CE}(q_c, E_{q_c,1}) = E_{q_c,2}$, where $|E_{q_c,2}| = 25$.
- **Ensemble:** This idea involves combining multiple retrievers into a single retriever. This can be accomplished by either concatenating the evidence sets $E_{q_c,j}$ or by aggregating the scores of the individual retrievers r_j into a single score. Formally:

$$\forall r \in R : E_{q_c} = \bigcup_{j=1}^{|R|} r_j(q_c, P) \quad (3.11)$$

$$\forall p \in P : \text{Score}(q_c, p) = f(r_1(q_c, p), r_2(q_c, p), \dots, r_{|R|}(q_c, p)) \quad (3.12)$$

Here, $f(\cdot)$ represents a function that combines the scores of the individual retrievers r_j , for example, using $\max(\cdot)$. R is a set of retrievers r .

- **Weighting:** This approach involves running two retrievers in parallel and multiplying their scores for each passage p . It can be seen as a sub-case of ensemble. Formally:

$$\forall p \in P : \text{Score}(q_c, p) = \prod_{j=1}^{|R|} \alpha_j \cdot r_j(q_c, p), \sum_{j=1}^{|R|} \alpha_j = 1 \quad (3.13)$$

Here, R represents a set of retrievers r_j , and α_j is the weight assigned to each retriever r_j . These weights α_j can be either fixed or trainable.

3. Open-domain QA Chatbot

Fine-Tuning: A potential issue that may arise for a retriever and a knowledge base P is the misalignment between the underlying formats of the passages p and the formats on which the retriever $r_{Ret,\theta}$ was trained. To address this discrepancy, fine-tuning can be employed to adjust the retriever based on the specific formats used in the given use case.

Given a dataset of tuples (q, p) , as described in Section 3.3.1, a retriever with trainable parameters $r_{Ret,\theta}$ can undergo fine-tuning. The first step involves creating the training data $\mathcal{TD} = \{\langle q_i, p_i^+, p_{i,1}^-, \dots, p_{i,n}^- \rangle\}_{i=1}^m$. In this dataset, we already have tuples that contain a question q_i and a positive passage p_i^+ . The primary objective is to sample n negative passages $p_{i,j}^-$ for each positive passage p_i^+ . Various methods can be employed [Karpukhin et al., 2020]:

1. Randomly sampling passages from the knowledge source P , where $p_{i,j}^- \in P$ and $p_{i,j}^- \neq p_i^+$.
2. Retrieving evidence passages using a high-performing Out-of-Domain (OOD) retriever $r_{Ret,OOD}$ and sampling from the evidence set $p_{i,j}^- \in E_{q_i}$, ensuring $p_{i,j}^- \neq p_i^+$.
3. Sampling passages from other tuples (q_k, p_k) in the dataset, where $p_k \neq p_i^+$.

Method (1) is the simplest but does not yield high-quality negative passages $p_{i,j}^-$. Method (2), while more complex, typically provides a higher-quality evidence set E_{q_i} than method (1). Method (3) can also be applied in-batch, known as *in-batch negatives*, during the training process. This means that given a batch B of $|B|$ -many (q_i, p_i^+) tuples, the negative passages $p_{i,j}^-$ are derived from the positive passages of the other tuples in the same batch B . Consequently, this results in $|B| \times |B| - 1$ negative passages $p_{i,j}^-$ per batch B per question q_i .

After the construction of the training dataset \mathcal{TD} , the retriever $r_{Ret,\theta}$ can be trained straightforwardly. The forward pass is simply the score prediction of the retriever for every passage, positive and negative, within a tuple of the \mathcal{TD} . Therefore the loss function can be defined as follows:

$$\mathcal{L}(q_i, p_i^+, p_{i,1}^-, \dots, p_{i,n}^-) = -\log \frac{e^{\text{Score}(q_i, p_i^+)}}{e^{\text{Score}(q_i, p_i^+)} + \sum_{j=1}^n e^{\text{Score}(q_i, p_{i,j}^-)}} \quad (3.14)$$

The calculated loss will be used in backpropagation to update θ .

The process of fine-tuning needs adaptation when using *synthetically generated data* due to the potential quality issues of the synthetic dataset. To address this, the dataset must be iteratively filtered during the training process. Given a synthetic dataset $\mathcal{TD}_f =$

3. Open-domain QA Chatbot

$\{(q_i, p_i^+)\}_{i=1}^m$, a retriever $r_{\text{Ret}, \theta}$ with trainable parameters θ , and a high-performing out-of-domain (OOD) retriever $r_{\text{Ret}, \text{OOD}}$, the following procedure is based on the proposed approach of PROMPTAGATOR [Dai et al., 2022b]:

1. Extend $\mathcal{T}\mathcal{D}_f$ by adding n negative passages p_i^- for every tuple using Method (2) from above with $r_{\text{Ret}, \text{OOD}}$.
2. Train $r_{\text{Ret}, \theta}$ for s iterations on $\mathcal{T}\mathcal{D}_f$ while applying in-batch negatives, as described in Method (3) above.
3. Filter $\mathcal{T}\mathcal{D}_f$ using $r_{\text{Ret}, \theta}$ and $r_{\text{Ret}, \text{OOD}}$ by removing all tuples (q_i, p_i^+) where p_i^+ is not in the Evidence set E_i , given a parameter k , for either $r_{\text{Ret}, \theta}$ or $r_{\text{Ret}, \text{OOD}}$.

Figure 3.4 illustrates the fine-tuning process with synthetic data. Steps (2) and (3) are repeated cyclically, with step (3) concluding each cycle. After s epochs, $\mathcal{T}\mathcal{D}_f$ is filtered once, retrained for s epochs, refiltered, and finally trained again for s epochs. This process may appear counterintuitive, as the retriever is trained on the dataset it's supposed to filter. However, experiments from related work [Dai et al., 2022b] have demonstrated that this approach still yields favorable results.

Alternative approaches to filtering synthetic datasets may rely solely on high-performing $r_{\text{Ret}, \text{OOD}}$. For a synthetic question q_s , the evidence set E_{q_s} is generated using $r_{\text{Ret}, \text{OOD}}$ with high values for k . If $E_{q_s} = \emptyset$, the tuple (q_s, p_s^+) is excluded from $\mathcal{T}\mathcal{D}_f$. This approach is simpler but may result in a performance loss. It can be improved by incorporating multiple different high-performing $r_{\text{Ret}, \text{OOD}}$ and combining via union their evidence sets E_{q_s} .

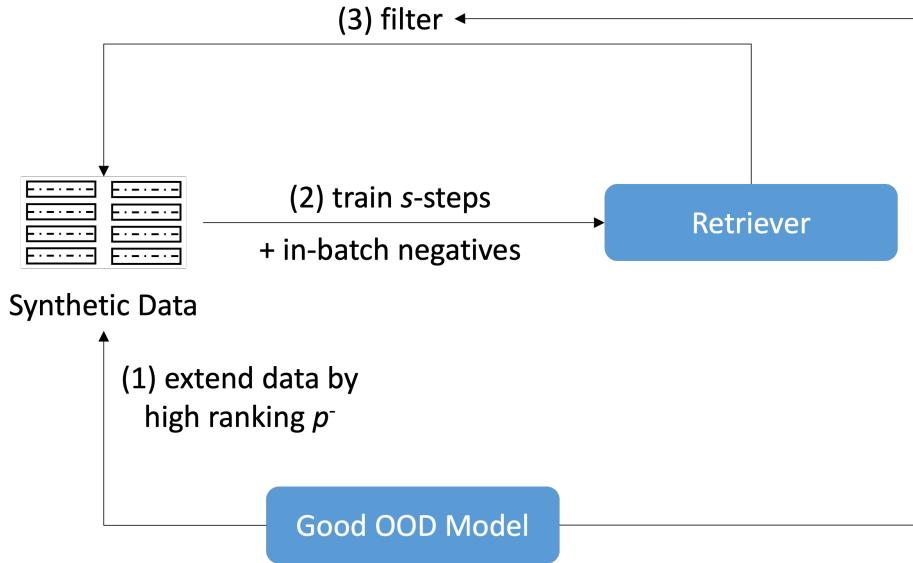


Figure 3.4.: Fine-Tuning Process for Retriever

3. Open-domain QA Chatbot

Metadata-Filtering: In certain search scenarios, it may be necessary to filter the knowledge source P by metadata. State-of-the-art open-domain datasets and benchmarks usually don't develop solutions for this problem. However, in real-world scenarios, this is a common user desire.

A filtering request can be understood as the already established intent-fulfilling retrieval task performed over P , extended by a consideration of the metadata M_d of the documents $d \in D$ from which p originates. As established in Definition 3.1, $p = (\text{content}, \text{UID}_p, \text{UID}_d)$ and $d = (C_d, M_d, \text{UID}_d)$. Therefore, $\forall p \in P : \exists! M_d$ such that $M_p = M_d$. The task of filtering can therefore be understood as follows:

$$m_{\text{Ret}}(p|q_c) = \begin{cases} p \in E_m, & \bigcap_{i=p_1}^{E_{q_c}} M_i \subset M_p \\ p \notin E_m, & \bigcap_{i=p_1}^{E_{q_c}} M_i \not\subset M_p \end{cases} \quad (3.15)$$

Here, a passage p is only added to the evidence set E_m if the intersection of all metadata of the passages, in the perfect intention-fulfilling evidence set E_{q_c} of a question q_c , is a subset of the metadata M_p of the passage p .

In the retrieval process itself, the two tasks of metadata-filtering (m_{Ret}) and scoring (p_{Ret}) can be executed sequentially or combined. This depends on the implementation of the actual index. Some approaches towards metadata-filtering in Conv QA are the following:

1. **Metadata Passage Integration:** In this approach, the metadata M_p is integrated into the *content* of a passage p . This enables a one-step approach to retrieval. The retriever will perform scoring and metadata filtering simultaneously. However, in this approach, the original form of the metadata M_p , a collection of key-value pairs, will not be utilized. Instead, the metadata will be added implicitly (e.g., embedding addition) or explicitly (e.g., attaching keywords) to the *content* of p , leading to a loss in the functionality of the metadata, which could cause issues like a not fully applied filtering.
2. **Separate Metadata Index:** In this approach, two indices are created and utilized, one representing the passages' *content* and the other representing the passages' metadata. During retrieval, the retriever follows a two-step approach. Ideally, the metadata-filtering task is executed first, and then the scoring task is performed on E_m : $E_{q_c} = p_{\text{Ret}}(m_{\text{Ret}}(q_c, P))$.
3. **Hierarchical Index:** The hierarchical index utilizes a scoring function p_{Ret} as the metadata filter. In this approach, in addition to the passage index, a secondary

metadata index is employed, typically representing the set of documents D . Instead of using M_d as key-value pairs, a document d is encoded in this index as *content*, resembling a passage p . This *content* may include a document summary, a concatenation of metadata keywords, or similar information. During retrieval, the retriever initially generates an evidence set E_m by executing $p_{\text{Ret}}(q_c, D)$, followed by the execution of $p_{\text{Ret}}(q_c, E_m)$ to form the ultimate evidence set E_{q_c} : $E_{q_c} = p_{\text{Ret}}(p_{\text{Ret}}(q_c, D))$. This approach carries the same issues as (1) but is useful when the metadata M_d are not simple entity or number-based values, but rather more complex information like a document summary.

3.3.4. Reader

The Reader component handles the *Response Generation* task. The goal of this component is to extract an answer a based on the evidence set E_{q_c} given a question q . The answer a can be a single token, a span of tokens, or a set of tokens. The answer extraction can be defined as follows: $r_{\text{Read}}(E_{q_c}, q) := a$. Depending on the desired utility of the Reader, the Reader can also incorporate *Context Query Understanding* capabilities and therefore receive, in addition to the evidence set E_{q_c} and question q , the history H as input. This leads to the following two variations:

$$r_{1\text{-Reader}}(E_{q_c}, q_c) := a \quad (3.16)$$

$$r_{k\text{-Reader}}(E_{q_{i,c}}, q_i, H) := a_i \quad (3.17)$$

Enhancing the theoretical model \mathbf{M} from Definition 3.6 with the system requirements defined in Section 3.1, the challenges for the Reader component can be broken down into the following:

1. Generate an answer a based on the evidence set E_{q_c} that satisfies the search intent $\mathcal{I}(q_c, a) = 1$.
2. As stated in Equation 3.8, there are three cases for the evidence set E_{q_c} in relation to the question q and the underlying search intent. The Reader component has to identify the case at hand and determine the final gold-evidence set $\hat{E}_{q_c} \subset E_{q_c}$, $\mathcal{I}(q_c, \hat{E}_{q_c}) = 1$.
3. Identify ambiguity in questions q_i and determine a corresponding clarification question a_i which steers the user towards a more specific question in the next turn $i+1$.

3. Open-domain QA Chatbot

4. Chatlike behaviour is expected for Conv QA. This may include the ability of *Contextual Query Understanding*, but also challenges like (3). In general, this challenge describes the model’s capability to generate a human-like conversation.

At the core of the Reader’s implementation is the evidence set E_{qc} . It is constructed using the methods outlined in Section 3.3.1. However, it hasn’t been discussed yet what exactly is a good format for a passage p in the evidence set E_{qc} for the Reader. Basically, all variations of the passages p can be broken down to the following scale:



Figure 3.5.: Scale of Passage implementations

Most datasets used for QA, Conv QA, or training of PrLM find themselves on the right end of this scale. So, the LLMs get trained on semantically correct and most of the time content-complete passages p . Exceptions are datasets for multi-hop QA where the passages are not content-complete, but the aggregation of multiple passages is needed. This evokes the question of how to construct the passages, keeping in mind how the used *Transformer* models have been trained. However, this thesis will not further evaluate the influence of different passage extraction variations, as there is currently no comprehensive research on this topic to the best of our knowledge.

The first challenge (1) of the Reader is closely related to the task MRC and classical QA. Let’s assume that the evidence set is the perfect evidence set that it takes to answer question q_c : $E_{qc} = \hat{E}_{qc}$. The model from Equation 3.16 can be easily trained with a supervised dataset. To fulfill this task, a dataset consisting of question-context-answer triples is necessary, where the context corresponds to \hat{E}_{qc} and therefore is a list of passages. Depending on the dataset, the model can even be trained to fulfill multiple question types. This can be achieved by varying the question types of the questions in the triples.

The second challenge (2) is more complicated, as in a modular implementation, the Retriever will always pass an evidence set E_{qc} with a fixed length $|E_{qc}| = k$ to the Reader. The Reader, therefore, has to incorporate its own mechanism of passage identification. This naturally leads to a trade-off that has to be made. Setting a small k -value for the Retriever may reduce the Hit-Ratio (HR), and therefore the correct passages won’t necessarily be included in the evidence set E_{qc} . On the other hand, the Reader can’t handle an arbitrary number of passages, due to the nature of the underlying *Transformer* models. Solutions to this problem may include:

- **Re-Ranker:** The idea of Re-Rankers was already introduced in Section 3.3.3. Basically, before handing E_{q_c} to the Reader, the evidence set is re-ranked by a second Retriever. This second Retriever can be a simple filter over the scores of the passages or an even more complex model than the first Retriever.
- **Compression:** A compressor can iterate over a larger set E_{q_c} than the Reader itself, as it only operates on one passage at a time. The goal of the compressor is to identify the most relevant information of a passage $p \in E_{q_c}$ given the question q_c . This compressed information will then create the new evidence set E_{q_c} . This approach reduces the number of tokens the evidence set E_{q_c} consists of and therefore enables larger k -values for the Retriever. It is similar to a Re-Ranker.
- **Multi-Retrieval:** The concept of Multi-Retrieval involves triggering multiple retrievals until a certain level of satisfaction is reached regarding the evidence set E_{q_c} . The specific approaches may vary and will not be further discussed in this thesis.
- **Trained Reader:** The core of this idea is to train a Reader on a dataset consisting of question-context-answer triples, where the context is a list of passages, which not necessarily is the gold-evidence set $E_{q_c} \neq \hat{E}_{q_c}$. The Reader will learn to identify the correct passages and therefore can be used to filter the evidence set E_{q_c} himself.
- **Parametric Knowledge:** This approach is the fundamental idea of RAG. The underlying *Transformer* of the Reader gets trained on the knowledge source P to store implicitly the knowledge. This may lead to a better understanding of the evidence set E_{q_c} and give the Reader the ability to access knowledge, which is not included in E_{q_c} (e.g., when $\hat{p} \notin E_{q_c}$). The fine-tuning itself is the common adoption of a PrLM to a small new knowledge source P . During fine-tuning the PrLM is tasked to predict the next token and therefore learn the knowledge source P specific language and knowledge.

Given the (3) challenge, the Reader must be able to determine ambiguous questions q_i and generate a clarification question a_i . This task consists of two steps:

1. **Identify Ambiguity:** Given a question q_i , the Reader $r_{\text{Read}}(\cdot)$ has to determine if the question q_i is ambiguous. A question is ambiguous if there are multiple answers that fulfill the question's search intent, but themselves satisfy different intents ($\mathcal{I}(q_i, a_{i,1}) = 1 \wedge \mathcal{I}(q_i, a_{i,2}) = 1 \wedge \mathcal{I}(a_{i,1}, a_{i,2}) = 0$). The ambiguity in q_i can either originate in a factoid ambiguity (e.g., $\mathcal{I}(q_i, p_1) = 1 \wedge \mathcal{I}(q_i, p_2) = 1 \wedge \mathcal{I}(p_1, p_2) = 0$) or a linguistic ambiguity (e.g., ambiguous cross-references).

3. Open-domain QA Chatbot

2. **Resolve Ambiguity:** Given an ambiguous question q_i , the Reader $r_{\text{Read}}(\cdot)$ has to generate a clarification question a_i that resolves the ambiguity via the user's next question q_{i+1} . Meaning that $\exists! a_{i+1} : \mathcal{I}(q_{c,i+1}, a_{i+1}) = 1$.

As can be observed, *Clarification Questions* can only be generated and resolved in a dialog and not in a single-turn QA scenario. Therefore, it is necessary to train a Reader on a conversational dataset. The dataset is a history H , providing for every turn a quadruple of the question-contexts-answer-ambiguity flag. The answer is an annotated clarification question in case of a positive ambiguity flag. Therefore, the *Transformer* model can be trained to identify ambiguity in a question given context and generate a clarification question.

This history-based dataset leads to the challenge (4) of a Reader in this thesis context. Challenge (4) is a more abstract requirement for the reader. It requires the language and generated text to linguistically mimic the language and behavior humans would use in a conversation. In order to fine-tune a *Transformer* model to be able to perform text generation which is more human conversation-like, a dataset of human conversation histories H is needed. This involves a diverse set of training and fine-tuning strategies, which will not be discussed in this thesis. As a reference for this topic, the reader is referred to the work of Touvron et al. [Touvron et al., 2023].

As observed, the Reader's ability to address challenges (1-4) is highly dependent on the training datasets and fine-tuning methods. The challenges defined in this thesis do not cover every possible requirement that someone might have for their Conv QA system. Naturally, larger and more capable models exhibit higher zero- and few-shot performance. Previous solutions mainly focused on fine-tuning LLMs to address challenges in the Reader (e.g., Noise in the Evidence Set or Negative Rejection). However, as this thesis aims to provide a solution for resource-constrained systems (see Section 3.1, Objective 5 of the QA system), the following approach will be presented based on simplifying the overall task of a Reader to address these challenges without fine-tuning.

Chain-, Chain-of-Thought- and Agent-Reader: Drawing inspiration from the work of Langchain [Langchain, 2023], Wei et al. [Wei et al., 2022] and Kuhn et al. [Kuhn et al., 2023], the diverse task of a Reader can be broken down into a sequential decision-based procedure. The core idea is to break down the multiple challenges a Reader has to solve into smaller tasks that are easier to solve independently. An easy implementation could involve breaking down the holistic task of the reader into simple sub-tasks that follow a sequential order. One sub-task is for example: *Is the question ambiguous?*. The Reader will be prompted to solve this sub-task, and depending on the result, other

sub-tasks will be triggered. This is a rigid activity diagram-like implementation and is similar to the concept of *Chains* in *Langchain*. This can be further improved by adding more flexibility, where the Reader is fine-tuned in a *CoT* manner to reason over a question, given the evidence and the history. The single challenges or sub-tasks can, therefore, be interpreted as reasoning steps. This requires a whole new dataset for fine-tuning. Extending this idea further, the Reader is developed into an agent, similar to the *ReAct* approach, where it can reason over a question given the evidence and the history, and also has the ability to take certain actions, such as retrieving more evidence or re-running retrieval with a variation of the contextualized question. These are just some ideas on how the Reader is further extended to increase the quality of the QA system. However, this thesis will not further evaluate this idea, as it is not the main focus.

3.4. Summary

The main contribution of this thesis is a holistic framework covering the entire Conv QA system, from a collection of documents to detailed approaches used for specific components. To the best of our knowledge, no related work has yet provided such a comprehensive perspective on the problem field of Conv QA.

An often overlooked aspect in related research work is the step of *Extraction*. This is due to the fact that most research work is based on already existing datasets. However, in real-world scenarios, the extraction of passages from documents is a crucial step. This thesis breaks down the extraction into simple operations that can be combined to form an extraction pipeline. The parameters of the pipeline can be adjusted to determine the final passage model.

The thesis does not elaborate on new approaches/algorithms to retrievers itself but rather identifies the problems and tasks the retriever component has in a Conv QA system. Given a retriever and a knowledge base, problems can arise from the retriever towards the data (e.g., misalignment between the underlying formats of the passages and the formats on which the retriever was trained). Therefore, the thesis introduces a fine-tuning method based on generated data using a LLM and multiple state-of-the-art optimizations. Other problems can arise from the retriever towards the evidence set, namely that the evidence set's quality is not sufficient. In order to address this, the thesis elaborates on the concept of *Mixture-of-Experts* as it is an approach that can be used to enhance the retrieval quality without the need to reinvent a new retriever.

Lastly, this thesis introduces the approach of breaking down the task of a reader into challenges. These challenges can be of any nature (e.g., identifying ambiguity in

3. Open-domain QA Chatbot

questions or assessing the completeness of the evidence set). Previous related work has not yet broken down the challenges towards a reader in such a manner to the best of our knowledge. Given these small and well-defined challenges, it's easy to identify possible solutions regarding the reader. In this thesis, multiple such solutions based on related work are introduced, and new ones for specific challenges are proposed, mainly focusing on fine-tuning. On top of that, the thesis introduces the idea of using these smaller, finite, and simpler challenges of the whole reader task to create a sequential decision-based procedure, which could enable smaller LLMs to solve the reader task without performance losses compared to larger LLMs.

4. Experimental Evaluation

The previous chapter, Chapter 3, layed out a holistic framework for Conv QA. This chapter evaluates the applicability of the established system in a real-world scenario. Section 4.1 describes the available data for the real-world scenario and delves into applied data augmentation techniques. Section 4.2 introduces the metrics used to evaluate the performance of the individual system components, as well as the complete Conv QA system. These metrics are selected based on those used in related work. Section 4.3 details the experimental setup, implementation specifics, and provides an implementation framework for similar use cases. Finally, Section 4.4 presents both quantitative and qualitative results from the experiments.

4.1. Data

To evaluate the proposed system architecture and framework discussed in Chapter 3, we will implement and assess its performance using the PDFs of the Heidelberg University's examination regulations (Examination Regulation (ER)). The university provides these regulations on two separate websites: one in German¹ and the other in English². It's important to note that only the German ER holds legal authority, and in cases of ambiguity between the English translation and the German original, the German ER version takes precedence.

Both websites offer structured access to the ER of nearly all faculties at the Heidelberg University. However, it's worth mentioning that the German website is regularly updated, while the English version primarily contains outdated ER. Since there is no centralized source for accessing all English ER, the decision was made to utilize both the outdated English ER and the current German ER for this thesis. This decision aligns with the primary objective of this thesis, which is to demonstrate the PoC of the introduced framework. Obtaining the latest English ER of the Heidelberg University is beyond the scope of this thesis.

¹[https://www.uni-heidelberg.de/de/studium/studienorganisation/downloadcenter/
studien-und-pruefungsordnungen](https://www.uni-heidelberg.de/de/studium/studienorganisation/downloadcenter/studien-und-pruefungsordnungen)

²https://www.uni-heidelberg.de/courses/download/examination_rules_regulations.html

4. Experimental Evaluation

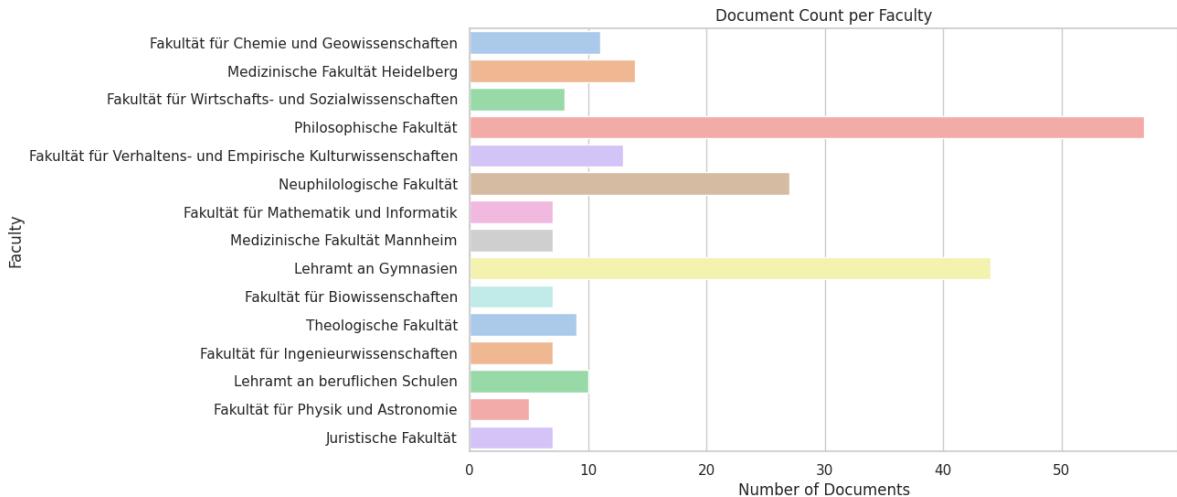


Figure 4.1.: Documents per Faculty of the German ER Dataset

As part of the experiments outlined in Section 4.3, the knowledge source of the German ER will be translated. Additional details on this process can be found in Section 4.3.1.

The two datasets (German/English) differ in their statistics. Therefore, Figure 4.1 displays the number of documents per faculty regarding the German ER. Figure 4.2 provides the same statistics for the English ER dataset. In general, the German dataset consists of 233 individual ER, while the English dataset contains only 151 individual ER. As observed, some faculties are overrepresented in both datasets, such as the philosophical faculty. The impact of the document distribution on the system will be discussed in Section 4.4.

While the previous statistics described the underlying documents of the PoC, what's even more interesting are the statistics related to the final knowledge source of passages that will be used throughout the PoC. Details on how these passages have been extracted from the PDFs will be provided in Section 4.3.1. In this section, we will focus on the statistical analysis.

Figure 4.3 illustrates the distribution of passages across documents in the German dataset, offering insights into the diverse influence that documents can have on the knowledge source in terms of the number of passages. These differences in passages are mainly due to variations in the length of the examination regulations (ER). The total number of passages in the German dataset is 39,039, resulting in an average of 167.55 passages per document. In comparison, the English dataset has an average of 212.94 passages per document. The statistics for the translated dataset are similar to those of the German dataset. An apparent difference between these two datasets is that the distribution of German passages per document exhibits a slight camel-like shape with a local maximum, while the English distribution follows a clear Gaussian pattern.

4. Experimental Evaluation

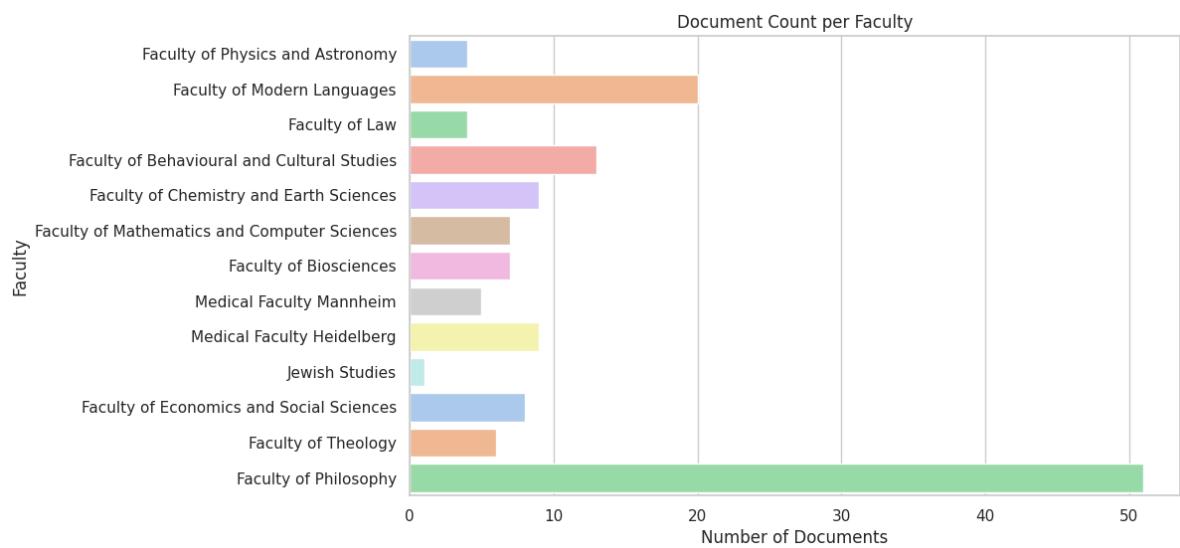


Figure 4.2.: Documents per Faculty of the English ER Dataset

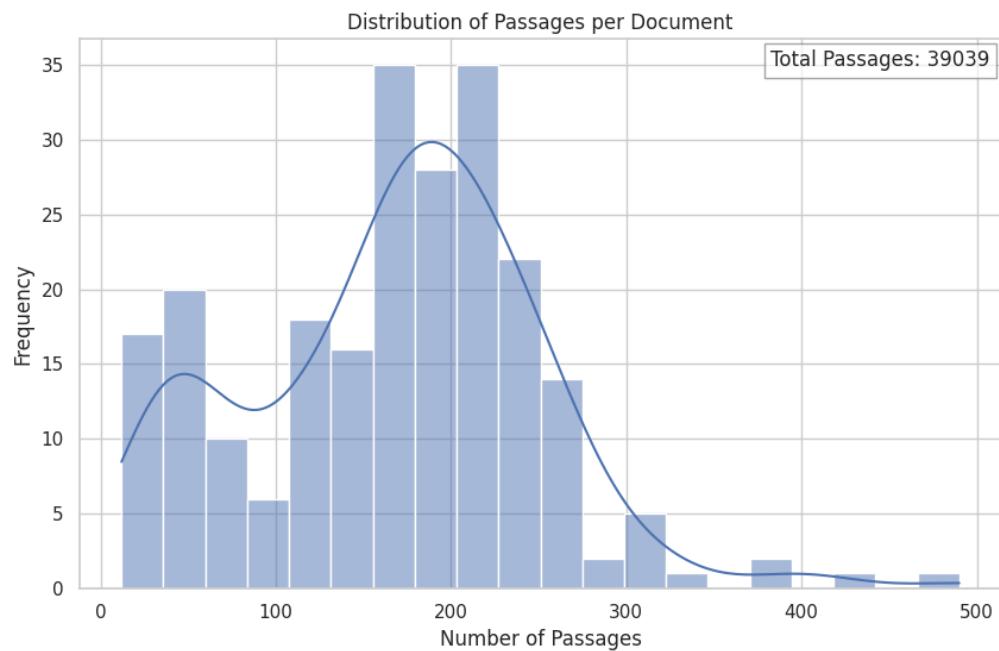


Figure 4.3.: Passages per Document of the German ER Dataset

4. Experimental Evaluation

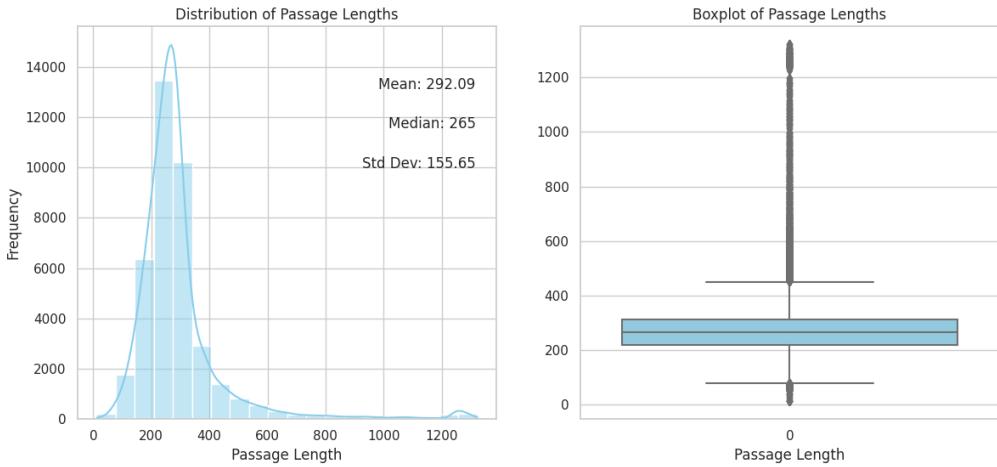


Figure 4.4.: Passage Length Distribution of the German ER Dataset

However, these distribution differences are not further evaluated, as we assume that they do not significantly impact the system’s quality.

Figure 4.5 presents the same statistics for the English dataset. As observed, the English knowledge source is smaller (31,659 vs. 38,642) than the German one, which is expected given that the English dataset contains only 151 ER compared to the 233 German ER. Still, the average number of passages per document is higher in the English dataset, as mentioned earlier. The total number of passages in the knowledge source indicates a highly closed-domain scenario. For comparison, MS MARCO [Nguyen et al., 2016], an often-used dataset for open-domain QA, has a knowledge source comprising 3.2 million passages, with an average length of 442 characters and passages ranging from 19 to 1,167 characters.

Figure 4.4 depicts the distribution of passage lengths in the German dataset. It’s evident that the majority of passages fall within the range of 200 to 350 characters. The same statistics for the English dataset are presented in Figure 4.6. As observed, they are quite similar. The shortest passage is 5 characters, as we applied a lower filter to the extracted passages. The longest passages are 1,300 characters long, also subjected to filtering. Considering statistics such as standard deviation, mean, and median, the English and German datasets do not significantly differ in terms of passage length.

4. Experimental Evaluation

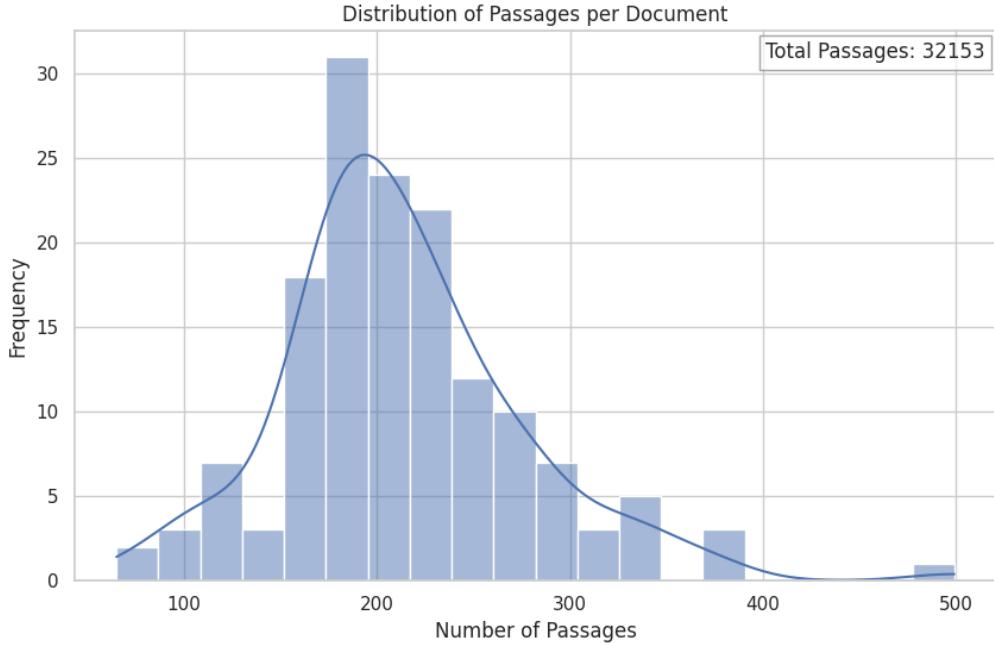


Figure 4.5.: Passages per Document of the English ER Dataset

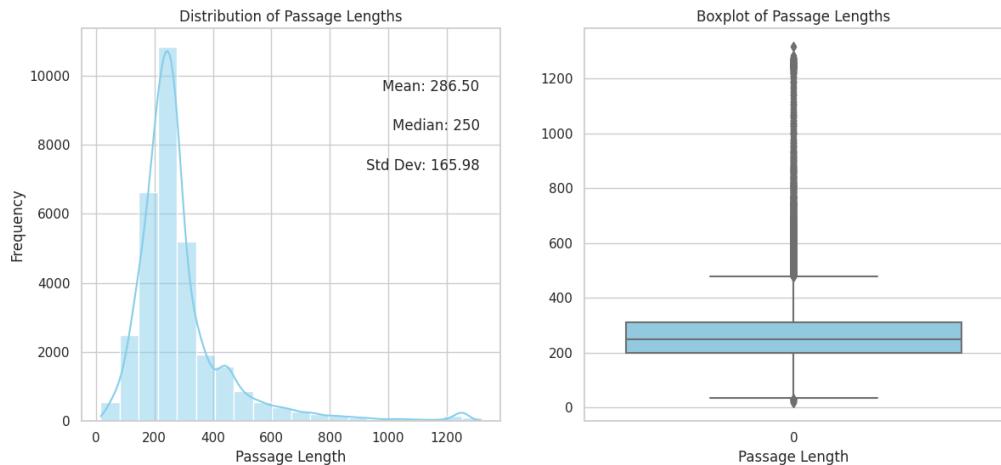


Figure 4.6.: Passage Length Distribution of the English ER Dataset

4.2. Evaluation Approach

The evaluation of a RAG system, especially in the context of Conv QA, is not trivial. Two known tools/frameworks for this task are RAGAS [Ragas, 2024] and TruLens [TruLens, 2024]. Additionally, the survey by Gao et al. [Gao et al., 2024] provides a holistic

4. Experimental Evaluation

perspective on the evaluation of RAG systems. In general, it is useful to evaluate the performance of each individual component (Retriever and Reader) separately in relation to the user's question. The outputs of the Retriever and Reader, namely the Evidence Set and the Answer, can be assessed in different aspects concerning each other and the user's question:

- **Answer Relevance:** How relevant is the Answer in relation to the Question?
- **Context Relevance:** How relevant is the Evidence Set in relation to the Question?
- **Faithfulness:** Are the facts in the Answer based on the facts of the passages in the Evidence Set?
- **Noise Robustness:** How robust is the Reader against noise in the Evidence Set and can it reflect that in the Answer?
- **Negative Rejection:** How well can the Reader determine if the Evidence Set does not contain the answer to the Question and reflect that in the Answer?

Figure 4.7 displays the composition of the specific entities and evaluation aspects.

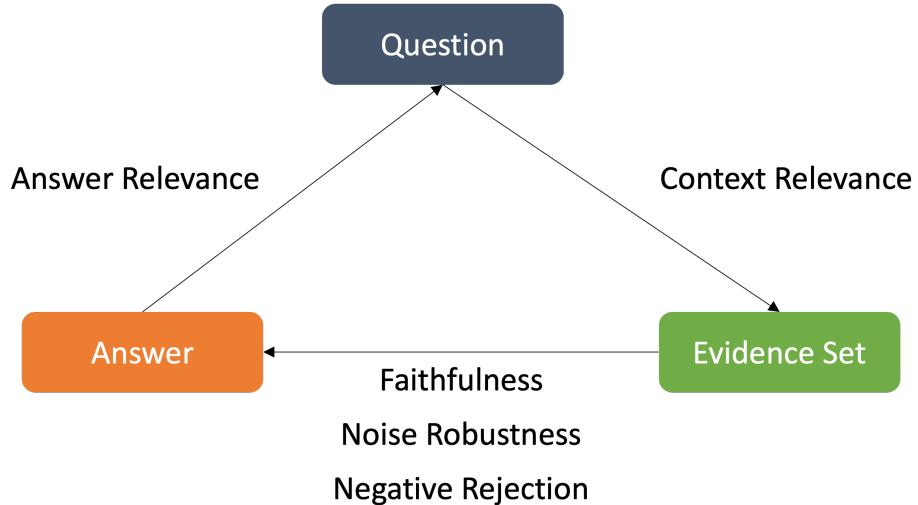


Figure 4.7.: Adapted Evaluation Aspects and Components by TruLens [TruLens, 2024]

Both the Reader and Retriever components will be evaluated offline using a synthetic dataset. This allows for the identification of the aforementioned aspects for the components in an isolated setting with only one turn, thus in a non-conversational setting.

4. Experimental Evaluation

Section 4.2.1 discusses metrics used for the Retriever evaluation, while Section 4.2.2 introduces metrics for the Reader evaluation.

To evaluate the overall end-to-end performance of the final ConRAG system, an online approach will be taken. This means that real conversations will be collected, and based on these collected data, evaluation will be conducted. Other approaches may include the generation of synthetic conversations, but this approach will not be covered in this thesis as it cannot be guaranteed with the used smaller LLMs to create high-quality conversational data. Additionally, to the best of our knowledge, there currently exists no state-of-the-art method for conversation generation based on a knowledge base. Section 4.2.3 elaborates on this end-to-end approach and the used metrics.

4.2.1. Retrieval Evaluation Metrics

Evaluating a Retriever largely depends on the use-case and the evaluation data available. Since the data introduced in Section 4.1 lacks a supervised dataset for $(question, passages)$ pairs, we will evaluate it using the synthetic dataset created, as also established in Section 4.1. This dataset consists of $(question, passages)$ pairs, where for every question, there is an exact matching passage. Therefore, this dataset is essentially a binary task, where a passage is either the correct one or not. An alternative approach would be a graded relevance task, where each passage has a certain relevance score in relation to the question. However, for our use-case, we opted for the simpler metrics HR@k and Mean Reciprocal Rank (MRR), instead of the Normalized Cumulative Gain (NDCG) used in benchmarks like BEIR [Thakur et al., 2021]. We chose these metrics because it's crucial for our system to retrieve the correct passage, and we don't have a relevance score for every passage in relation to every question.

Given a pair of (q, \hat{p}) , where \hat{p} corresponds to the correct passage and $\forall q, \exists! \hat{p} \in P$ and a retriever model $p_\eta(p|q) = \text{Score}(q, p)$ (as defined in Definition 3.7) that assigns a score to every passage $p \in P$ in relation to the question q is used. We can rank all passages based on their relevance to q . Each passage receives a rank $r_{q,p}$ based on its score in relation to q . These passages are then ranked in descending order of $r_{q,p}$, and the top k passages are added to the retrieved set R_q .

- **HR@k** This metric calculates the proportion of questions for which the correct passage is retrieved within the top k retrieved passages.

$$\text{HR}@k = \frac{1}{|Q|} \sum_{q \in Q} \begin{cases} 1 & \text{if } \hat{p}_q \in R_{q,k} \\ 0 & \text{otherwise} \end{cases} \in [0, 1] \quad (4.1)$$

4. Experimental Evaluation

The HR@k is a straightforward metric that provides a value between 0 and 1, with a higher value indicating the percentage of cases where the correct passage was retrieved within the top k passages.

- **MRR** This metric computes the mean reciprocal rank of the correct passage. It is similar to HR@k but considers the position of the correct passage \hat{p}_q within the ranking R_q .

$$\text{MRR} = \frac{1}{|Q|} \sum_{q \in Q} \frac{1}{r_{q,\hat{p}_q}} \in [0, 1] \quad (4.2)$$

In an ideal system, the MRR would be 1, indicating that the correct passage is always retrieved in the first position ($r_{q,\hat{p}} = 1$) for all $q \in Q$.

4.2.2. Reader Evaluation Metrics

Evaluating the task of answer generation, particularly the MRC aspect of the reader component, presents challenges similar to those discussed for the retrieval task evaluation in Section 4.2.1. For automatic and manual evaluation, we will utilize the synthetic dataset generated in Section 4.1. This dataset comprises triples of (question, passages, answer), where *answer* refers to a gold answer that has been syntactically generated.

In the context of a triple (q, \hat{p}, \hat{a}) , where \hat{p} corresponds to the correct passage and \hat{a} corresponds to the correct answer in relation to a question q , we employ a reader model $r_k\text{-Reader}(E_{q_c}, q_c, H) := a_i$ (as defined in Definition 3.16) to predict the answer \hat{a} given the question q and the passage \hat{p} . The predicted answer a' is then evaluated using the following metrics:

- **BLUE-1:** This precision-oriented metric compares the occurrence of unigrams (words $w \in \hat{a}$) in the predicted answer a' and the gold answer \hat{a} .

$$\text{BLUE-1} = \frac{\sum_{w \in a'} \min(\text{count}_{a'}(w), \text{count}_{\hat{a}}(w))}{\sum_{w \in a'} \text{count}_{a'}(w)} \in [0, 1] \quad (4.3)$$

Here, $\text{count}_{a'}(w)$ represents the number of occurrences of the word w in the predicted answer a' . Bilingual Evaluation Understudy (BLEU) is particularly useful for evaluating extractive questions [Papineni et al., 2002].

- **ROUGE-L:** This recall-oriented metric, especially Recall-Oriented Understudy for Gisting Evaluation (ROUGE)-L, compares the longest common subsequence

4. Experimental Evaluation

(LCS) between the predicted answer a' and the gold answer \hat{a} .

$$R_{LCS} = \frac{LCS(\hat{a}, a')}{|\hat{a}|} \quad (4.4)$$

$$P_{LCS} = \frac{LCS(\hat{a}, a')}{|a'|} \quad (4.5)$$

$$\text{ROUGE-L} = \frac{(1 + \beta^2)R_{LCS}P_{LCS}}{R_{LCS} + \beta^2P_{LCS}} \in [0, 1] \quad (4.6)$$

Here, β is a parameter to balance between precision and recall. Rouge operates similarly to BLEU but focuses on lexical matching [Lin, 2004].

- **F1-BERTscore:** BERTscore is a seq-2-seq-model-based evaluation metric for comparing two text fragments: x , which is the reference, and \hat{x} , which is the prediction. In this context, the predicted answer a' is compared to the gold answer \hat{a} . Essentially, the score between two tokens, a'_i and \hat{a}_i , is calculated as the inner product of their respective BERT embeddings: $\text{BERT}(a'_i)^T \text{BERT}(\hat{a}_i)$. For simplicity, we'll use the $\text{BERT}(a'_i) \rightarrow a'_i$ in the following equations. The final scores of F1-BERTscore are weighted by the inverse document frequency (idf) of each word-piece token:

$$P_{BERT} = \frac{\sum_{a'_j \in a'} \text{idf}(a') \max_{\hat{a}_i \in \hat{a}} (\hat{a}_i^T a'_j)}{\sum_{a'_j \in a'} \text{idf}(a')} \quad (4.7)$$

$$R_{BERT} = \frac{\sum_{\hat{a}_i \in \hat{a}} \text{idf}(\hat{a}_i) \max_{a'_j \in a'} (\hat{a}_i^T a'_j)}{\sum_{\hat{a}_i \in \hat{a}} \text{idf}(\hat{a}_i)} \quad (4.8)$$

$$F1_{BERT} = \frac{2P_{BERT}R_{BERT}}{P_{BERT} + R_{BERT}} \quad (4.9)$$

The advantage of F1-BERTscore lies in its reliance on semantic matching between the gold answer \hat{a} and the predicted answer a' rather than mere lexical matching [Zhang et al., 2020].

- **Accuracy** This metric calculates the proportion of questions for which the predicted answer a' matches the gold answer \hat{a} .

$$\text{Accuracy} = \frac{1}{|Q|} \sum_{q \in Q} \begin{cases} 1 & \text{if } a' = \hat{a} \\ 0 & \text{otherwise} \end{cases} \in [0, 1] \quad (4.10)$$

This metric is useful for evaluating question, answer realtions, where there is only

4. Experimental Evaluation

one correct answer. In order to define if an answer is correct or not, the following approaches will be used:

- **LLM-based:** A LLM can be prompted to determine a binary value (0 or 1) indicating whether the underlying message of a' and \hat{a} matches. The prompt used is based on the work of [Kamalloo et al., 2023]:

```
Question: q
Gold Answer:  $\hat{a}$ 
Predicted Answer:  $a'$ 
Is the predicted answer correct? Yes/No
```

This approach is especially useful for evaluating generative questions, as it allows for semantic matching.

- **Human-based:** In this approach, a human evaluator is asked to assign a binary value of 0 or 1, indicating whether there is a match in the underlying message of a' and \hat{a} . The evaluator assigns 1 if the answer a' covers the information from \hat{a} and 0 if it does not. The evaluator is provided with the question q , the important passage \hat{p} , the gold answer \hat{a} , and the generated answer a' . This approach is particularly useful for evaluating generative questions and closely resembles real-world applications. In addition to indicating accuracy, the evaluator gives a binary value of 0 or 1 to indicate whether the generated gold answer is correct and another binary value, if the question is an *extractive* question, which is answerable given the context. For this thesis, two evaluators will assess 50 randomly sampled datapoints with the context present in the evidence set and another 50 where the context was dropped. They will assign a binary value to each answer. The inter-rater agreement will be calculated using Cohen’s Kappa [Cohen, 1960].

4.2.3. End-to-End RAG Evaluation

The most challenging aspect of evaluation within the context of the system developed here is assessing ConRAG as a holistic system. To conduct the evaluation, a dataset of conversations needs to be generated. This will be achieved through interactions of human evaluators engaging in conversations with the system. No automatic metrics will be utilized. Instead, all evaluation aspects will be assessed by human evaluators. There exist evaluation approaches using LLMs or the embeddings of DPR in order to automatically evaluate certain aspects, but those won’t be utilized for this limited evaluation dataset. The evaluation process involves the following steps:

4. Experimental Evaluation

1. A human evaluator initiates a conversation with a question, either based on their own intuition or using provided questions from the synthetic dataset. Evaluators are encouraged to pose context-dependent questions, which can be answered given a specific passage of the POs or not at all. Complex multi passage questions should be avoided, as this complicates the evaluation process.
2. The evaluator continues the conversation, asking follow-up questions or engaging in a general discussion with the Conv QA-System for 4-10 turns.
3. After the conversation, the evaluator has access to the evidence set E_{q_e} for each turn and all other passages P in the index. The evaluator assesses the answer provided by the system and the retrieved evidence set E_{q_e} . They are required to provide the following information for every turn of the conversation:
 - *Context Relevance*: The evaluator assigns a binary accuracy value (0 or 1), indicating if the evidence set contains the required information to answer the question. A score of 1 indicates that the evidence set contains the required information, while 0 indicates that the evidence set does not contain the required information.
 - *Noise Robustness*: To have insights on the noise robustness, the evaluator has to indicate via number the position in the evidence set E_{q_e} of the relevant passage. For example, 0 indicates that the first passage in the evidence set is the relevant passage, and 4 indicates that the fifth passage in the evidence set is the relevant passage. If multiple passages are relevant, the evaluator has to indicate the position of the most relevant passage.
 - *Faithfulness*: The evaluator assigns a binary value (0 or 1), indicating whether the answer is based on the facts of the evidence set. A score of 1 indicates that the answer is based on the facts of the evidence set, while 0 indicates that the answer is not based on the facts of the evidence set. When the provided evidence set does not indicate any answer and the system provides an answer indicating that it does not have any information on the topic, the evaluator has to assign a score of 1.
 - *Negative Rejection*: This can be evaluated indirectly via the combination of Context Relevance, Faithfulness, and Answer Relevance.
 - *Answer Relevance*: The evaluator assigns a value between 0 and 1 to the system’s answer. A score of 1 indicates a correct answer to the given question in relation to the conversation history, while 0 represents an incorrect answer. A positive negative rejection is indicated by a score of 1.

4. Experimental Evaluation

- *Reason for Incorrectness:* The evaluator selects one of the following reasons for incorrectness, only if the answer was incorrect:
 - *Metadata Error:* Irrelevant passages in the evidence set could have been excluded with better metadata filtering.
 - *Wrong Crop:* The evidence set contains the correct passage, but crucial information is missing due to a flawed crop in the text corpus.
 - *Coreference Error:* The system failed to resolve a coreference problem.
 - *Not all Evidence Found:* The system failed to retrieve all necessary passages.
 - *Correct Evidence Not Found:* The system failed to retrieve the correct passage.
 - *Overlapping Evidence:* Multiple passages from different documents contain the exact same information.
 - *Hallucination:* The system provided an answer not based on the evidence set.
 - *No Negative Rejection:* The system failed to indicate insufficient information to answer the question.
 - *Nonsense Answer:* The system provided an answer with syntactical, semantic errors, or other random occurrences.
 - *Wrong Answer:* The system provided an incorrect answer given the evidence set and context.
- 4. All results will be grouped by turn, enabling deeper insights into the performance and errors of the system.

To mitigate bias, human manual evaluations are conducted by two separate evaluators. Each evaluator engages in 10 conversations across four different system settings. To ensure comparability across settings, both evaluators aim to conduct similar conversations with all four systems for all 10 interactions. Consequently, a total of 80 conversations will be conducted.

4.3. Experimental Setup and Implementation

This section outlines the experimental setup and provides details on the implementation of the System Architecture and Framework introduced in Chapter 3. First and foremost,

4. Experimental Evaluation

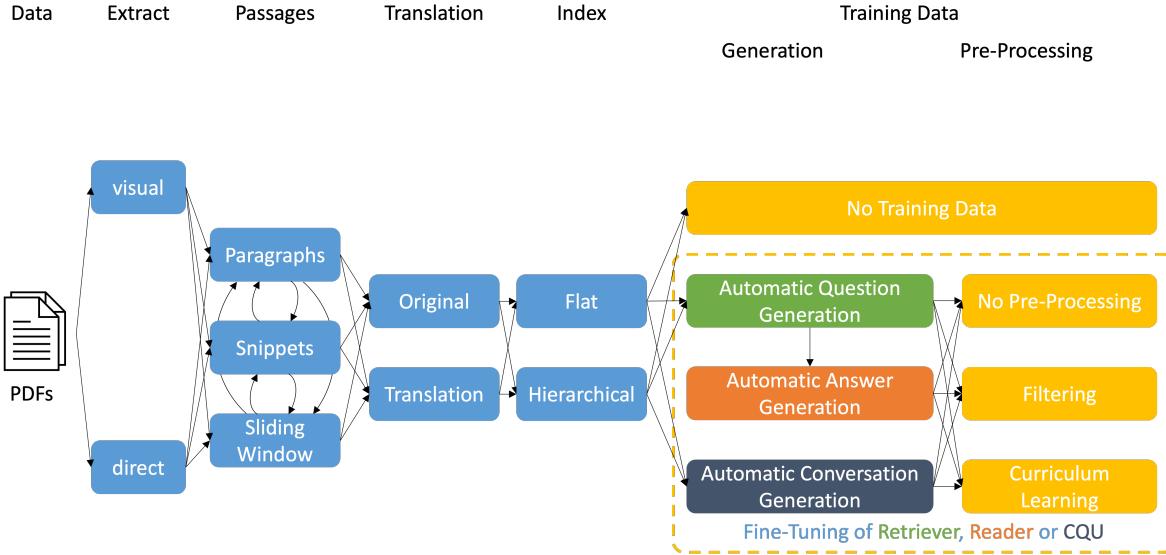


Figure 4.8.: Possible Implementations of the Extract Pipeline

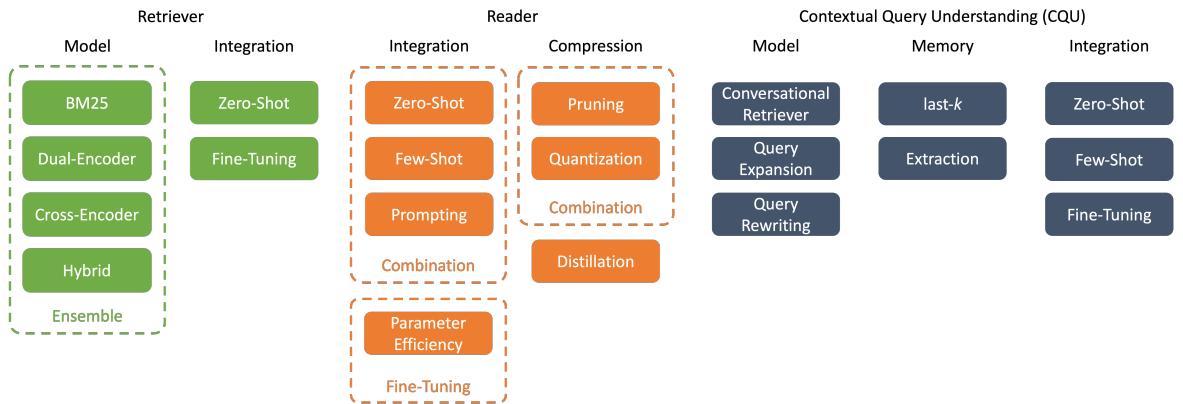


Figure 4.9.: All Components of the System Architecture

we apply the theoretically outlined model (M) to create an implementation decision map. Figure 4.8 illustrates a decision tree showcasing the possible implementations of the *Extract* pipeline. Figure 4.9 displays the main components, namely the *Retriever*, *Reader*, and *CQU*, in a similar decision matrix, where each column represents a decision to be made during implementation. Opting not to make a decision is also a valid choice. Developing an implementation can then be easily achieved using a decision tree, as shown in Figure A.1. For the implementation and benchmarking, resources of the BwUniCluster³ have been used. For further details on the model parameters used for all tasks, please refer to the Appendix A.3.1.

³<https://wiki.bwhpc.de/e/BwUniCluster2.0>

4. Experimental Evaluation

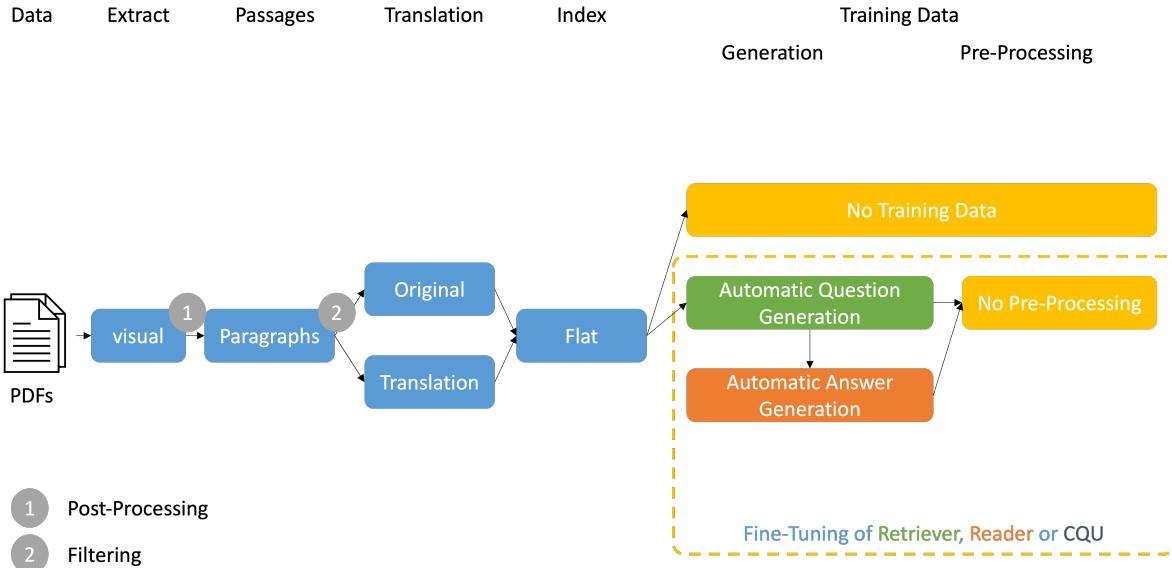


Figure 4.10.: Implemented Extraction Pipeline

4.3.1. Extraction

Figure 4.10 illustrates the extraction pipeline implemented for this thesis' PoC, resulting in the creation of the *document model* of the passages forming the knowledge source. This pipeline comprises the following steps:

1. *Extract*: Visual extraction using the Google Cloud Vision API for PDF OCR⁴, followed by post-processing.
2. *Passages*: Extraction of paragraphs using the NLTK tokenizer-based Text Splitter by Langchain⁵, followed by filtering.
3. *Translation*: Retaining the original data and providing translations from German to English and English to German using the Google Cloud Translation API⁶.
4. *Index*: A flat index, with each passage extended to include the title of the document.
5. *Training Data*: Application of Automatic Question Generation and Automatic Answer Generation using the Llama2-7b-chat and LeoLama-7b models quantized using GPTQ.

⁴<https://cloud.google.com/vision/docs/pdf>

⁵https://python.langchain.com/docs/modules/data_connection/document_transformers/text_splitters/split_by_token#nltk

⁶<https://cloud.google.com/translate/docs/overview>

4. Experimental Evaluation

Extraction Pipeline: In step (1), when given a PDF, the textual content of every page is extracted as plain text using OCR without paragraph awareness. This choice was made after simple qualitative experiments, which revealed that using direct methods leads to an unclean text corpus for the given PDFs. In the OCR, line breaks are detected and inserted as characters like \n. The textual content of separate pages is concatenated using a linebreak character \n. For post-processing, all \n characters will be replaced by spaces " ". This process results in a fully concatenated text corpus for every PDF.

In step (2), the NLTK-based Text Splitter receives a text corpus and identifies sentences in a first step based on punctuation. This list of sentences is then combined recursively to ensure it does not exceed the desired maximum length of 240 tokens. This choice of token length was made based on reference works and their chosen token lengths (see Section 2.4). It's also important to consider the input token sizes of the later-implemented Reader components. If an identified sentence itself has more than 240 tokens, e.g., 400, it will still be kept as a 400-token-long passage. Misidentifying sentences can occur quickly, e.g. text originally corresponding to a table, which cannot be easily split into sentences:

```
Coding reference Appendix 1: Semester 1 30 CPS A03-16-3 Semester  
2 30 CP Semester 3 30 CPS Key qualifications: Patient Orientation,  
Consultation, Moderation / Presentation, English, Interdisciplinary  
Collaboration 4 CP 2 CP 4 CP Scientific Writing 1 Thematic Area  
I: Scientific Principles and Methods Scientific Writing ...
```

To ensure that passages do not exceed a maximum character length, identified passages will be truncated at the next empty space after reaching 1300 characters. For comparison and to understand the impact of this filtering on the knowledge source, Figure 4.11 displays the distribution of passage lengths before filtering, compared to the distribution after filtering in Figure 4.4, which illustrates the influence of filtering. The filtering process removes outliers at 1300 characters, a decision that aligns with the later components of the system.

In Step (3), every passage and all in Step (5) generated questions of the German and English ER dataset are translated using the Google Translate API. This creates a two more dataset in addition to the existing two, which are directly based on the document's language.

The final Step (4) toward the *knowledge source* is the index creation. To simplify the later RAG system, the decision was made not to use a hierarchical index structure. However, this could lead to problems when users need to find specific information within one document and want to perform metadata-filtering. To address this, all passages have

4. Experimental Evaluation

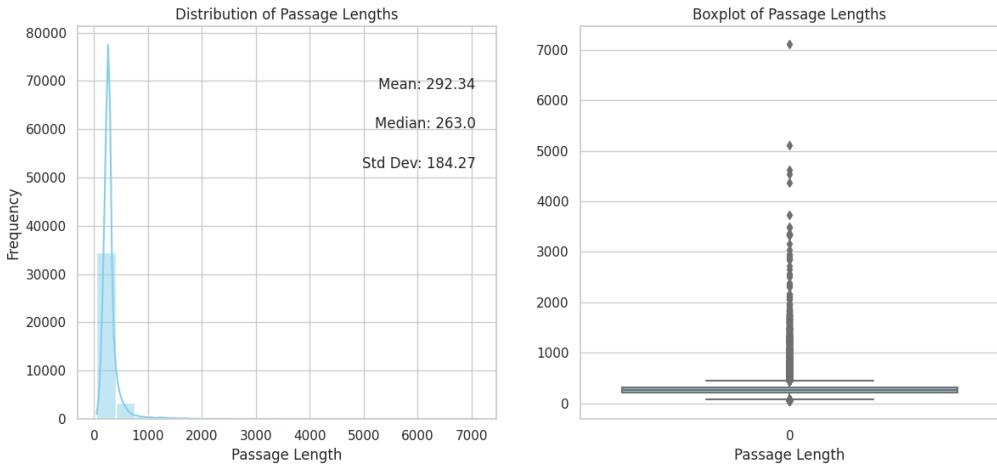


Figure 4.11.: Passage Length Distribution of the German ER Dataset before Filtering

the title of the document from which they were extracted concatenated to them in the following way:

$\{\text{passage}\} + \text{ - } + \{\text{documentName}\}$

This approach enables the retriever and reader to identify passages and their corresponding documents in a single step, as opposed to hierarchical index implementations that necessitate multiple steps and potentially multiple databases for embeddings and metadata of the passages. This method has been chosen because, in this use-case, the only utilized and easy extractable metadata of a document is a documents title (e.g., *Examination Regulation for the Master Data and Computer Science*, with some even including the date). Other use-cases may require more complicated approaches. For the translated datasets, the document name will be translated, as described in Step (3).

Data Augmentation: Step (5) of the extraction pipeline involves data augmentation. For *Question Generation* based on a given passage, the few-shot approach from *PROMPTAGATOR* [Dai et al., 2022b] has been employed. For the English ER dataset, the following prompt is used to generate a question, given a passage and i examples:

System Prompt:

You are an assistant generating one question given a context. Please start your generated question with the indicator 'Generated Question:'

Here are some examples:

for example in examples:

{i}-Example:

4. Experimental Evaluation

Context: {passage}

Generated Question: {question}

User:

Generate the question based on the following context:

{passage}

The snippet in the System Prompt is looped four times to accommodate the number of examples. The examples represent different question types, which themself indicate different answer types, in order to add some variety to the generated data. The used examples can be found in the appendix A.2.1. The Prompt is generally designed to be used with a chat fine-tuned model, as it showed better results in a qualitative evaluation. For English question generation, the *Llama2-7B-Chat-GPTQ*⁷ is employed. It's a quantized version of the original Llama2-7B-Chat Model by Meta [Touvron et al., 2023]. For the German ER, the same prompt translated to English is used. For the German tasks, the *Leo-Hessianai-7B-Chat-GPTQ*⁸ model is used, which is a quantized version of a German language fine-tuned Llama2-7B-Chat Model [Plüster, 2023].

It's worth noting that the English Llama2-7B-Chat was capable of directly understanding the system prompt and generating questions with the indicator **Generated Question**, which made using the generated text easier. However, the German LLM struggled with this Few-Shot task and often generated multiple responses that didn't start with the indicator **Generierte Frage**. Additionally, the model sometimes included answers after the questions. To extract the single question from the generated text in the German LLM output, the text is cropped after a : and ?, so the string enclosed by these characters is used as the question for a given passage.

The choice of Llama2-based models was based on their new state-of-the-art results in multiple benchmarks, as demonstrated by these models in the open-source LLM niche [Touvron et al., 2023]. Additionally, the work of Plüster et al. [Plüster, 2023] on LeoLama opens the opportunity to compare the performance of German and English LLMs from the same family.

For *Answer Generation*, the previously generated questions per passage are utilized. The LLM is prompted in the following manner:

System Prompt:

You are an assistant who generates gold answers for a given question and context.

⁷<https://huggingface.co/TheBloke/Llama-2-7b-Chat-GPTQ>

⁸<https://huggingface.co/TheBloke/leo-hessianai-7B-chat-GPTQ>

4. Experimental Evaluation

User:

Generate the answer based on the following question and context:

Question: {question}

Context: {passage}

This prompt was applied in both German and English to generate gold answers for all datasets. A random selection of 2000 samples was drawn from the original German and English ER datasets and their translated versions, resulting in four selections of 2000 question-context pairs each⁹. On this smaller selection, *gpt-3.5-turbo* was employed with the mentioned prompt to generate high-quality question-context-answer triples. Additionally, *leo-hessianai-7B-chat-GPTQ* and *Llama2-7B-Chat-GPTQ* (See Section 4.3.3) were used to generate gold answers over the entire dataset. The sampled high-quality gold-answer sets will be employed in Section 4.4.3 to compare the gold-answer generation quality of the smaller models against the larger *gpt-3.5-turbo* model. In order to extend the triples to contain not only the correct context, but also other high probable contexts per question, the large DPR retriever (See Section 4.3.2) was used per triple to retrieve the top 10 passages. Those were appended to the context of the triple. This will be used in the reader evaluation later.

4.3.2. Retriever

The implemented retrievers can be found in Figure 4.12. Given the limited availability of (synthetic) data suitable for fine-tuning, the decision was made to exclusively utilize top-performing OOD zero-shot retrievers. Fine-tuning a model on such a small dataset could lead to a high risk of overfitting. To assess retriever performance, the synthetic dataset will be used for benchmarking. Among the top-performing retrievers from the BEIR benchmarks, the following three retrievers have been chosen:

1. **BM25:** This employs the standard lexical-based best match algorithm with hyperparameters $k_1 = 1.5$, $b = 0.75$, and $\epsilon = 0.25$. It is implemented using the open-source project *rank-bm25*¹⁰.
2. **BM25 + CE:** This uses a BM25 retriever in combination with a Cross-Encoder re-ranker, which serves as the baseline established in BEIR [Thakur et al., 2021]. This combination continues to perform as the state-of-the-art. The Cross-Encoder

⁹This sampling approach will also be employed later to mitigate the costs associated with the OpenAI API

¹⁰https://github.com/dorianbrown/rank_bm25

4. Experimental Evaluation

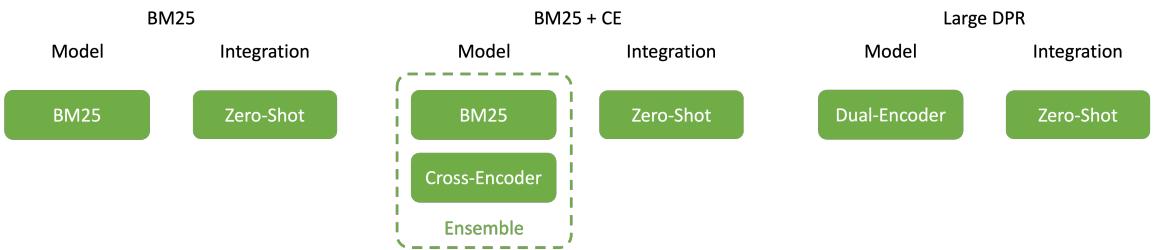


Figure 4.12.: Implemented Retrievers

utilized is the same as the one in the BEIR baseline implementation: *ms-marco-MiniLM-L-6-v2* [Wang et al., 2020], specifically the implementation provided on Hugging Face ¹¹. It has close to 17 million parameters. For German indices, a German version of *ms-marco-MiniLM-L-6* is used: *ms-marco-MiniLM-L-6-en-de-v1* ¹².

3. **Large DPR:** For the large DPR, the embedding model from OpenAI, *text-embedding-ada-002* ¹³, is employed. The parameter size of this model is estimated to be around 350 million parameters[Muennighoff, 2022]¹⁴.

These retrievers have been implemented and evaluated based on the synthetic data, and the results are presented in Section 4.4.

4.3.3. Reader

The implemented readers can be found in Figure 4.13. Due to the limited hardware resources and high-quality data, the decision was made to only use zero-shot pre-trained LLMss for the reader component and not fine-tune any. The following readers have been chosen:

1. **gpt-3.5-turbo:** GPT 3.5 Turbo is part of the OpenAI chat completion family¹⁵. Details on the parameters are not available. However, the latest published research by OpenAI indicates 175B parameters for the model on which GPT 3.5 Turbo was built. In general, the model is an LLM trained using human feedback on chat-like completion tasks [Ouyang et al., 2022]. It accepts up to 4,096 tokens as input.
2. **Llama2-7B-Chat-GPTQ:** This is a quantized version of the original Llama2-7B-Chat Model by Meta [Touvron et al., 2023]. The model is an LLM trained on

¹¹<https://huggingface.co/cross-encoder/ms-marco-MiniLM-L-6-v2>

¹²<https://huggingface.co/cross-encoder/msmarco-MiniLM-L6-en-de-v1>

¹³<https://platform.openai.com/docs/models/embeddings>

¹⁴As OpenAI does not publicly provide information on their model specs, there exist only estimates.

¹⁵<https://platform.openai.com/docs/models/gpt-3-5>

4. Experimental Evaluation

chat-instruction datasets [Ouyang et al., 2022]. The model is available on Hugging Face¹⁶. The model is quantized using GPTQ [Muennighoff, 2022]. However, the provider did not report any benchmarks indicating the accuracy drop of the GPTQ version compared to the original. It accepts up to 4,096 tokens as input.

3. **leo-hessianai-7B-chat-GPTQ:** *Linguistisch Erweitertes Offenes Language Model* (LeoLM) is a German tasks fine-tuned version of the original Llama2-7B-Chat model. The original Llama2 is fine-tuned on a large corpus of German language texts, and multiple adjustments are applied to overcome the problem of forgetting. This results in a maximum input token length of 8,000 tokens. In addition to original German chat instruction-based datasets, automatically translated English benchmark datasets are used as well. Benchmarking results indicate an increase in German language capabilities while maintaining English language task performance [Plüster, 2023].

Generally, **gpt-3.5-turbo** acts as a higher-end benchmark. It is a state-of-the-art LLM with high popularity in the media and developer community. **Llama2-7B-Chat-GPTQ** should be a hardware resource-efficient alternative pre-trained model, especially the quantized version that can run on consumer hardware. **leo-hessianai-7B-chat-GPTQ** is representative of a German native model, providing further insights into language dependencies and the performance of LLMs, as the use case involves multiple indices in German and English.

Additional approaches, such as evaluating different fine-tuning strategies, have been excluded, as they would open up a vast problem space that cannot be adequately covered within the scope of this thesis. Therefore, the thesis predominantly concentrates on a comprehensive PoC to demonstrate and evaluate one possible implementation approach of the framework introduced in Section 3, utilizing the mentioned LLMs.

However, to address specific questions in the evaluation (See Section 4.4.3), such as *How many contexts can the model receive and still determine the correct answer?*, the reader prompt was structured to accept multiple contexts as inputs to the question:

System Prompt:

You are an assistant helping a student with examination regulations of the Heidelberg University. Generate a concise answer based on the provided contexts. If the answer isn't in the contexts, state 'I can't answer this question'. If there's ambiguity, ask clarifying questions.

¹⁶<https://huggingface.co/TheBloke/Llama-2-7b-Chat-GPTQ>

4. Experimental Evaluation

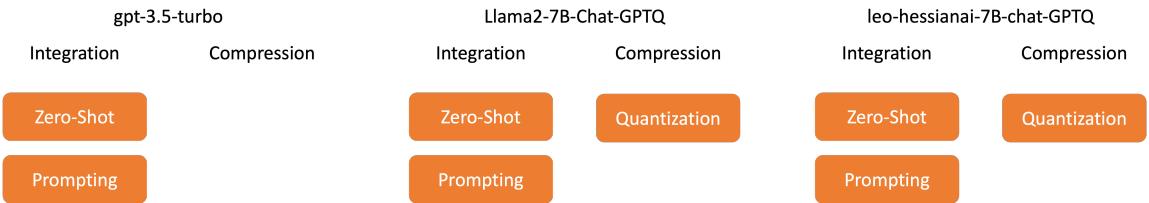


Figure 4.13.: Implemented Readers

{previous chat history}

User:

Please generate an answer for the following question:

{question}

Use the following contexts:

for context in contexts:

{context}

Assistant:

This prompt is applied consistently across all models. In cases where the model is applied on a German dataset, a German version of this prompt is utilized.

4.3.4. CQU

Similar to the Reader and Retriever components, the implementation of the CQU component can be observed in Figure 4.14. Generally, the CQU component utilizes a LLM for query rewriting based on the last- k turns of the conversation history. k is thereby not a fixed value, but rather the number of tokens for every turn will be calculated and only the last- k turns until which the prompt does not overstep the token limit of the model, will be used. This follows the *Rewriting Prompt* implementation by Mao et al. [Mao et al., 2023]. Their work demonstrated that LLMs can achieve state-of-the-art benchmarks for CQU. Mao et al. also explored other approaches to prompting beyond simple rewriting, which, while outperforming simple rewrite, are only applicable to DPR-based retrievers.



Figure 4.14.: Implemented CQU

4. Experimental Evaluation

The simple rewrite prompt is implemented as follows:

System Prompt:

As an assistant, you're tasked to decontextualize and reformulate the current question considering a multi-turn information-seeking dialogue. Please provide only the rewritten question after 'Rewrite:'. It will be used for passage retrieval. Here are some examples: for example in example:

{example}-Example:

for turn in turns:

{turn}-Turn:

Question:

{question}

Answer:

{answer}

Question:

{question}

Rewrite:

{rewritten last user question}

{previous chat history}

User:

Please rewrite the following question given this conversation history:

Question:

{question}

Rewrite:

In terms of models, the same models as for the Reader component will be utilized: *Llama2-7B-Chat-GPTQ* for English conversations and *leo-hessianai-7B-chat-GPTQ* for German conversations. The CQU component itself will not be evaluated since there is no dataset available. For evaluations of this approach on other datasets and LLMs, the reader is referred to the work of Mao et al. [Mao et al., 2023].

4.4. Experimental Results

This section evaluates the implemented system and the introduced framework, using metrics introduced in Section 4.2. The evaluation covers the following components: *Extraction*, especially synthetic data generation (Section 4.4.1), *Retriever* (Section 4.4.2),

4. Experimental Evaluation

Reader (Section 4.4.3), and the complete Conv QA system (Section 4.4.4). The evaluation is based on the synthetic dataset described in Section 4.3.1, generated from the dataset introduced in Section 4.1.

4.4.1. Evaluation Limitations

The outlined limitations primarily pertain to the Reader Evaluation. The Retriever Evaluation should be minimally affected, mainly by the usefulness of the generated questions. This section serves as a guide for future work, particularly when dealing with synthetic data generation. For more robust insights into RAG and the reader component in a conversational scenario, refer to Section 4.4.4, which is based on real-world data rather than synthetically generated data.

LLM Performance: As detailed in Section 4.3.1, the two main indices, namely *indexEnglish* and *indexGerman*, are based on the corresponding ERs. During the generation of the question-context datasets, qualitative issues arose when using *leo-hessianai-7B-chat-GPTQ*. As indicated by System Prompt 4.3.1, the model was prompted to start the generated answer after a defined indication text. In multiple cases, the German model *leo-hessianai-7B-chat-GPTQ* was unable to fulfill this request, necessitating post-processing as described in Section 4.3.1. This issue was not present when using the English model *Llama2-7B-Chat-GPTQ*. This observation highlights a zero-shot performance gap between the English model *Llama2-7B-Chat-GPTQ* and the German model *leo-hessianai-7B-chat-GPTQ*. Consequently, the subsequent evaluations will concentrate on the *indexEnglish* and *translated-indexEnglish* datasets.

Gold-Answer-Generation Performance: The *Human-based* reader’s accuracy evaluation (Table 4.1) shows a low rate of correct gold answers, approximately 22.25%. The main reason for incorrectness is hallucination by smaller LLMs, leading to mostly incorrect gold answers. This impacts the performance of different evaluation metrics for the reader. The *F1-BERTscore* operates only on the gold answer \hat{a} and the predicted answer a' , potentially lacking sufficient insights. In contrast, the *LLM-based* accuracy considers a triplet of question q , gold answer \hat{a} , and predicted answer a' , making it more robust to hallucination within gold answers. Across all models and indices, there is a gap of approximately 20.4% between *Human-based* and *LLM-based* results. The *short* versions of the indices, which use gold answers generated by *gpt-3.5-turbo* with less hallucination, show a gap of 9.11% when filtered for only useful gold answers. Therefore, the evaluation will focus on the *short* versions of the indices for the reader.

4. Experimental Evaluation

Generated Question Usefulness: As observed in Section 4.3.1, the approach for synthetic question generation aims to create various question types, including extractive and conceptual questions. However, the *Human-based* Evaluation reveals that only 51.25% of the generated questions are extractive, meaning they can be answered with information from the passage they are based on. This poses a significant challenge for automatically evaluating the reader and retriever components, as only a portion of the question-context tuples can be assessed in such constellations. Others may not make sense as the desired information might not be in the answer or even provided in the knowledge base. Despite this, in Figure 4.15d, retrievers show diminishing returns in performance improvement at approximately $k = 300$, reaching HRs between 0.71 and 0.96. However, for future work, the entire evaluation pipeline using synthetic data must be reconsidered, especially for evaluating complex question types.

Negative Rejection Evaluation: The chosen approach to automatically evaluate negative rejection, as described in Section 4.4.3, faces two major challenges. Firstly, multiple passages are identical across the knowledge base, rendering the removal of the gold passage from the evidence set redundant. Secondly, as discussed in the **Generated Question Usefulness** section, the removal of the gold passage may still leave sufficient information in the evidence set to answer the question. This does not create an ideal scenario for evaluating negative rejection capabilities. A more suitable approach would be to manually generate a set of questions that are unanswerable by the knowledge base and then evaluate the models on this set. However, this approach may not be feasible with synthetic data.

BLUE and ROUGE Metrics: The applicability of *BLUE* and *ROUGE* metrics in the context of generative readers for question-answering needs to be questioned. Applying these metrics to evaluate the reader in the given setup leads to random results. This is because these metrics are designed to determine a vocabulary-based overlap between two texts, which is only applicable in the case of extractive readers, not generative readers as used in RAG.

4.4.2. Retrieval Results

The evaluation based on the described metrics (See Section 4.2) delivers multiple insights:

1. Overall, the out-of-domain zero-shot performance of the large DPR outperforms

4. Experimental Evaluation

BM25 for all k -values, indices and languages. This can be observed in Figure 4.15d.

2. The ensemble of *BM25-CE* outperforms the regular *BM25* retriever for three out of four indices (all except for the *indexGerman*) when considering small values of k . Figure 4.15a compares the MRR between the standard *BM25* retriever and *BM25-CE* retriever with the BM25 k -value set to 500. As the *Reader* component ideally takes as small a set of passages as possible as input, the *BM25-CE* retriever is the better choice. The internal k -value of the first stage *BM25* retriever can be as large as possible, as it only improves the performance of the Re-Ranker. However, depending on the dataset and use-case, an increasing k -value may not necessarily improve performance, as it seems to exist an upper performance limit in the cross-encoder component.
3. Figure 4.15b compares different *BM25-CE* Re-Ranker settings over the *indexEnglish* index. It can be observed, that an increase in the first-stage Retriever's k -value does at some point not longer increase the performance of the Re-Ranker. This is true for the MRR and HR.
4. In the comparison of *BM25* and *DPR* retrievers for different languages, the final choice of the *Retriever* component should consider the best MRR value for English and German datasets for $k \leq 10$. Figure 4.15c compares *BM25-CE-500* with the large *DPR* retriever over all datasets. The DPR retriever significantly outperforms the *BM25-CE-500* retriever for all datasets, indicating a high multi-language zero-shot performance of the *text-embedding-ada-002* model by OpenAI. Nevertheless, considering the difference in parameters (350 million vs. 17 million), the *BM25-CE-500* retriever is a good alternative.

The influence of the Retriever component will be further evaluated in Section 4.4.4. Here the bottlenecks of the system will be identified and discussed. A full set of evaluation results can be found in the tables in Appednix A.3.2.

4.4.3. Reader Results

The evaluation of the Reader is based on the metrics introduced in Section 4.2. Some adjustments have been made as follows: Firstly, whenever *gpt-3.5-turbo* is evaluated on a dataset, it is assessed on a random sample of 1,000 elements based on that dataset. Secondly, when calculating the LLM-based accuracy, it is computed on a random sample set of 1000 elements of a model's predictions. These adjustments aim to mitigate

4. Experimental Evaluation

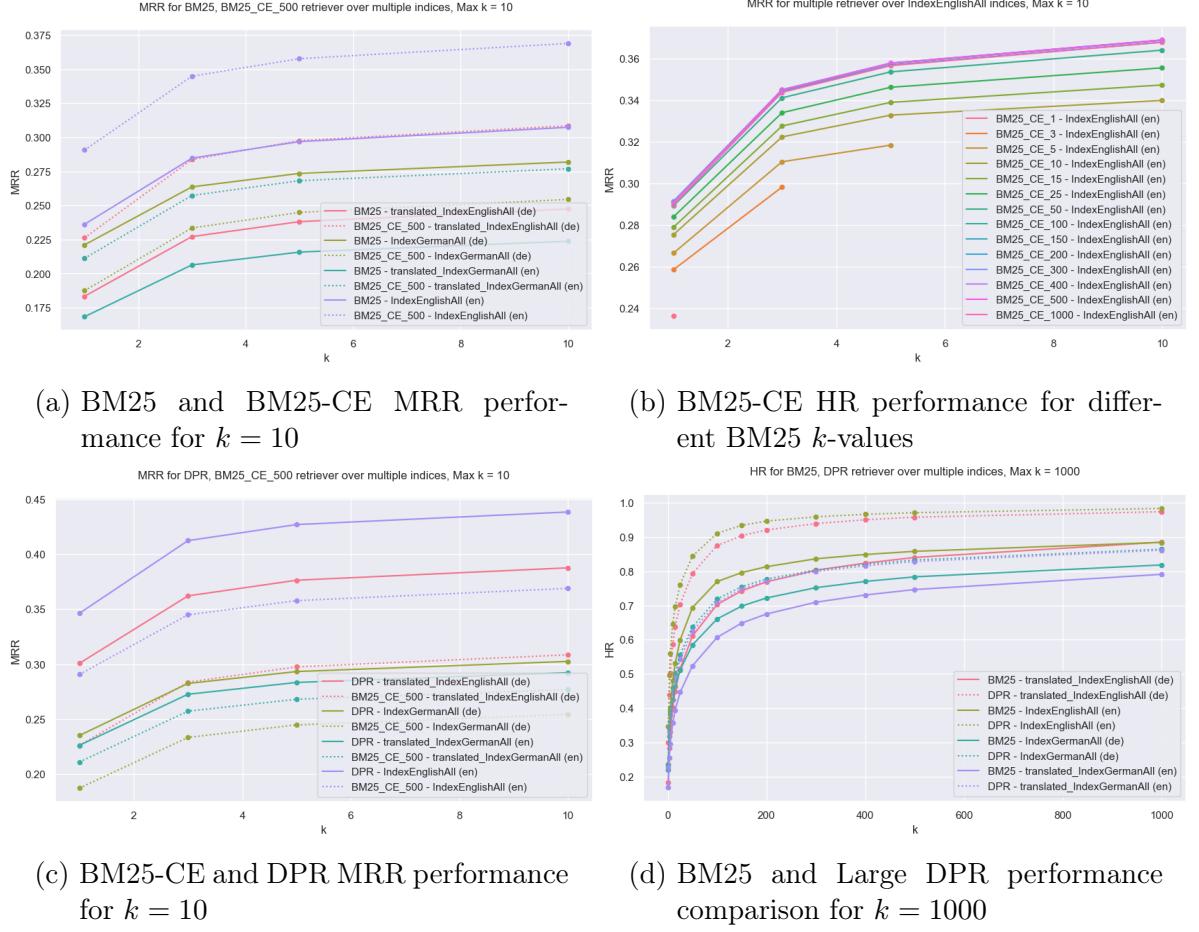


Figure 4.15.: Multiple Retriever Performance Comparisons

the costs associated with the OpenAI API. Thirdly, two separate datasets have been generated based on the four already introduced indices: *indexEnglish*, *indexGerman*, *translated indexEnglish*, and *translated indexGerman*. The *short* version of the datasets corresponds to 2,000 random samples, where *gpt-3.5-turbo* is used to generate the gold-answers in the data augmentation stage. The *filtered* version of the datasets drops question-context pairs if the gold-passage is not included in the evidence set of the DPR retriever with $k = 1,000$ and generates the gold-answers using *Llama2-7B-Chat-GPTQ* for *indexEnglish* and *translated indexGerman*, and *leo-hessianai-7B-chat-GPTQ* for *indexGerman* and *translated indexEnglish*. In order to test the *Negative Rejection*, with a randomness of 20% the correct passage was removed from the evidence set of the DPR retriever. Those evaluation datasets will be referred to as *dropped* datasets.

Table 4.1 presents the evaluation results for the reader component based on filtered versions of the datasets. Specifically, the *gpt-3.5-turbo* model is assessed on the *indexEnglish* and *translated indexEnglish* datasets, while the *leo-hessianai-7B-chat-GPTQ* and *Llama2-7B-Chat-GPTQ* models are evaluated on the *indexGerman* and *translated*

4. Experimental Evaluation

indexGerman datasets. The Cohen’s Kappa coefficient indicates a moderate agreement between the two evaluators, scoring 0.72.

Index	<i>gpt-3.5-turbo</i>		<i>leo-hessianai-7B-chat-GPTQ</i>		<i>Llama2-7B-Chat-GPTQ</i>	
	accuracy	accuracy qat*	accuracy	accuracy qat*	accuracy	accuracy qat*
IndexEnglish	all	0.68 (34.5/10)	-	0.2367 (13/10)	0.24	0.3571 (16/14)
	dropped	0.56 (40)			0.52	0.5334 (30)
translated IndexEnglish	all	0.62 (39/10.5)	0.2	0.2367 (13/10)	-	-
	dropped	0.58 (44)	0.44	0.6786 (28.5)		

* *qat* stands for *question and answer true*, meaning only those datapoints are considered where the question is extractive and the gold answer is correct.
(x/y) Number of *extractive questions* and *correct gold answer* in parentheses.

Table 4.1.: Human-based Reader Evaluation for context length of 3 and no shuffled contexts

The following insights can be drawn from the evaluation results, keeping in mind the evaluation limitations mentioned in Section 4.4.1:

1. In Figure 4.16a and Figure 4.16b, the accuracy of the *gpt-3.5-turbo* model increases as the size of the evidence set (k) grows. Conversely, the performance of smaller LLMs like *leo-hessianai-7B-chat-GPTQ* and *Llama2-7B-Chat-GPTQ* declines with larger k values, possibly due to their limited capacity to process extensive evidence sets.
2. Figure 4.16a compares *LLM-based* accuracy across both *indexEnglish* and *translated indexEnglish* without shuffled contexts, while Figure 4.16b compares *F1-BERTScore* for the same. Both metrics exhibit a similar trend, although *F1-BERTScore* shows a notable disparity across languages.
3. Figures 4.16d and 4.16d illustrate that context order has minimal impact for a maximum context length of 10, irrespective of language.
4. Regarding *Negative Rejection*, the *Human-based* evaluation suggests a performance of 0.675 for *gpt-3.5-turbo* and 0.5334 for *Llama2-7B-Chat-GPTQ* on the *indexEnglish* dataset, and 0.6364 for *gpt-3.5-turbo* and 0.6786 for *leo-hessianai-7B-chat-GPTQ* on the translated dataset. Due to limitations, automated evaluations will be excluded. This indicates a relatively moderate *Negative Rejection* capability.
5. Overall, *gpt-3.5-turbo* outperforms smaller models consistently across languages and settings. For instance, at $k = 10$ without shuffled contexts, the gap between *gpt-3.5-turbo* and *Llama2-7B-Chat-GPTQ* is 0.63 for *indexEnglish* and 0.82 for *translated IndexEnglish*.

4. Experimental Evaluation

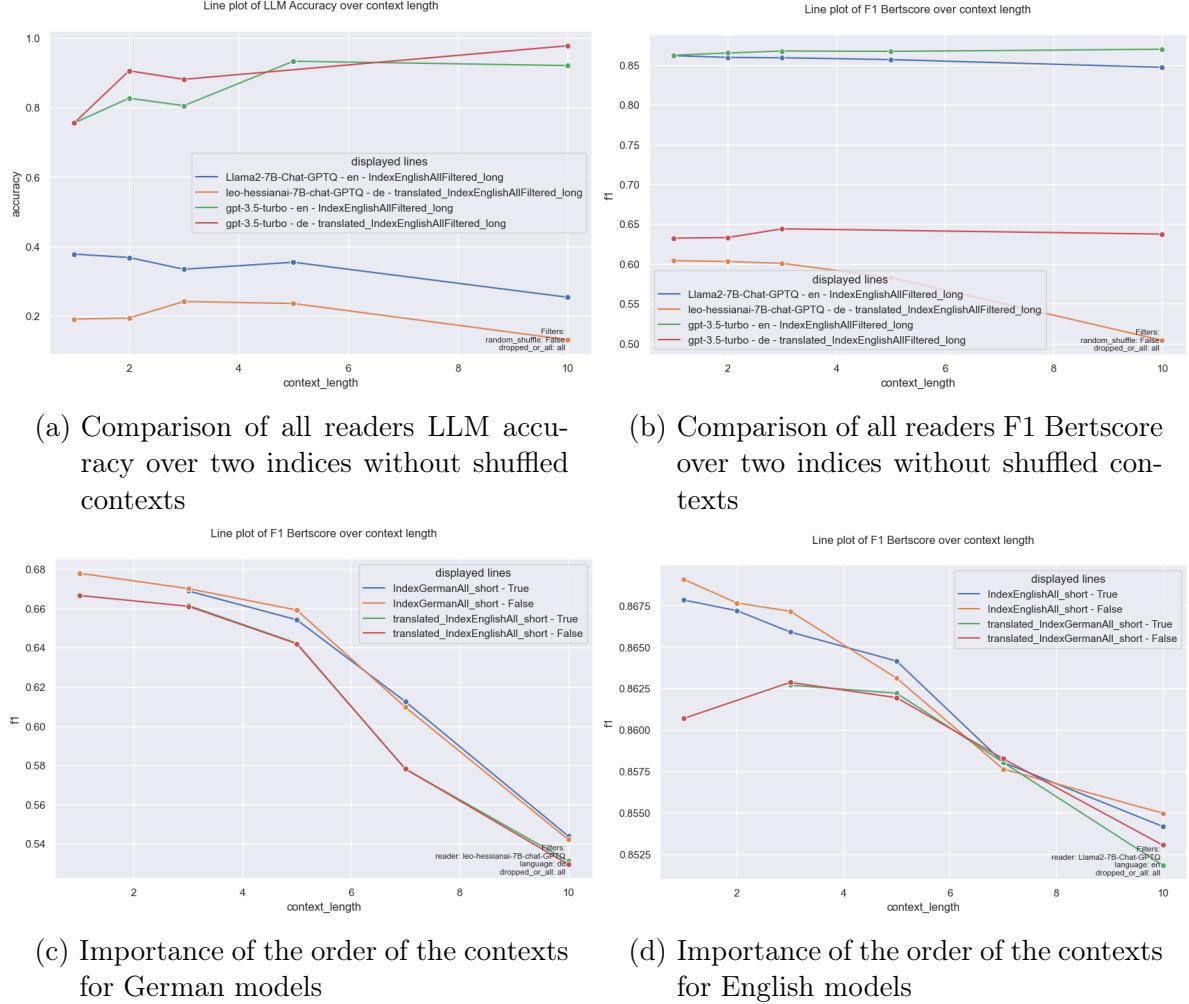


Figure 4.16.: Multiple Reader Performance Comparisons

- Language gap analysis indicates no significant performance difference for *gpt-3.5-turbo* across languages in this setting. However, the smaller models *leo-hessianai-7B-chat-GPTQ* and *Llama2-7B-Chat-GPTQ* show a gap of approximately 0.1 accuracy between English and German models, consistent across automatic and *Human-based* evaluations.

As discussed in Section 4.4.1, we strongly recommend readers to refer to the evaluation results presented in the next Section 4.4.4. This is because the evaluation of the reader component itself, based on synthetic data, is severely constrained.

4.4.4. Conversational Question Answering Results

The evaluation procedure for the Conv QA system is detailed in Section 4.2.3. Four system settings are assessed: *gpt-3.5-turbo-de*, *gpt-3.5-turbo-en*, *leo-hessian-7B-chat-de*,

4. Experimental Evaluation

and *Llama2-7B-chat-en*. The suffixes *de* and *en* denote the language of the prompts (German and English, respectively), and indicate the corresponding index used (*IndexGerman* for German and *translated IndexGerman* for English settings). For precise model inference parameters, please consult Table A.1 in the appendix.

Question Type	GPT-3.5-de		GPT-3.5-en		Llama2-7B-chat-en		leo-hessian-7B-chat-de	
	avg. AR*	Count	avg. AR*	Count	avg. AR*	Count	avg. AR*	Count
Abstract	0.478	18	0.421	14	0.200	31	0.258	12
Complex	0.327	11	0.263	8	0.000	5	0.167	3
Factoid	0.435	31	0.477	26	0.323	40	0.154	52
List	0.550	4	0.467	6	0.600	2	0.225	4
Yes/No	0.500	2	0.200	2	0.067	3	0.420	5

* AR stands for Answer Relevance

Table 4.2.: Answer Relevance per Model, Language and Question Type

Answer Relevance: Table 4.2 provides insights into answer relevance based on the number of datapoints per question type. Qualitative analysis is feasible for abstract and factoid questions, as they have a sufficient count. Notably, *gpt-3.5-turbo-de* and *gpt-3.5-turbo-en* outperform the other models, exhibiting a 0.43 and 0.375 improvement, respectively, in German and English settings. However, *leo-hessian-7B-chat* shows better performance on complex questions compared to simple factoid questions, while *Llama2-7B-chat-en* excels in factoid questions. Given the data in Table 4.3, it's clear that the larger model consistently outperforms the smaller models, achieving an average of 12.25% more error-free turns. However, even the best-performing model, *gpt-3.5-turbo-de*, achieves only 20.59% error-free turns, indicating room for improvement across all models. When considering the average answer relevance in error-free turns, all models demonstrate similar performance. Additionally, Figure 4.18 underscores that smaller models significantly underperform compared to *gpt-3.5-turbo*. Notably, in system settings with smaller models, almost half of the errors stem from weak LLM performance, contrasting with *gpt-3.5-turbo*, where model-related errors comprise only around 18%

Model-Language	Number of Error-free turns	avg. Answer Relevance
GPT-3.5-de	14 (20.59 %)*	0.7786
GPT-3.5-en	10 (17.54 %)*	0.85
Llama2-7B-chat-en	3 (3.37 %)*	0.7333
leo-hessian-7B-chat-de	8 (10.26 %)*	0.825

* The percentage compared to all turns

Table 4.3.: Error-free turns per Model and Language

4. Experimental Evaluation

of total errors. It's noteworthy that the *Nonsense Answer* error occurs only for the smaller models, indicating weaker consistent text generation performance compared to larger models. Additionally, comparing the total occurrences of *Wrong Answer* errors, it's evident that the smaller models have nearly three times as many occurrences of this error type in total numbers as the larger LLM.

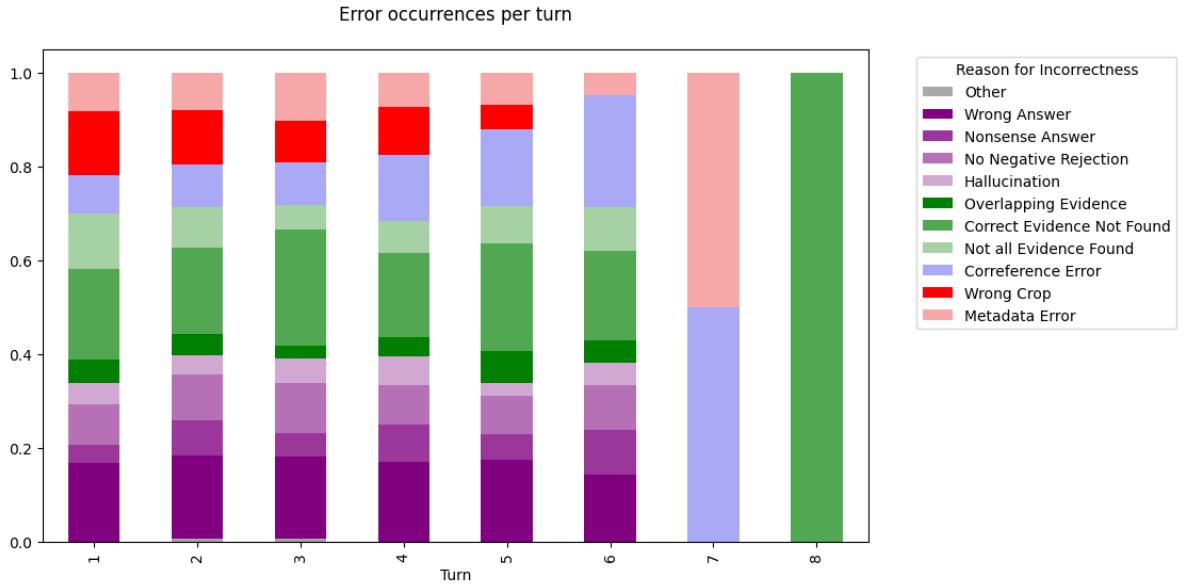


Figure 4.17.: Error Type Distribution per Turn

Context Relevance: The system exhibits a significant weakness in the Retriever component combined with the Extraction pipeline. Figure 4.18 illustrates that approximately 18.5% of errors in each system setting stem directly from Extraction issues, such as incorrect cropping or metadata filtering. Moreover, a major drawback, accounting for around 46.5% of errors in settings with the robust LLM *gpt-3.5-turbo*, is evidence retrieval performance. This can be attributed to Extraction or the Retriever component itself. Table 4.4 further underscores these challenges, indicating that the Retriever found relevant passages in only 49% of *Factoid* questions and a mere 28% in *Complex* questions. Although the exact bottleneck is unclear from this evaluation, the Retrieval evaluation in Section 4.4.2 suggests improvements are needed in the Retriever Component. Additionally, approximately 11.6% of errors result from coreference errors in the Conversational QA component, particularly notable in longer conversations as shown in Figure 4.17.

Faithfulness: Regarding faithfulness and avoidance of hallucination, Table 4.5 illustrates that *gpt-3.5-turbo* exhibits hallucination in approximately 10% of cases. In contrast, *Llama2-7B-chat* and *leo-hessian-7B-chat* demonstrate higher rates, with around 27% and 38% respectively. This significant drawback of the smaller models results in a

4. Experimental Evaluation

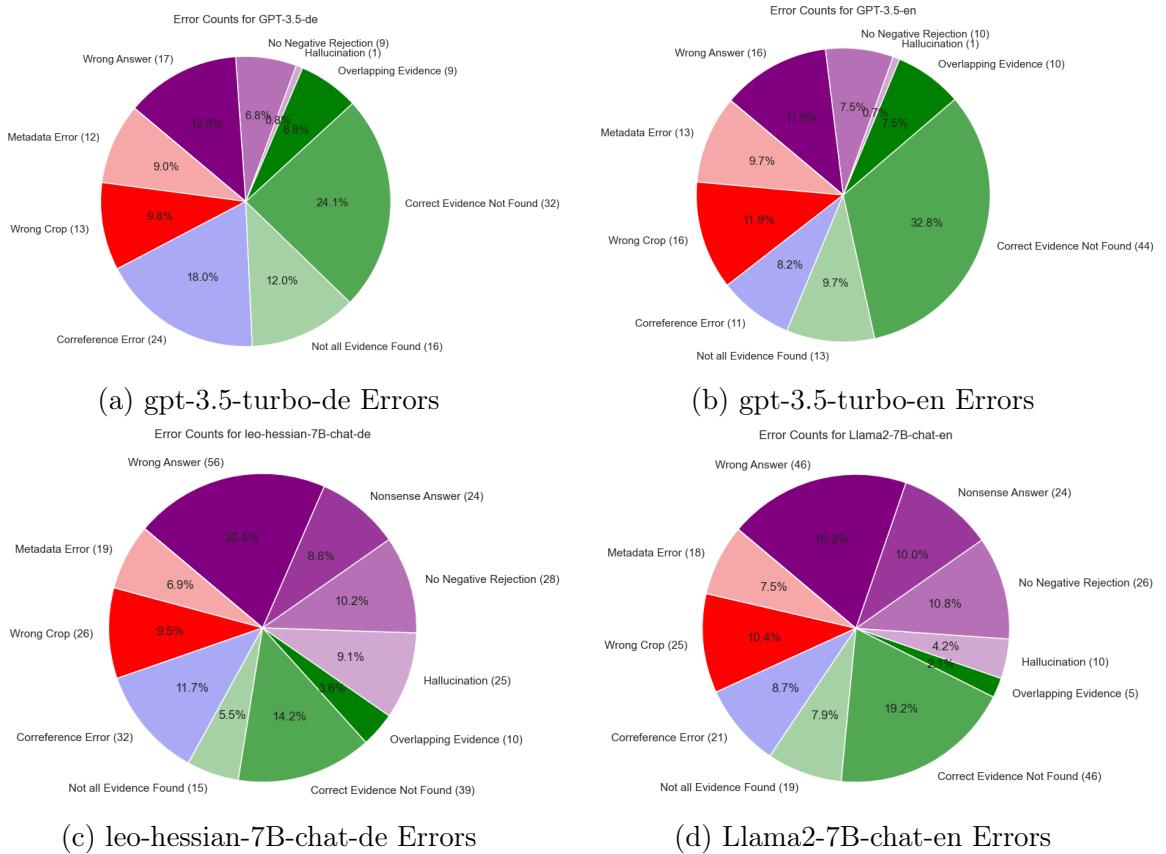


Figure 4.18.: Error Distribution per Model

notable proportion of incorrect answers, which is particularly undesirable in applications where factual accuracy is crucial.

Question Type	avg. Context Relevance	Count
Abstract	0.404040	99
Complex	0.282051	39
Factoid	0.492308	195
List	0.833333	18
Yes/No	0.428571	14

Table 4.4.: Average Context Relevance per Question Type

4. Experimental Evaluation

Model Language	avg. Faithfulness	Count	avg. AR1*	Count	avg. AR>1**	Count	Negative Rejection	Specificity
GPT-3.5-de	0.892	93	0.692	26	0.654	13		0.607
GPT-3.5-en	0.908	98	0.550	20	0.613	8		0.826
Llama2-7B-chat-en	0.726	95	0.395	19	0.373	15		0.128
leo-hessian-7B-chat-de	0.620	92	0.300	23	0.200	24		0.467

* $AR1$ stands for *Answer Relevance* where the *Context Position* = 1

** $AR>1$ stands for *Answer Relevance* where the *Context Position* > 1

Table 4.5.: Multiple Metrics per Model and Language

Noise Robustness: Table 4.5 illustrates the noise robustness of the system with an evidence set size of 10 for all models. Except for *leo-hessian-7B-chat*, the average answer relevance remains consistent regardless of whether the relevant context is in the first position or elsewhere in the evidence set. This suggests robustness against noise in evidence sets of length 10. However, for *leo-hessian-7B-chat*, there is a notable 10% drop in average answer relevance when the relevant context is not in the first position, indicating a lack of robustness against noise in evidence sets of length 10. This aligns with findings in Section 4.4.3, where this model showed reduced performance with increasing k in the evidence set.

Negative Rejection: The negative rejection performance varies significantly between the models. As shown in Table 4.5, *Llama2-7B-chat* exhibits the lowest negative rejection performance with a specificity of 0.128, while *gpt-3.5-turbo-en* demonstrates the highest specificity at 0.826. However, beyond the general trend of better performance with larger models, no further insights can be gleaned from this evaluation.

5. Conclusions and Future Work

This thesis serves as an introduction guide for developing conversational question-answering systems tailored to specific use cases. Beginning with an extensive overview of the field in Section 2, we delve into the essential algorithms and concepts necessary for constructing such systems.

Focusing on practical implementations, Section 3 narrows down the implementation landscape by focusing on conversational retrieval-augmented generators for question-answering. We outline a system comprising four distinct components: Extract, CQU, Retriever, and Reader. Each component addresses its own well-defined problem domain, shedding light on the major challenges associated with them. Notable contributions include the detailed breakdown of the often-overlooked Extract component into various operations and the identification of challenges faced by the Retriever component, both with the knowledge base and the evidence set for the Reader. Additionally, we dissect the tasks of the Reader component into smaller challenges and propose multiple existing and novel solutions to tackle them. Ultimately, the main contribution lies in providing a decision map for constructing a conversational question-answering system, based on the diverse components and implementation possibilities.

To validate the feasibility of the decision map, Chapter 4 presents the construction and evaluation of a default system primarily based on zero-shot components. Various evaluation approaches are introduced, implemented, and tested within this framework.

Specifically, the component-wise evaluation using synthetic data provides insights into the feasibility of synthetic datasets for metric-based component evaluation. Challenges arise with complex or abstract questions, particularly applicable to non-factoid questions, making human-based end-to-end evaluation preferable at this stage. Future work could explore using LLM agents to automate evaluation for real-world, use-case-specific datasets.

Furthermore, the performance of two smaller LLMs, *leo-hessian-7B-chat* and *Llama2-7B-chat*, is compared against *gpt-3.5-turbo*. Issues with the smaller LLMs, such as consistency in text generation quality and suitability as Readers, highlight considerations for their use in conversational question-answering systems.

Moreover, in the PoC utilizing state-of-the-art LLMs, the Retriever and Extractor

5. Conclusions and Future Work

components emerged as the main bottlenecks for Conv QA systems. Even in search scenarios with clear search intent and minimal ambiguity, the Retriever struggled to retrieve the correct passages from the knowledge base. This issue likely stems from challenges within both the Extractor and Retriever components. There is a clear need for further research in this area to address the difficulties in retrieving all relevant contexts from the knowledge base. Additionally, there is a pressing need for a universal solution capable of efficiently extracting textual or multimodal information from PDF documents, which remains a significant challenge in the field.

Additionally, as system complexity increases, so does inference time and computational intensity, especially with LLM-based approaches. Balancing computational resources may involve using combinations of smaller LLMs or alternative approaches for different subtasks.

In Section 4.4.1, we discuss the identified limitations of the chosen approaches, which is necessary reading for anyone interested in building their own conversational question-answering system and how to evaluate it.

We recommend any reader interested in building their own conversational question-answering system to start with this thesis as a foundational guide. They can then proceed to build a basic system similar to Section 4.3 and evaluate it in an end-to-end manner as described in Section 4.2.3. This approach provides the quickest way to understand the specific challenges of their use case and how to further iterate and improve their system.

Regarding future work, we propose the following directions for research:

- **Practical Feasibility of Conversational Question Answering Systems:** Evaluate the costs and user experience associated with implementing conversational question-answering systems, particularly considering the significant expenses linked with utilizing LLMs. This evaluation should delve into whether it is financially viable to construct a search system relying on LLMs, especially given the escalating costs associated with the growing complexity of utilized approaches and the increasing number of LLM inference runs required.
- **Constrained System Design:** Develop a search system that strikes a balance between resource utilization and search result quality. This endeavor may entail amalgamating multiple system components, implementing post-processing techniques on the evidence set, and orchestrating a series of reader components to create a system capable of addressing diverse question types while maintaining

5. Conclusions and Future Work

efficient resource usage. Relying solely on LLM-based approaches, such as RAG, may not always be the optimal solution for all use cases and setups. It might be more resource-efficient to employ a combination of smaller LLMs for different subtasks, or to explore techniques like quantization and incorporating extractive readers as post-processing steps, among others.

- **Improving the Retriever Component:** Enhance the performance of the retriever component, specifically focusing on its ability to locate all relevant contexts within the evidence set. This endeavor could entail investigating novel methodologies, such as employing agent-like retrievers or capitalizing on the structural elements of documents to direct the extraction and retrieval process. E.g. aligning the retrieval process with the search strategies employed by humans, who often hierarchically navigate through documents to pinpoint relevant passages, may offer valuable insights for enhancing retriever performance.
- **PDF Data Extraction:** Develop a universal solution for extracting textual or multimodal information from PDFs. This includes addressing the challenges of existing extraction operations, which yield different document models and may produce inconsistent passages.

A. Appendix A

A.1. Implementation Trees

A.1.1. Example

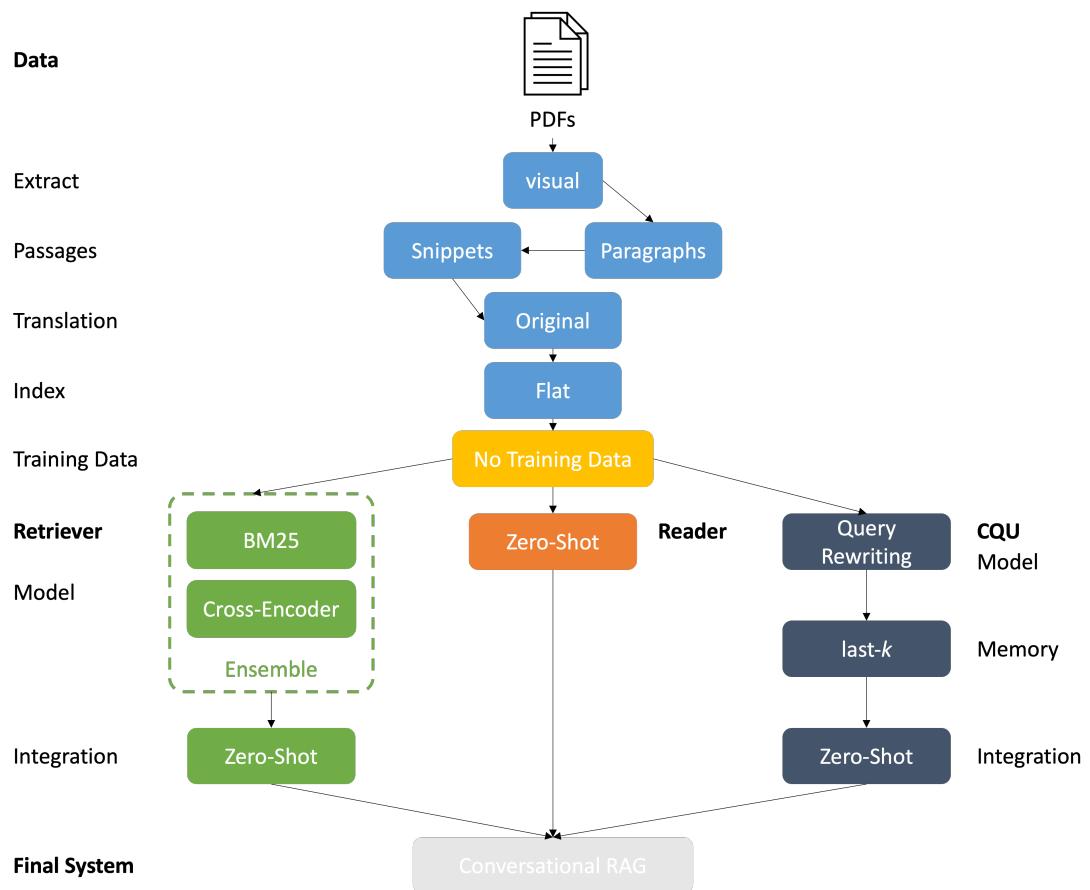


Figure A.1.: Example Implementation Zero-Shot Baseline

A.2. Data Augmentation Implementation

A.2.1. Examples for Few-Shot Prompt

1. Example:

Context: The master's program in Data and Computer Science includes an application area. Annex 3 lists the possible application areas. Upon request, the examination board can also approve a different application area. - Master Data and Computer Science

Question: I study Data and Computerscience. Can I choose application areas that are not listed?

2. Example:

Context: (5) The final failure in a mandatory module leads to the loss of the examination claim. In elective mandatory modules, if provided for in the module handbook, the failure can be compensated for by the successful completion of another elective mandatory module or another performance within the respective module. § 4 Paragraph 2 remains unaffected. - Bachelor Physics

Question: When does the failure of an exam lead to the termination of the study?

3. Example:

Context: (3) The application must be made in writing to the examination board. It is the responsibility of the applicant to provide the necessary information about the performance to be recognized. The burden of proof for the existence of a significant difference in academic achievements lies with Heidelberg University; the obligations of the applicant, especially according to Sentence 1 and Sentence 2, remain unaffected. The burden of proof for the existence of equivalence in non-academic achievements lies with the applicant. - Bachelor Philosophy

Question: What would happen, if I can't provide information on my academic achievements?

4. Example:

Context: The intermediate examination consists of successful participation in the exercises for beginners in the subjects Civil Law, Public Law, and

A. Appendix A

Criminal Law. The partial performances of the exercise (homework and supervisory work under examination conditions) must in principle be performed in the exercise of a semester; § 4 paragraph 5 remains unaffected. - Civil Law

Question: What are the components of the intermediate examination in Civil Law, Public Law, and Criminal Law?

A.3. Evaluation

A.3.1. Model Configurations

Application	Model-Language	<i>top k</i>	<i>top p</i>	<i>temperature</i>	<i>repetition penalty</i>
CQU	gpt-3.5-turbo-de	-*	1	1	0
	gpt-3.5-turbo-en	-*	1	1	0
	leo-hessian-7B-chat-de	60	0.95	0.9	1.2
	Llama2-7B-chat-en	1	0.1	0.3	1.4
Chat	gpt-3.5-turbo-de	-*	1	1	0
	gpt-3.5-turbo-en	-*	1	1	0
	leo-hessian-7B-chat-de	10	0.6	0.4	0
	Llama2-7B-chat-en	40	0.8	0.5	1.5

* OpenAI does not expose the *top k* parameter for the GPT-3.5-turbo models.

Table A.1.: Parameters used for Models in End-to-End Evaluation

A.3.2. Retriever

A. Appendix A

index metric retriever k	IndexEnglish						IndexGerman		
	BM25	DPR	HR DPR	MR DPR	BM25	DPR	HR DPR	BM25	DPR
1	0.236152	0.346251	0.236152	0.346251	0.221038	0.235225	0.221038	0.235225	0.235225
3	0.347355	0.496136	0.284705	0.412255	0.317794	0.342329	0.263626	0.282504	0.282504
5	0.400344	0.560026	0.296833	0.426853	0.360997	0.389765	0.273503	0.293309	0.293309
10	0.479333	0.645523	0.307341	0.438297	0.423759	0.459314	0.281863	0.302535	0.302535
15	0.530269	0.696031	0.311357	0.442275	0.462978	0.500956	0.284947	0.305811	0.305811
25	0.598848	0.760644	0.314831	0.445564	0.511582	0.556726	0.287416	0.308630	0.308630
50	0.693761	0.843895	0.317537	0.447936	0.584736	0.637499	0.289478	0.310920	0.310920
100	0.769780	0.910700	0.318649	0.448917	0.660721	0.718856	0.290569	0.312093	0.312093
150	0.795991	0.933902	0.318865	0.449109	0.698262	0.754966	0.290878	0.312391	0.312391
200	0.813112	0.946265	0.318964	0.449181	0.721834	0.777181	0.291015	0.312520	0.312520
300	0.836275	0.958705	0.319060	0.449232	0.751814	0.803585	0.291139	0.312629	0.312629
400	0.848770	0.966037	0.319096	0.449253	0.770087	0.820341	0.291192	0.312677	0.312677
500	0.857960	0.970819	0.319117	0.449264	0.783486	0.832689	0.291222	0.312705	0.312705
1000	0.884140	0.983089	0.319155	0.449282	0.818224	0.864172	0.291272	0.312751	0.312751

Table A.2.: Retriever comparison BM25 and DPR over index-English and indexGerman

A. Appendix A

index metric retriever k	translated IndexEnglish						translated IndexGerman		
	BM25	DPR	HR	BM25	MRR	DPR	BM25	HR	DPR
1	0.183295	0.300796	0.183295	0.300796	0.168289	0.226307	0.168289	0.226307	
3	0.283473	0.439959	0.227119	0.362059	0.255251	0.332098	0.206333	0.272683	
5	0.331586	0.502232	0.238078	0.376290	0.296455	0.379140	0.215718	0.283417	
10	0.402085	0.586524	0.247408	0.387556	0.356575	0.445857	0.223716	0.292320	
15	0.449383	0.638004	0.251129	0.391613	0.394480	0.488156	0.226698	0.295654	
25	0.514843	0.703348	0.254438	0.394939	0.446412	0.543634	0.229318	0.298465	
50	0.610612	0.793301	0.257162	0.397506	0.523755	0.625181	0.231502	0.300771	
100	0.702672	0.874887	0.258492	0.398697	0.606892	0.708026	0.232692	0.301962	
150	0.742715	0.903959	0.258823	0.398938	0.648592	0.747902	0.233035	0.302293	
200	0.768902	0.920427	0.258974	0.399033	0.675069	0.770949	0.233188	0.302426	
300	0.802740	0.938948	0.259113	0.399109	0.709398	0.798929	0.233329	0.302542	
400	0.823531	0.950183	0.259174	0.399142	0.730504	0.815699	0.233391	0.302590	
500	0.840124	0.957656	0.259211	0.399159	0.746355	0.828033	0.233426	0.302618	
1000	0.884777	0.973681	0.259276	0.399182	0.790784	0.861399	0.233490	0.302667	

Table A.3.: Retriever comparison BM25 and DPR over translated IndexEnglish and translated IndexGerman

A. Appendix A

index	metric	IndexEnglish						IndexGerman					
		BM25-	DPR	HR	MRR	BM25-	DPR	BM25-	DPR	HR	MRR	BM25-	DPR
retriever	CE-	BM25-	DPR	BM25-	DPR	BM25-	DPR	CE-	BM25-	DPR	CE-	BM25-	DPR
k		CE-	CE-	CE-	CE-	CE-	CE-	CE-	CE-	CE-	CE-	CE-	CE-
1		0.289670	0.290665	0.346251	0.289670	0.290665	0.346251	0.183031	0.187395	0.235225	0.183031	0.187395	0.235225
3		0.412333	0.413632	0.496136	0.343724	0.344858	0.412255	0.286573	0.292441	0.342329	0.228314	0.233467	0.282504
5		0.468588	0.470151	0.560026	0.356579	0.357755	0.426853	0.334973	0.342475	0.389765	0.239329	0.244864	0.293309
10		0.553105	0.553665	0.645523	0.367850	0.368904	0.438297	0.405514	0.414301	0.459314	0.248710	0.254398	0.302535

Table A.4: Retriever comparison BM25-CE versions and DPR
over IndexEnglish and IndexGerman

A. Appendix A

index metric	translated IndexEnglish						translated IndexGerman					
	HR	MRR	BM25- DPR	BM25- CE-	BM25- CE-	HR	DPR	BM25- DPR	BM25- CE-	BM25- CE-	MRR	
retriever	BM25- CE-	DPR	BM25- CE-	CE-	1000	500	500	1000	500	500	DPR	
1	0.226215	0.300796	0.226215	0.300796	0.210515	0.211157	0.226307	0.210515	0.211157	0.226307	0.226307	
3	0.357346	0.439959	0.283679	0.362059	NaN	0.315852	0.332098	NaN	0.257302	0.272683	0.272683	
5	0.418429	0.502232	0.297588	0.376290	0.363084	0.363435	0.379140	0.267510	0.268156	0.283417	0.283417	
10	0.500264	0.586524	0.308518	0.387556	0.430064	0.429218	0.445857	0.276427	0.276925	0.292320	0.292320	

Table A.5.: Retriever comparison BM25-CE versions and DPR over translated IndexEnglish and translated IndexGerman

Bibliography

- [Anil et al., 2023] Anil, R., Dai, A. M., Firat, O., Johnson, M., Lepikhin, D., Passos, A. T., Shakeri, S., Taropa, E., Bailey, P., and Chen, Z. (2023). Palm 2 technical report. arXiv. <https://arxiv.org/abs/2305.10403>.
- [AutoGPTQ, 2023] AutoGPTQ (2023). Autogptq. <https://github.com/PanQiWei/AutoGPTQ>. Accessed: 2023-10-02.
- [Brown et al., 2020] Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., and Askell, A. (2020). Language models are few-shot learners. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, NIPS’20, Red Hook, NY, USA. Curran Associates Inc.
- [Chen and Ma, 2021] Chen, H. and Ma, D. (2021). Improving Out-of-Domain Question Answering with Mixture of Experts. Technical report, Department of Computer Science, Stanford University.
- [Cohen, 1960] Cohen, J. (1960). A Coefficient of Agreement for Nominal Scales. *Educational and Psychological Measurement*, 20(1):37–46.
- [Dai et al., 2022a] Dai, Z., Chaganty, A. T., Zhao, V. Y., Amini, A., Rashid, Q. M., Green, M., and Guu, K. (2022a). Dialog inpainting: Turning documents into dialogs. In Chaudhuri, K., Jegelka, S., Song, L., Szepesvári, C., Niu, G., and Sabato, S., editors, *International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA*, volume 162 of *Proceedings of Machine Learning Research*, pages 4558–4586. PMLR.
- [Dai et al., 2022b] Dai, Z., Zhao, V. Y., Ma, J., Luan, Y., Ni, J., Lu, J., Bakalov, A., Guu, K., Hall, K. B., and Chang, M.-W. (2022b). Promptagator: Few-shot Dense Retrieval From 8 Examples. arXiv. <http://arxiv.org/abs/2209.11755>.
- [Dalton et al., 2020] Dalton, J., Xiong, C., and Callan, J. (2020). TREC cast 2019: The conversational assistance track overview. arXiv. <https://arxiv.org/abs/2003.13624>.

Bibliography

- [Dasigi et al., 2021] Dasigi, P., Lo, K., Beltagy, I., Cohan, A., Smith, N. A., and Gardner, M. (2021). A dataset of information-seeking questions and answers anchored in research papers. In Toutanova, K., Rumshisky, A., Zettlemoyer, L., Hakkani-Tur, D., Beltagy, I., Bethard, S., Cotterell, R., Chakraborty, T., and Zhou, Y., editors, *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4599–4610, Online. Association for Computational Linguistics.
- [Devlin et al., 2019] Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. (2019). BERT: Pre-training of deep bidirectional transformers for language understanding. In Burstein, J., Doran, C., and Solorio, T., editors, *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)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- [Dimitrakis et al., 2020] Dimitrakis, E., Sgontzos, K., and Tzitzikas, Y. (2020). A survey on question-answering systems over linked data and documents. *Journal of Intelligent Information Systems*, 55(2):233–259.
- [Ding et al., 2022] Ding, Y., Huang, Z., Wang, R., Zhang, Y., Chen, X., Ma, Y., Chung, H., and Han, S. (2022). V-doc : Visual questions answers with documents. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 21460–21466, Los Alamitos, CA, USA. IEEE Computer Society.
- [Elgohary et al., 2019] Elgohary, A., Peskov, D., and Boyd-Graber, J. (2019). Can You Unpack That? Learning to Rewrite Questions-in-Context. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 5918–5924, Hong Kong, China. Association for Computational Linguistics.
- [Etezadi and Shamsfard, 2023] Etezadi, R. and Shamsfard, M. (2023). The state of the art in open domain complex question-answering: a survey. *Applied Intelligence*, 53(4):4124–4144.
- [Farea et al., 2022] Farea, A. A. S., Yang, Z., Duong, K., Perera, N., and Emmert-Streib, F. (2022). Evaluation of question answering systems: Complexity of judging a natural language. arXiv. <http://arxiv.org/abs/2209.12617>.
- [Feng, 2021] Feng, S. (2021). DialDoc 2021 Shared Task: Goal-Oriented Document-grounded Dialogue Modeling. In *Proceedings of the 1st Workshop on Document-*

Bibliography

- grounded Dialogue and Conversational Question Answering (DialDoc 2021)*, pages 1–7, Online. Association for Computational Linguistics.
- [Ferrucci, 2012] Ferrucci, D. A. (2012). Introduction to “This is Watson”. *IBM Journal of Research and Development*, 56(3.4):1:1–1:15.
- [Frantar and Alistarh, 2022] Frantar, E. and Alistarh, D. (2022). Optimal brain compression: A framework for accurate post-training quantization and pruning. In Koyejo, S., Mohamed, S., Agarwal, A., Belgrave, D., Cho, K., and Oh, A., editors, *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA*.
- [Frantar and Alistarh, 2023] Frantar, E. and Alistarh, D. (2023). Sparsegpt: massive language models can be accurately pruned in one-shot. In *Proceedings of the 40th International Conference on Machine Learning*, ICML’23. JMLR.org.
- [Frantar et al., 2023] Frantar, E., Ashkboos, S., Hoefer, T., and Alistarh, D. (2023). GPTQ: Accurate Post-Training Quantization for Generative Pre-trained Transformers. arXiv. <http://arxiv.org/abs/2210.17323>.
- [Gangi Reddy et al., 2021] Gangi Reddy, R., Iyer, B., Sultan, M. A., Zhang, R., Sil, A., Castelli, V., Florian, R., and Roukos, S. (2021). Synthetic target domain supervision for open retrieval qa. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR ’21, page 1793–1797, New York, NY, USA. Association for Computing Machinery.
- [Gao et al., 2023] Gao, J., Xiong, C., Bennett, P., and Craswell, N. (2023). *Neural Approaches to Conversational Information Retrieval*, volume 44 of *The Information Retrieval Series*. Springer.
- [Gao et al., 2024] Gao, Y., Xiong, Y., Gao, X., Jia, K., Pan, J., Bi, Y., Dai, Y., Sun, J., Guo, Q., and Wang, M. (2024). Retrieval-Augmented Generation for Large Language Models: A Survey. arXiv. <http://arxiv.org/abs/2312.10997>.
- [Gholami et al., 2022] Gholami, A., Kim, S., Zhen, D., Yao, Z., Mahoney, M., and Keutzer, K. (2022). *A Survey of Quantization Methods for Efficient Neural Network Inference*, pages 291–326.
- [Gholami and Noori, 2021] Gholami, S. and Noori, M. (2021). Zero-shot open-book question answering. arXiv. <https://arxiv.org/abs/2111.11520>.

Bibliography

- [Green et al., 1961] Green, B. F., Wolf, A. K., Chomsky, C., and Laughery, K. (1961). Baseball: an automatic question-answerer. In *Papers presented at the May 9-11, 1961, western joint IRE-AIEE-ACM computer conference*, IRE-AIEE-ACM '61 (Western), pages 219–224, New York, NY, USA. Association for Computing Machinery.
- [Gu et al., 2023] Gu, Y., Dong, L., Wei, F., and Huang, M. (2023). Knowledge Distillation of Large Language Models. arXiv. <http://arxiv.org/abs/2306.08543>.
- [Guo et al., 2021] Guo, M., Zhang, M., Reddy, S., and Alikhani, M. (2021). Abg-coqa: Clarifying ambiguity in conversational question answering. In Chen, D., Berant, J., McCallum, A., and Singh, S., editors, *3rd Conference on Automated Knowledge Base Construction, AKBC 2021, October 4-8, 2021*, Online.
- [Gupta et al., 2020] Gupta, S., Rawat, B. P. S., and Yu, H. (2020). Conversational Machine Comprehension: a Literature Review. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 2739–2753, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- [Gururangan et al., 2020] Gururangan, S., Marasović, A., Swayamdipta, S., Lo, K., Beltagy, I., Downey, D., and Smith, N. A. (2020). Don’t Stop Pretraining: Adapt Language Models to Domains and Tasks. *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8342–8360.
- [Guu et al., 2020] Guu, K., Lee, K., Tung, Z., Pasupat, P., and Chang, M.-W. (2020). Realm: retrieval-augmented language model pre-training. In *Proceedings of the 37th International Conference on Machine Learning*, ICML’20. JMLR.org.
- [Hao et al., 2022] Hao, T., Li, X., He, Y., Wang, F. L., and Qu, Y. (2022). Recent progress in leveraging deep learning methods for question-answering. *Neural Computing and Applications*, 34(4):2765–2783.
- [Harabagiu et al., 2003] Harabagiu, S. M., Maiorano, S. J., and Pașca, M. A. (2003). Open-domain textual question-answering techniques. *Natural Language Engineering*, 9(3):231–267.
- [He et al., 2021] He, P., Liu, X., Gao, J., and Chen, W. (2021). Deberta: decoding-enhanced bert with disentangled attention. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net.
- [Hinton et al., 2015] Hinton, G., Vinyals, O., and Dean, J. (2015). Distilling the Knowledge in a Neural Network. arXiv. <http://arxiv.org/abs/1503.02531>.

Bibliography

- [Houlsby et al., 2019] Houlsby, N., Giurgiu, A., Jastrzebski, S., Morrone, B., de Laros-silhe, Q., Gesmundo, A., Attariyan, M., and Gelly, S. (2019). Parameter-efficient transfer learning for NLP. In Chaudhuri, K. and Salakhutdinov, R., editors, *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning Research*, pages 2790–2799. PMLR.
- [Hu et al., 2022] Hu, E. J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., and Chen, W. (2022). Lora: Low-rank adaptation of large language models. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net.
- [Huang et al., 2022] Huang, Y., Chen, Y., Yu, Z., and McKeown, K. (2022). In-context Learning Distillation: Transferring Few-shot Learning Ability of Pre-trained Language Models. arXiv. <http://arxiv.org/abs/2212.10670>.
- [Huggingface, 2023a] Huggingface (2023a). Peft. <https://huggingface.co/docs/peft/index>. Accessed: 2023-10-01.
- [Huggingface, 2023b] Huggingface (2023b). Quantize transformers models. https://huggingface.co/docs/transformers/main_classes/quantization. Accessed: 2023-10-02.
- [Izacard and Grave, 2021] Izacard, G. and Grave, E. (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*, pages 874–880, Online. Association for Computational Linguistics.
- [Izacard et al., 2023] Izacard, G., Lewis, P. S. H., Lomeli, M., Hosseini, L., Petroni, F., Schick, T., Dwivedi-Yu, J., Joulin, A., Riedel, S., and Grave, E. (2023). Atlas: Few-shot learning with retrieval augmented language models. *J. Mach. Learn. Res.*, 24:251:1–251:43.
- [Jiang et al., 2023] Jiang, Y., Chan, C., Chen, M., and Wang, W. (2023). Lion: Adversarial distillation of proprietary large language models. In *Conference on Empirical Methods in Natural Language Processing*.
- [Johnson et al., 2021] Johnson, J., Douze, M., and Jégou, H. (2021). Billion-scale similarity search with gpus. *IEEE Transactions on Big Data*, 7(3):535–547.

Bibliography

- [Jurafsky and Martin, 2023] Jurafsky, D. and Martin, J. H. (2023). *Speech and Language Processing*. Palo Alto, 3 edition.
- [Kamalloo et al., 2023] Kamalloo, E., Dziri, N., Clarke, C. L. A., and Rafiei, D. (2023). Evaluating open-domain question answering in the era of large language models. In Rogers, A., Boyd-Graber, J. L., and Okazaki, N., editors, *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 5591–5606. Association for Computational Linguistics.
- [Karpukhin et al., 2020] Karpukhin, V., Oguz, B., Min, S., Lewis, P., Wu, L., Edunov, S., Chen, D., and Yih, W.-t. (2020). Dense passage retrieval for open-domain question answering. In Webber, B., Cohn, T., He, Y., and Liu, Y., editors, *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6769–6781, Online. Association for Computational Linguistics.
- [Khashabi et al., 2022] Khashabi, D., Kordi, Y., and Hajishirzi, H. (2022). Unifiedqa-v2: Stronger generalization via broader cross-format training. arXiv. <http://arxiv.org/abs/2202.12359>.
- [Khattab and Zaharia, 2020] Khattab, O. and Zaharia, M. (2020). Colbert: Efficient and effective passage search via contextualized late interaction over bert. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR ’20, page 39–48, New York, NY, USA. Association for Computing Machinery.
- [Kuhn et al., 2023] Kuhn, L., Gal, Y., and Farquhar, S. (2023). CLAM: Selective Clarification for Ambiguous Questions with Generative Language Models. arXiv. <http://arxiv.org/abs/2212.07769>.
- [Langchain, 2023] Langchain (2023). Building applications with llms through compositability. <https://github.com/langchain-ai/langchain>. Accessed: 2023-09-22.
- [Langchain, 2023] Langchain (2023). Question answering | langchain. https://python.langchain.com/docs/use_cases/question_answering/. Accessed: 2023-09-24.
- [Larsson, 2002] Larsson, S. (2002). *Issue-based Dialogue Management*. PhD thesis, Department of Linguistics, Goeteborg University.
- [Lester et al., 2021] Lester, B., Al-Rfou, R., and Constant, N. (2021). The Power of Scale for Parameter-Efficient Prompt Tuning. In *Proceedings of the 2021 Conference*

Bibliography

on Empirical Methods in Natural Language Processing, pages 3045–3059, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

- [Lewis et al., 2020a] Lewis, M., Liu, Y., Goyal, N., Ghazvininejad, M., Mohamed, A., Levy, O., Stoyanov, V., and Zettlemoyer, L. (2020a). BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In Jurafsky, D., Chai, J., Schluter, N., and Tetreault, J., editors, *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.
- [Lewis et al., 2020b] Lewis, P., Perez, E., Piktus, A., Petroni, F., Karpukhin, V., Goyal, N., Küttler, H., Lewis, M., Yih, W.-t., and Rocktäschel, T. (2020b). Retrieval-augmented generation for knowledge-intensive nlp tasks. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, NIPS’20, Red Hook, NY, USA. Curran Associates Inc.
- [Li et al., 2022a] Li, J., Xu, Y., Lv, T., Cui, L., Zhang, C., and Wei, F. (2022a). Dit: Self-supervised pre-training for document image transformer. In *Proceedings of the 30th ACM International Conference on Multimedia*, MM ’22, page 3530–3539, New York, NY, USA. Association for Computing Machinery.
- [Li et al., 2022b] Li, S., Chen, J., Shen, Y., Chen, Z., Zhang, X., Li, Z., Wang, H., Qian, J., Peng, B., and Mao, Y. (2022b). Explanations from Large Language Models Make Small Reasoners Better. arXiv. <http://arxiv.org/abs/2210.06726>.
- [Li and Liang, 2021] Li, X. L. and Liang, P. (2021). Prefix-tuning: Optimizing continuous prompts for generation. In Zong, C., Xia, F., Li, W., and Navigli, R., editors, *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 1: Long Papers), Virtual Event, August 1-6, 2021*, pages 4582–4597. Association for Computational Linguistics.
- [Lin, 2004] Lin, C.-Y. (2004). ROUGE: A Package for Automatic Evaluation of Summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.
- [Ling et al., 2023] Ling, C., Zhao, X., Lu, J., Deng, C., Zheng, C., Wang, J., Chowdhury, T., Li, Y., Cui, H., and Zhang, X. (2023). Domain Specialization as the Key to Make Large Language Models Disruptive: A Comprehensive Survey. arXiv. <http://arxiv.org/abs/2305.18703>.

Bibliography

- [Liu et al., 2022] Liu, H., Tam, D., Muqeeth, M., Mohta, J., Huang, T., Bansal, M., and Raffel, C. (2022). Few-shot parameter-efficient fine-tuning is better and cheaper than in-context learning. In Koyejo, S., Mohamed, S., Agarwal, A., Belgrave, D., Cho, K., and Oh, A., editors, *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*.
- [Liu et al., 2023a] Liu, X., Zheng, Y., Du, Z., Ding, M., Qian, Y., Yang, Z., and Tang, J. (2023a). Gpt understands, too. *AI Open*.
- [Liu et al., 2021] Liu, Y., Hashimoto, K., Zhou, Y., Yavuz, S., Xiong, C., and Yu, P. (2021). Dense hierarchical retrieval for open-domain question answering. In Moens, M.-F., Huang, X., Specia, L., and Yih, S. W.-t., editors, *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 188–200, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- [Liu et al., 2019] Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L., and Stoyanov, V. (2019). Roberta: A robustly optimized BERT pretraining approach. arXiv. <http://arxiv.org/abs/1907.11692>.
- [Liu et al., 2023b] Liu, Z., Oguz, B., Zhao, C., Chang, E., Stock, P., Mehdad, Y., Shi, Y., Krishnamoorthi, R., and Chandra, V. (2023b). LLM-QAT: Data-Free Quantization Aware Training for Large Language Models. arXiv. <http://arxiv.org/abs/2305.17888>.
- [Liusie et al., 2022] Liusie, A., Qian, M., Li, X., and Gales, M. J. F. (2022). University of cambridge at TREC cast 2022. In Soboroff, I. and Ellis, A., editors, *Proceedings of the Thirty-First Text REtrieval Conference, TREC 2022, online, November 15-19, 2022*, volume 500-338 of *NIST Special Publication*. National Institute of Standards and Technology (NIST).
- [Luo et al., 2022] Luo, M., Hashimoto, K., Yavuz, S., Liu, Z., Baral, C., and Zhou, Y. (2022). Choose your qa model wisely: A systematic study of generative and extractive readers for question answering. In *Proceedings of the 1st Workshop on Semiparametric Methods in NLP: Decoupling Logic from Knowledge*, pages 7–22.
- [Ma et al., 2023] Ma, X., Fang, G., and Wang, X. (2023). LLM-Pruner: On the Structural Pruning of Large Language Models. arXiv. <http://arxiv.org/abs/2305.11627>.

Bibliography

- [Magic, 2023] Magic, N. (2023). Deepsparse. <https://github.com/neuralmagic/deepsparse>. Accessed: 2023-10-02.
- [Mao et al., 2023] Mao, K., Dou, Z., Mo, F., Hou, J., Chen, H., and Qian, H. (2023). Large language models know your contextual search intent: A prompting framework for conversational search. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 1211–1225.
- [Mathew et al., 2021] Mathew, M., Tito, R., Karatzas, D., Manmatha, R., and Jawahar, C. V. (2021). Document Visual Question Answering Challenge 2020. arXiv. <http://arxiv.org/abs/2008.08899>.
- [McDonald et al., 2022] McDonald, T., Tsan, B., Saini, A., Ordonez, J., Gutierrez, L., Nguyen, P., Mason, B., and Ng, B. (2022). Detect, Retrieve, Comprehend: A Flexible Framework for Zero-Shot Document-Level Question Answering. arXiv. <http://arxiv.org/abs/2210.01959>.
- [Meuschke et al., 2023] Meuschke, N., Jagdale, A., Spinde, T., Mitrovic, J., and Gipp, B. (2023). A benchmark of PDF information extraction tools using a multi-task and multi-domain evaluation framework for academic documents. In Sserwanga, I., Goulding, A., Sandy, H. M., Du, J. T., Soares, A. L., Hessami, V., and Frank, R. D., editors, *Information for a Better World: Normality, Virtuality, Physicality, Inclusivity - 18th International Conference, iConference 2023, Virtual Event, March 13-17, 2023, Proceedings, Part II*, volume 13972 of *Lecture Notes in Computer Science*, pages 383–405. Springer.
- [Mishra and Jain, 2016] Mishra, A. and Jain, S. K. (2016). A survey on question-answering systems with classification. *Journal of King Saud University - Computer and Information Sciences*, 28(3):345–361.
- [Muennighoff, 2022] Muennighoff, N. (2022). SGPT: GPT Sentence Embeddings for Semantic Search. arXiv. <http://arxiv.org/abs/2202.08904>.
- [Nakano et al., 2022] Nakano, R., Hilton, J., Balaji, S., Wu, J., Ouyang, L., Kim, C., Hesse, C., Jain, S., Kosaraju, V., and Saunders, W. (2022). WebGPT: Browser-assisted question-answering with human feedback. arXiv. <http://arxiv.org/abs/2112.09332>.
- [Nassiri and Akhloufi, 2023] Nassiri, K. and Akhloufi, M. (2023). Transformer models used for text-based question-answering systems. *Applied Intelligence*, 53(9):10602–10635.

Bibliography

- [Naveed et al., 2023] Naveed, H., Khan, A. U., Qiu, S., Saqib, M., Anwar, S., Usman, M., Akhtar, N., Barnes, N., and Mian, A. (2023). A Comprehensive Overview of Large Language Models. arXiv. <http://arxiv.org/abs/2307.06435>.
- [Neelakantan et al., 2022] Neelakantan, A., Xu, T., Puri, R., Radford, A., Han, J. M., Tworek, J., Yuan, Q., Tezak, N. A., Kim, J. W., and Hallacy, C. (2022). Text and code embeddings by contrastive pre-training. arXiv. <https://arxiv.org/abs/2201.10005>.
- [Nguyen et al., 2016] Nguyen, T., Rosenberg, M., Song, X., Gao, J., Tiwary, S., Majumder, R., and Deng, L. (2016). MS MARCO: A human generated machine reading comprehension dataset. In Besold, T. R., Bordes, A., d’Avila Garcez, A. S., and Wayne, G., editors, *Proceedings of the Workshop on Cognitive Computation: Integrating neural and symbolic approaches 2016 co-located with the 30th Annual Conference on Neural Information Processing Systems (NIPS 2016), Barcelona, Spain, December 9, 2016*, volume 1773 of *CEUR Workshop Proceedings*. CEUR-WS.org.
- [Ni et al., 2022] Ni, J., Qu, C., Lu, J., Dai, Z., Hernandez Abrego, G., Ma, J., Zhao, V., Luan, Y., Hall, K., Chang, M.-W., and Yang, Y. (2022). Large dual encoders are generalizable retrievers. In Goldberg, Y., Kozareva, Z., and Zhang, Y., editors, *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 9844–9855, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- [Nishida et al., 2018] Nishida, K., Saito, I., Otsuka, A., Asano, H., and Tomita, J. (2018). Retrieve-and-Read: Multi-task Learning of Information Retrieval and Reading Comprehension. In *Proceedings of the 27th ACM International Conference on Information and Knowledge Management*, pages 647–656.
- [OpenAI, 2023] OpenAI (2023). Chatgpt retrieval plugin. <https://github.com/openai/chatgpt-retrieval-plugin>. Accessed: 2023-09-22.
- [Ouyang et al., 2022] Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C. L., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., and Ray, A. (2022). Training language models to follow instructions with human feedback. In Koyejo, S., Mohamed, S., Agarwal, A., Belgrave, D., Cho, K., and Oh, A., editors, *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*.

Bibliography

- [Owoicho et al., 2022] Owoicho, P., Dalton, J., Aliannejadi, M., Azzopardi, L., Trippas, J. R., and Vakulenko, S. (2022). TREC cast 2022: Going beyond user ask and system retrieve with initiative and response generation. In Soboroff, I. and Ellis, A., editors, *Proceedings of the Thirty-First Text REtrieval Conference, TREC 2022, online, November 15-19, 2022*, volume 500-338 of *NIST Special Publication*. National Institute of Standards and Technology (NIST).
- [Papineni et al., 2002] Papineni, K., Roukos, S., Ward, T., and Zhu, W.-J. (2002). Bleu: a Method for Automatic Evaluation of Machine Translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.
- [Pereira et al., 2023] Pereira, J. A., do Nascimento Fidalgo, R., de Alencar Lotufo, R., and Nogueira, R. F. (2023). Visconde: Multi-document QA with GPT-3 and neural reranking. In Kamps, J., Goeuriot, L., Crestani, F., Maistro, M., Joho, H., Davis, B., Gurrin, C., Kruschwitz, U., and Caputo, A., editors, *Advances in Information Retrieval - 45th European Conference on Information Retrieval, ECIR 2023, Dublin, Ireland, April 2-6, 2023, Proceedings, Part II*, volume 13981 of *Lecture Notes in Computer Science*, pages 534–543. Springer.
- [Plüster, 2023] Plüster, B. (2023). Leolm: Ein impuls für deutschsprachige llm-forschung | laion. <https://laion.ai/blog-de/leo-lm/>. Accessed: 2023-10-17.
- [PyPDF2, 2023] PyPDF2 (2023). Welcome to pypdf2 — pypdf2 documentation. <https://pypdf2.readthedocs.io/en/3.0.0/index.html>. Accessed: 2023-09-25.
- [Raffel et al., 2020] Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., and Liu, P. J. (2020). Exploring the limits of transfer learning with a unified text-to-text transformer. *J. Mach. Learn. Res.*, 21(1).
- [Ragas, 2024] Ragas (2024). Core concepts | ragas. <https://docs.ragas.io/en/stable/concepts/index.html>. Accessed: 2024-01-11.
- [Rajpurkar et al., 2016] Rajpurkar, P., Zhang, J., Lopyrev, K., and Liang, P. (2016). SQuAD: 100,000+ questions for machine comprehension of text. In Su, J., Duh, K., and Carreras, X., editors, *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2383–2392, Austin, Texas. Association for Computational Linguistics.

Bibliography

- [Rastogi et al., 2020] Rastogi, A., Zang, X., Sunkara, S., Gupta, R., and Khaitan, P. (2020). Schema-guided dialogue state tracking task at dstc8. arXiv. <https://arxiv.org/abs/2002.01359>.
- [Reddy et al., 2019] Reddy, S., Chen, D., and Manning, C. D. (2019). CoQA: A conversational question answering challenge. *Transactions of the Association for Computational Linguistics*, 7:249–266.
- [Roberts et al., 2020] Roberts, A., Raffel, C., and Shazeer, N. (2020). How Much Knowledge Can You Pack Into the Parameters of a Language Model? In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 5418–5426, Online. Association for Computational Linguistics.
- [Robertson and Zaragoza, 2009] Robertson, S. and Zaragoza, H. (2009). The Probabilistic Relevance Framework: BM25 and Beyond. *Foundations and Trends in Information Retrieval*, 3:333–389.
- [Sachan et al., 2023] Sachan, D. S., Lewis, M., Yogatama, D., Zettlemoyer, L., Pineau, J., and Zaheer, M. (2023). Questions are all you need to train a dense passage retriever. *Transactions of the Association for Computational Linguistics*, 11:600–616.
- [Serban et al., 2016] Serban, I. V., García-Durán, A., Gulcehre, C., Ahn, S., Chandar, S., Courville, A., and Bengio, Y. (2016). Generating Factoid Questions With Recurrent Neural Networks: The 30M Factoid Question-Answer Corpus. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 588–598, Berlin, Germany. Association for Computational Linguistics.
- [Shen et al., 2023] Shen, S., Hou, L., Zhou, Y., Du, N., Longpre, S., Wei, J., Chung, H. W., Zoph, B., Fedus, W., and Chen, X. (2023). Flan-moe: Scaling instruction-finetuned language models with sparse mixture of experts. arXiv. <https://arxiv.org/abs/2305.14705>.
- [Shewale, 2024] Shewale, R. (2024). Chatgpt statistics. <https://www.demandsage.com/chatgpt-statistics/>. Accessed: 2024-02-20.
- [Tavakoli et al., 2021] Tavakoli, L., Zamani, H., Scholer, F., Croft, W., and Sanderson, M. (2021). Analyzing clarification in asynchronous information-seeking conversations. *Journal of the Association for Information Science and Technology*, 73.

Bibliography

- [Thakur et al., 2021] Thakur, N., Reimers, N., Rücklé, A., Srivastava, A., and Gurevych, I. (2021). BEIR: Heterogenous benchmark for zero-shot evaluation of information retrieval models. In *Proceedings of the 2021 Neural Information Processing Systems (NeurIPS-2021): Track on Datasets and Benchmarks*.
- [Thoppilan et al., 2022] Thoppilan, R., De Freitas, D., Hall, J., Shazeer, N., Kulshreshtha, A., Cheng, H.-T., Jin, A., Bos, T., Baker, L., and Du, Y. (2022). LaMDA: Language Models for Dialog Applications. arXiv. <http://arxiv.org/abs/2201.08239>.
- [Tito et al., 2021] Tito, R., Karatzas, D., and Valveny, E. (2021). Document Collection Visual Question Answering. volume 12822, pages 778–792. arXiv:2104.14336 [cs].
- [Touvron et al., 2023] Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., and Bhosale, S. (2023). Llama 2: Open Foundation and Fine-Tuned Chat Models. arXiv. <https://arxiv.org/abs/2307.09288>.
- [Treviso et al., 2023] Treviso, M., Lee, J.-U., Ji, T., Aken, B. V., Cao, Q., Ciosici, M. R., Hassid, M., Heafield, K., Hooker, S., and Raffel, C. (2023). Efficient Methods for Natural Language Processing: A Survey. *Transactions of the Association for Computational Linguistics*, 11:826–860.
- [TruLens, 2024] TruLens (2024). Rag triad - trulens. https://www.trulens.org/trulens_eval/core_concepts_rag_triad/. Accessed: 2024-01-11.
- [Voorhees, 1999] Voorhees, E. M. (1999). The TREC-8 question answering track report. In Voorhees, E. M. and Harman, D. K., editors, *Proceedings of The Eighth Text REtrieval Conference, TREC 1999, Gaithersburg, Maryland, USA, November 17-19, 1999*, volume 500-246 of *NIST Special Publication*. National Institute of Standards and Technology (NIST).
- [Voskarides et al., 2020] Voskarides, N., Li, D., Ren, P., Kanoulas, E., and de Rijke, M. (2020). Query Resolution for Conversational Search with Limited Supervision. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 921–930.
- [Wang et al., 2020] Wang, W., Wei, F., Dong, L., Bao, H., Yang, N., and Zhou, M. (2020). Minilm: deep self-attention distillation for task-agnostic compression of pre-trained transformers. In *Proceedings of the 34th International Conference on Neural*

Bibliography

Information Processing Systems, NIPS’20, Red Hook, NY, USA. Curran Associates Inc.

[Wang, 2022] Wang, Z. (2022). Modern Question Answering Datasets and Benchmarks: A Survey. arXiv. <https://arxiv.org/abs/2206.15030>.

[Wang et al., 2019] Wang, Z., Ng, P., Ma, X., Nallapati, R., and Xiang, B. (2019). Multi-passage BERT: A globally normalized BERT model for open-domain question answering. In Inui, K., Jiang, J., Ng, V., and Wan, X., editors, *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 5878–5882, Hong Kong, China. Association for Computational Linguistics.

[Wei et al., 2022] Wei, J., Wang, X., Schuurmans, D., Bosma, M., Ichter, B., Xia, F., Chi, E. H., Le, Q. V., and Zhou, D. (2022). Chain-of-thought prompting elicits reasoning in large language models. In Koyejo, S., Mohamed, S., Agarwal, A., Belgrave, D., Cho, K., and Oh, A., editors, *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*.

[White et al., 2023] White, J., Fu, Q., Hays, S., Sandborn, M., Olea, C., Gilbert, H., Elnashar, A., Spencer-Smith, J., and Schmidt, D. C. (2023). A prompt pattern catalog to enhance prompt engineering with chatgpt. arXiv. <https://arxiv.org/abs/2302.11382>.

[Yang et al., 2019] Yang, J., Lin, S., Wang, C., Lin, J., and Tsai, M. (2019). Query and answer expansion from conversation history. In Voorhees, E. M. and Ellis, A., editors, *Proceedings of the Twenty-Eighth Text REtrieval Conference, TREC 2019, Gaithersburg, Maryland, USA, November 13-15, 2019*, volume 1250 of *NIST Special Publication*. National Institute of Standards and Technology (NIST).

[Yang et al., 2018] Yang, Z., Qi, P., Zhang, S., Bengio, Y., Cohen, W., Salakhutdinov, R., and Manning, C. D. (2018). HotpotQA: A dataset for diverse, explainable multi-hop question answering. In Riloff, E., Chiang, D., Hockenmaier, J., and Tsujii, J., editors, *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2369–2380, Brussels, Belgium. Association for Computational Linguistics.

[Zaib et al., 2022] Zaib, M., Zhang, W. E., Sheng, Q. Z., Mahmood, A., and Zhang, Y. (2022). Conversational question-answering: a survey. *Knowl. Inf. Syst.*, 64(12):3151–3195.

Bibliography

- [Zamani et al., 2020] Zamani, H., Mitra, B., Chen, E., Lueck, G., Diaz, F., Bennett, P. N., Craswell, N., and Dumais, S. T. (2020). Analyzing and learning from user interactions for search clarification. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR ’20, page 1181–1190, New York, NY, USA. Association for Computing Machinery.
- [Zamani et al., 2023] Zamani, H., Trippas, J. R., Dalton, J., and Radlinski, F. (2023). Conversational Information Seeking. arXiv. <http://arxiv.org/abs/2201.08808>.
- [Zhang et al., 2023a] Zhang, J., Zhang, H., Zhang, D., Liu, Y., and Huang, S. (2023a). Beam Retrieval: General End-to-End Retrieval for Multi-Hop Question Answering. arXiv. <http://arxiv.org/abs/2308.08973>.
- [Zhang et al., 2023b] Zhang, Q., Chen, S., Xu, D., Cao, Q., Chen, X., Cohn, T., and Fang, M. (2023b). A Survey for Efficient Open Domain Question Answering. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 14447–14465, Toronto, Canada. Association for Computational Linguistics.
- [Zhang et al., 2020] Zhang, T., Kishore, V., Wu, F., Weinberger, K. Q., and Artzi, Y. (2020). BERTscore: Evaluating text generation with BERT. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.
- [Zhao et al., 2023] Zhao, W. X., Zhou, K., Li, J., Tang, T., Wang, X., Hou, Y., Min, Y., Zhang, B., Zhang, J., and Dong, Z. (2023). A survey of large language models. arXiv. <https://arxiv.org/abs/2303.18223>.
- [Zhu et al., 2021] Zhu, F., Lei, W., Wang, C., Zheng, J., Poria, S., and Chua, T. (2021). Retrieving and reading: A comprehensive survey on open-domain question answering. arXiv. <https://arxiv.org/abs/2101.00774>.
- [Zhu et al., 2023] Zhu, X., Li, J., Liu, Y., Ma, C., and Wang, W. (2023). A Survey on Model Compression for Large Language Models. arXiv. <https://arxiv.org/abs/2308.07633v3>.
- [Zhu et al., 2022] Zhu, Y., Liu, N., Xu, Z., Liu, X., Meng, W., Wang, L., Ou, Z., and Tang, J. (2022). Teach less, learn more: On the undistillable classes in knowledge distillation. In Koyejo, S., Mohamed, S., Agarwal, A., Belgrave, D., Cho, K., and Oh,

Bibliography

A., editors, *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022.*