

# Improving RAG Systems via Sentence Clustering and Reordering

Marco Alessio  
marco.alessio@isti.cnr.it  
ISTI-CNR  
Pisa, Italy

Guglielmo Faggioli  
guglielmo.faggioli@unipd.it  
University of Padua  
Padua, Italy

Nicola Ferro  
nicola.ferro@unipd.it  
University of Padua  
Padua, Italy

Franco Maria Nardini  
francomaria.nardini@isti.cnr.it  
ISTI-CNR  
Pisa, Italy

Raffaele Perego  
raffaele.perego@isti.cnr.it  
ISTI-CNR  
Pisa, Italy

## ABSTRACT

Large Language Models (LLMs) have gained noteworthy importance and attention across different domains and fields in recent years. Information Retrieval (IR) is one of the domains they impacted the most, as witnessed by the recent increase in the number of IR systems incorporating generative models. Specifically, Retrieval Augmented Generation (RAG) is the emerging paradigm that integrates existing knowledge from large-scale document corpora into the generation process, enabling the model to generate more coherent, contextually relevant, and accurate text across various tasks. Such tasks include summarization, question answering, and dialogue systems. Recent studies have highlighted the significant positional dependence exhibited by RAG systems. Such studies observed how the placement of information within the LLM input prompt drastically affects the generated output. We ground our study on this property by investigating alternative strategies for ordering sentences within the LLM prompt to improve the average quality of the generated responses in the user and conversational system dialogues. We propose the architecture of an end-to-end RAG-based conversational assistant and empirically evaluate our strategies using the TREC CAsT 2022 collection. Our experiments highlight significant differences between distinct arrangement strategies. By employing an evaluation methodology based on RankVicuna, we show that our best approach achieves improvements up to 54% in terms of overall response quality over baseline methods.

## ACM Reference Format:

Marco Alessio, Guglielmo Faggioli, Nicola Ferro, Franco Maria Nardini, and Raffaele Perego. 2024. Improving RAG Systems via Sentence Clustering and Reordering. In . ACM, New York, NY, USA, 9 pages. <https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

## 1 INTRODUCTION

Retrieval Augmented Generation (RAG) is an emerging paradigm in the field of Artificial Intelligence (AI) to enhance the accuracy

and reliability of generative models by exploiting external data sources. In recent years, RAG has gained noteworthy importance and attention across different domains and fields [8] as it allows to combine the strengths of Information Retrieval (IR) systems and generative models to overcome each other's limitations.

RAG can improve the output of a generative model in several ways. First, it allows the generation process to be grounded on information from trusted knowledge sources incorporated in the provided prompt, thus avoiding or at least mitigating the well-known Large Language Model (LLM) hallucination problem, i.e., when the model generates contents not factually true or that do not concern the prompted text [9, 10, 41]. Second, RAG allows for continuous knowledge updates and integration of domain-specific information: the LLM can successfully respond to facts and topics not covered in its training data; moreover, it is easily adapted to different scenarios and contexts, without retraining or fine-tuning the entire model using datasets that might be unavailable or limited in scope or size. Finally, grounding the generation process on external knowledge incorporated in the input permits linking the output to verifiable external documents, thus enhancing trustworthiness and transparency [9, 10, 41].

Current RAG systems, however, suffer of some drawbacks highlighted in the literature. One of these issues originates from the notable positional sensitivity shown by LLMs. The placement of information within the input prompt significantly impacts the resulting output. Previous research [15, 32, 33] has highlighted biases towards “primacy” and “recency”, suggesting that generative models tend to prioritize information placed at the beginning or end of the input while neglecting the central portion.

In this paper, we advance over previous studies by investigating the positional bias in the context of RAG-based conversational systems. Specifically, we propose a novel strategy for arranging sentences within the input prompt of the LLM to improve the average quality of the generated responses over simpler methods. Our approach is based on the intuition that as coherent, fluent, and well-structured text are critical factors for successful communication between human beings, the same should also apply to LLMs: among all the possible arrangements of the input, those having sentences with similar meaning placed closer in the LLM prompt should generate, on average, better quality output. Therefore, we propose an end-to-end RAG architecture to test our hypothesis. The components of this architecture allow us to precisely identify which sentences are likely useful for answering user queries. To this end, we cluster sentences by their similarity and we define

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).  
Conference'17, July 2017, Washington, DC, USA

© 2024 Copyright held by the owner/author(s). Publication rights licensed to ACM.  
ACM ISBN 978-x-xxxx-xxxx-x/YY/MM  
<https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

alternative strategies for ordering them both inter and intra-cluster. In this way, we can study the effect on the generated response of these alternatives for prompting the generative LLM. To our knowledge, this is the first work that explicitly considers this aspect and allows us to fine-tune in a principled way the ordering of input sentences provided to the generative component of a RAG system. We compare our proposed approach against competitive baselines that represent the solutions employed by current RAG systems. We experimentally evaluate the performance of our proposed approach using the TREC Conversational Assistance Track (CAsT) 2022 collection [23], which allows us to compare the results that different arrangement strategies can achieve in a widely accepted Conversational Search (CS) scenario. Results highlight remarkable differences among the tested sentence placement strategies, with improvements up to 8.66% w.r.t. the best baseline and 54.94% w.r.t. random ordering.

The remainder of this work is organized as follows: Section 2 surveys the current state-of-the-art about RAG systems and quality evaluation for their responses. Section 3 details the architecture of our RAG system. Section 4 and Section 5 detail the results of an experimental analysis, which aims to highlight how the ordering of clusters and sentences affects the quality of the generated response. Finally, Section 6 draws some conclusions and outlines future directions and extensions of our research.

## 2 RELATED WORK

In the following, we survey the main works dealing with LLM positional dependencies and the difficulties of RAG systems in conciliating internal and external knowledge. Then, we analyze the challenges related to the evaluation of the quality of RAG responses and to the use of an “LLM-as-a-judge”.

### 2.1 Retrieval Augmented Generation

RAG enhances LLMs by retrieving additional information from an external knowledge source, enabling them to successfully answer queries beyond the scope of the training data. At the same time, RAG mitigates the hallucination problem, which is generating factually incorrect text, by referencing the provided external knowledge.

The RAG paradigm is organized into two main stages: retrieval and generation. Upon receiving a query from the user, the relevant information is retrieved from an external knowledge source. This task is undertaken by a standard IR pipeline that outputs a ranked list of documents. Afterwards, in the generation phase, the LLM synthesizes the response to answer the user query using the information carried by the selected documents.

Despite its clear advantages, RAG has drawbacks and limitations, which spark several challenges. First, RAG systems employ the external knowledge as their main source of information, disregarding the internal knowledge memorized within the LLM [17, 30]. This, in turn, may determine a decrease in the quality of the generated output when the provided content is not high-quality [30]. It is not uncommon for RAG to obtain worse outputs w.r.t. what the LLM can achieve in the closed-book scenario, i.e., without supplying retrieved results [30]. In this line, it has been observed that the LLM produces better results without injecting external knowledge

when the topic popularity is very high [17]. In general, state-of-the-art LLMs provide good quality responses for a wide range of questions but require assistance from an IR system when the internal knowledge of the model lacks information about the current topic. This phenomenon is likely to occur if the topic is not very popular, requires exceptional expertise, or when scaling the number of parameters of the generative model produces little to no effect [17]. Another challenge lies in the significant positional dependence [15, 32, 33] exhibited by LLMs, whereby the placement of information within the input prompt drastically affects the generated output. Prior research [15] has identified “primacy” and “recency” biases, indicating the tendency of generative models to focus toward information positioned either at the beginning or the end of the input while disregarding the central part. Therefore, the performance degrades significantly when LLMs should rely on information in the middle of its input context, showing a characteristic U-shaped performance curve [15]. This, in turn, means that most state-of-the-art generative models do not use effectively their longer contexts w.r.t. smaller and earlier counterparts. These phenomena can be observed both in open-source, e.g., Llama [35, 36] by Meta, and closed-source, e.g., GPT-4 [20] by OpenAI, models. It is not advisable to directly input all the retrieved information to the LLM for generating the response. Redundant information and very long contextual data can interfere with the generation quality, leading to repetitive, disjointed, or incoherent outputs [8]. Therefore, the retrieved content is typically further processed before being given in input to the LLM [38]. A recent work in this direction systematically examines the retrieval strategy of RAG systems [6]. The authors consider multiple retrieval factors affecting the generation process, such as the relevance of the passages in the prompt context, their position, and their number. One counter-intuitive finding is that the retriever’s highest-scoring documents that are not directly relevant to the query, e.g., do not contain the answer, negatively impact the effectiveness of the LLM. Moreover, the authors discover that adding random documents in the prompt improves the LLM accuracy by up to 35%.

In this work, we rely on the intuition that the use of coherent, fluent, and well-structured inputs can improve RAG and we propose an end-to-end architecture for selecting and structuring the external information included in the LLM prompt for response generation.

### 2.2 Quality Evaluation

Another line of research is how to evaluate the overall quality of the generation output. Despite human assessment providing the most accurate and reliable measure for evaluating model performance, the high time and cost requirements severely limit the application. Therefore, there exists an ever-increasing demand for automated evaluation techniques that consistently align with human judgements while offering enhanced efficiency and cost-effectiveness.

In this paper, we focus on textual-based generative models. Classical automatic evaluation metrics, such as BLEU [24], ROUGE [13], and METEOR [1], are designed to quantify the degree of similarity between a candidate text and one or more reference texts, by assessing their n-grams matching. The simplicity and explainability, along with the good correlation with human judgements, make these metrics widely used as baselines. However, these metrics exhibit several

limitations [40]: firstly, they cannot account for lexical diversity; secondly, they penalize variations in the semantic ordering of words; thirdly, they struggle to capture and match paraphrases effectively; lastly, they inadequately account for distant dependencies within the text. With the advent of word embeddings [18, 25] and neural models [7, 28, 35, 36, 42] based on Transformers [37], new learned metrics [2, 40] have been developed. For example, BERTScore [40] can capture the semantic similarity between the candidate and reference texts employing the contextual embeddings generated by an encoder model, such as BERT [7].

In recent years, the rapid advancements of LLMs showing remarkable performance across many tasks have gained considerable interest in their potential application also as annotators and evaluators. Due to their training using Reinforcement Learning from Human Feedback (RLHF), these models demonstrate significant human alignment. Many research have investigated leveraging state-of-the-art LLMs to automatically produce assessments serving as proxies for human judgments, a paradigm known as “LLM-as-a-judge”.

Furthermore, in recent years LLMs have gained popularity also as evaluators. For example, Zheng et al. [42] assessed the quality of conversations with various LLMs, both open and closed source, employing GPT-4 [20] as judge. They experimented with various prompts and different approaches, such as single answer grading and pairwise comparisons both between responses and against a reference text. GPT-3.5 Turbo and GPT-4 [20] have been employed as listwise rerankers [32, 33] for the TREC Deep Learning 2019 and 2020 [4, 5] and BEIR [34] experimental collections, obtaining state-of-the-art performance [32]. The same LLMs have also been employed as teacher models to fine-tune smaller open-source student models, such as Llama and Vicuna [16, 31] (i.e.: RankVicuna [27]).

In this work, we rely on state-of-the-art assessment methods and evaluate the quality of the responses generated by the different methods using RankVicuna [27].

### 3 THE PROPOSED RAG ARCHITECTURE

Generative models exhibit strong biases towards information positioned at the start or the end of the input while disregarding the middle part [15]. This phenomenon motivates our research effort to determine how the order of the input sentences provided to a RAG-based conversational system affects the quality of the generated output and, in turn, the optimal ordering strategy to achieve the best response. This section describes each method and all variations considered in our experiments.

The architecture of our proposed RAG system is illustrated in Figure 1. It includes an IR pipeline, which retrieves top- $k$  documents  $D = \{d_1, d_2, \dots, d_k\}$  in response to each user utterance  $q$ . The retrieved documents are then processed by additional components responsible for splitting them into sentences, identifying the most relevant sentences, clustering such sentences based on their semantic similarity, and ordering them according to the various strategies analyzed. Finally, the selected—re-ordered—sentences are provided as input to the LLM for response generation. These components are the focus of our research. Their functionalities are detailed in the remainder of this section.

#### 3.1 Document Pre-processing and Splitting

As observed in literature [11, 29], the entire text of a relevant document rarely contains meaningful knowledge to satisfy the user information need expressed by a query  $q$ . In most cases, only one or a few portions of the document are relevant to the query, while the remaining parts contain irrelevant information. The proposed architecture aims to precisely identify the key information in the retrieved documents, i.e., the sentences, to reduce the noise in the prompt used for response generation.

Hereinafter, we consider sentences in the documents as the atomic units of information. Our pipeline, illustrated in Figure 1 works as follows. First, for each query  $q$  we consider only the top- $k$  documents  $\{d_1, d_2, \dots, d_k\}$  retrieved by the IR system. Then, a state-of-the-art co-reference resolution model is applied to all documents to replace pronouns and other generic terms within a sentence with the fully specified entity mentioned in a previous sentence. This allows us to remove the contextual dependencies among sentences in a document so they can be considered self-explanatory. The third step splits each document  $d_i$  into a sequence of sentences  $\{s_{i,1}, s_{i,2}, \dots, s_{i,n_i}\}$ . Afterwards, near-duplicate removal is employed to the sentences originated by all documents by discarding sentences with a Jaccard similarity  $\geq 0.9$  between their Bag-of-Words (BoW) representations<sup>1</sup>.

#### 3.2 Sentence Selection

After the first pre-processing phase, we obtain a sentence candidate set for each query to be included in the LLM prompt of our RAG system (see Figure 1). Since the cardinality of this set can be large and not all the sentences are useful for answering the query, we employ the BERT-based cross encoder answer-in-the-sentence classifier<sup>2</sup> developed by Lajewska and Balog [12] to rank the candidate sentences according to their predicted usefulness to (at least partially) answer the query and we retain the top- $n$  ranked sentences thus discarding the remaining ones. As a possible limitation, please note that the model by Lajewska and Balog [12] employed have been trained on queries and passages used in our experiments. Therefore, it is very likely that the model performs significantly better on our data w.r.t. any other model, ensuring that top-ranked sentences are indeed relevant to the query. Even though such a model is not available in a real practical scenario, this choice is justified by our research effort being focused exclusively on comparing the ordering strategy for sentences in the LLM input rather than on the absolute results achievable by our RAG system.

#### 3.3 Sentence Clustering and Ordering

The previous steps of the pipeline constrain the number of sentences per query while increasing their expected utility in answering the query. Furthermore, they allow us to control other noise sources, such as the number or the variable length of the retrieved documents. Therefore, we can assess how the positional bias affects the generation process. We highlight again that the positional bias of

<sup>1</sup>This step is particularly important in our setting because the CAsT 2022 corpus contains a multitude of near-duplicate documents. In particular, the same Wikipedia article is often replicated in documents retrieved from the KILT and MS-MARCO collections.

<sup>2</sup>The model named “squad\_snippets\_unanswerable” is available at [https://iai.group/downloads/emnlp2023-answerability\\_prediction](https://iai.group/downloads/emnlp2023-answerability_prediction).

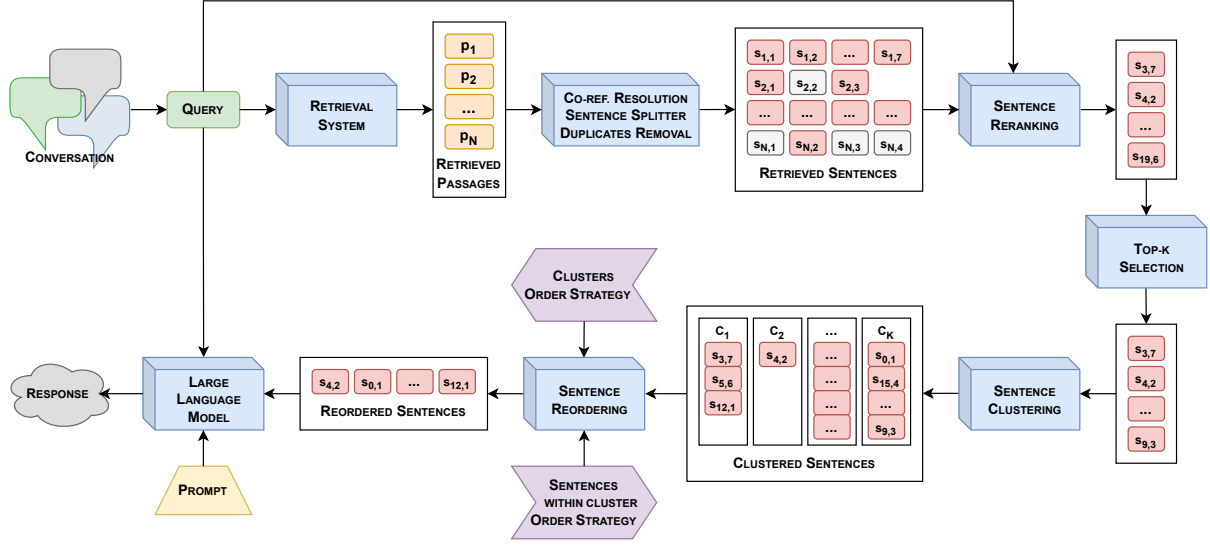


Figure 1: Architecture of our proposed RAG system.

LLM has already been observed in prior research [15, 32, 33]. However, it has been considered exclusively as a limitation of LLMs and RAG systems. Our research moves a step forward by investigating the best ordering strategy to maximize, on average, the quality of the generated responses over a testing query set  $Q$ . We believe that logically organized text where sentences with akin meanings are positioned closer in the LLM prompt should, on average, yield superior output quality. Consequently, our sentence ordering strategies exploit the similarities among sentences selected by the sentence selection step. To measure semantic inter-sentence similarity, we resort to the contextualized embeddings generated with the *tct-colbert* model<sup>3</sup> [14]. We generate the representation of the  $n$  selected sentences for each query and measure their pair-wise cosine similarity. Then, we progressively aggregate the most similar sentences by employing a hierarchical clustering algorithm. The maximum value of Silhouette statistic is used as the criteria to determine the optimal clustering among all possible. As a result, for each query  $q \in Q$ , the top- $n$  sentences are grouped in a variable number  $N_c \geq 1$  of clusters, each composed of one or more sentences with similar semantic meaning. To devise different strategies for ordering input sentences, we leverage the above clustering that allows us to study the impact of sentence placement variations occurring in both inter and intra-clusters.

More formally, given a query, the set  $S$  of the  $n$  previously selected sentences, and the prompt  $p$ , we aim to find the ordering  $ord^*$  of  $S$  such that:

$$ord^* = \underset{ord}{\operatorname{argmax}} \sum_{q \in Q} s(q, LLM(p, q, ord(S))),$$

where  $ord(S)$  is a sentence ordering strategy that returns an ordering of the sentences in  $S$ ,  $LLM(p, q, ord(S))$  is the response generated by the LLM used for prompt  $p$ , query  $q$  and sentence ordering

$ord(S)$ , and, finally,  $s(q, r)$  is a scoring function evaluating the perceived quality of the generated response  $r = LLM(p, q, ord(S))$  for query  $q$ .

The order of clusters and the order of the sentences within the same cluster uniquely determine the possible global ordering of the  $n$  sentences we consider for inputting the LLM. Our experimental assessment will evaluate six different ordering strategies for placing the clusters of sentences in the input, and four different methods for ordering sentences within the same cluster. Cluster placements consider different aspects, such as the clusters' cardinality and similarity to the query. The ordering tested includes the random one and those obtained by decreasing/increasing the value of each aspect. Finally, the U-shaped order suggested in [15] is also tested. Regarding the ordering within clusters, we consider random order, order by reranker score, visiting order, and the clustering aggregation order.

## 4 EXPERIMENTAL EVALUATION

We can now formulate the research questions we aim to answer with our experimental framework.

**Research Questions.** Given the sentence selection and clustering steps discussed above, the two main aspects to consider for defining our ordering strategies  $ord(\cdot)$  are the order of placement in the LLM prompt of the clusters and of the sentences within the same cluster. They uniquely determine the global ordering  $ord(\cdot)$  of the top- $n$  sentences given in input to the LLM for response generation. Our research questions assess which is the best solution among these alternatives considered. Specifically,

RQ1 What is the best cluster ordering strategy?

RQ2 What is the best ordering strategy for sentences within the same cluster?

RQ3 Can our proposed strategy enhance the effectiveness of the RAG system w.r.t. baseline methods?

<sup>3</sup>[https://huggingface.co/castorini/tct\\_colbert-v2-hnp-msmarco](https://huggingface.co/castorini/tct_colbert-v2-hnp-msmarco)

**Experimental Settings.** We experiment with the TREC CAsT 2022 dataset, a standard experimental collection for CS [23]. This choice is due to prior research that released additional datasets, models, and human judgments for this benchmark [11, 12]. The corpus is composed of three documents collections, MS-MARCO v2 [19], KILT [26], and Washington Post v4, which are subdivided into 106M short documents. CAsT 2022 includes 18 information needs (topics) and 205 user utterances (queries), with an average length of 11.39 user utterances per topic. The number of utterances for which relevance judgements are provided is 163.

For our experiments, as the retrieval system, we employ as the output of the retrieval pipeline the best-performing run originally submitted to TREC CAsT 2022<sup>4</sup> [39]. This allows us to focus exclusively on the following steps of our pipeline. In all our experiments, we consider only the top-20 retrieved documents, leaving the investigation about the implications of this choice and possible alternatives as future work. To provide meaningful results, all queries where  $Precision@20 < 0.2$ , that is, having at most 3 relevant passages in the top-20 results, are discarded<sup>5</sup>, ensuring that enough relevant information is retrieved to answer the considered queries successfully.

Furthermore, in the steps of the pipeline where the query text is needed, i.e., sentence ranking and response generation, we employed the manually rewritten text for every query. This allows us to account for the possible bias introduced by different query rewriting approaches. Future developments will investigate the relationship between query rewriting approaches and RAG solutions.

For co-reference resolution at the document level, i.e., removing co-references across different sentences in the “document processing” step, we use the “F-Coref” model<sup>6</sup> [21] based on the “LingMess” architecture [22]. After this step, we use the well-known SpaCy Python library to divide each document into a sequence of independent sentences.

In the following section, we report two different metrics for each comparison. The former is the average score of every approach when assessing all 10 random permutations using RankVicuna. The latter, instead, is a pairwise metric, assessing the number of queries for which the first approach obtains higher/the same/lower score w.r.t. the other one. This information should better highlight the differences and provide a more comprehensive view than a single average value.

**Response Generation.** For the response generation, we employ Vicuna 7B<sup>7</sup> [42], a LLM based on Llama 2 [35, 36] fine-tuned on 125K user conversations with ChatGPT gathered using public APIs from the ShareGPT.com website. To ensure the reproducibility of our experiments, we set the temperature of the LLM to 0.

**Quality Evaluation.** To evaluate the quality of the generated responses, we employ RankVicuna [27] to perform listwise ranking between all responses being compared. To mitigate the positional bias intrinsic in RankVicuna, we assess 10 different random permutations of the same responses, averaging the results obtained.

<sup>4</sup>The run is identified as “udinfo\_mi\_b2021” from the “udel\_fang” group, University of Delaware (USA)

<sup>5</sup>The number of queries considered in these experiments is 115 out of 163 evaluated in the official relevance judgments.

<sup>6</sup><https://huggingface.co/biu-nlp/f-coref>

<sup>7</sup><https://huggingface.co/lmsys/vicuna-7b-v1.5>

**Table 1: Comparisons between the six approaches proposed for RQ1: “What is the best ordering strategy for clusters?”. In the top half, each row reports three numbers, which are the wins for the approach in the column label, the ties, and the wins for the approach in the row label, respectively. In the bottom half, the overall results are reported.**

	A vs.	B vs.	C vs.	D vs.	E vs.	F vs.
<b>A</b>	—	56-4-51	57-2-52	<b>62-4-45</b>	51-1-59	55-2-54
<b>B</b>	51-4-56	—	47-8-56	<b>61-4-46</b>	52-5-54	52-2-57
<b>C</b>	52-2-57	56-8-47	—	<b>58-3-50</b>	55-0-56	59-4-48
<b>D</b>	45-4-62	46-4-61	50-3-58	—	44-1-66	47-1-63
<b>E</b>	59-1-51	54-5-52	56-0-55	<b>66-1-44</b>	—	57-5-49
<b>F</b>	54-2-55	57-2-52	48-4-59	<b>63-1-47</b>	49-5-57	—
<b>Overall</b>	261-13	269-23	258-17	<b>310-13</b>	251-12	270-14
<b>Avg. Score</b>	0.5723	0.5844	0.5510	<b>0.6219</b>	0.5736	0.5969

This is a reasonable trade-off between evaluation accuracy and the computational runtime required. For each assessment, we assign  $\frac{N+1-i}{N}$  points to the  $i$ -th ranked response, where  $1 \leq i \leq N$  and  $N$  is the number of responses being compared. Furthermore, we also evaluate the number of wins and ties between pairs of responses considered. Whether a valid judgment from the LLM can not be determined, the entire comparison is discarded from the evaluation.

#### 4.1 RQ1: Order of Clusters

For the first experiment, we evaluate the effects of different ordering of the clusters while keeping the order of sentences within the same cluster (based on the clustering aggregation order) fixed. We test six different strategies for ordering clusters: clusters selected in random order (strategy A); clusters selected in descending order of cardinality (strategy B); clusters selected in ascending order of similarity with the query<sup>8</sup> (strategy C); clusters selected in descending order of similarity with the query (strategy D); clusters selected in descending order by similarity with the query using a ping-pong layout from top to bottom (strategy E)<sup>9</sup>; clusters selected by similarity with the query in descending order, using a ping-pong layout from bottom to top (strategy F)<sup>10</sup>.

As shown in Table 1, sorting the clusters in descending order by their similarity with the query (strategy D) is the clear winner in this comparison, in terms of both score and pairwise wins. This approach performs 18.77%, 15.24%, 20.16%, 23.51%, and 14.81% better than other options. This figures suggest that the LLM used to generate the responses exhibit a much stronger “primacy” rather than “recency” biases, as highlighted by option C being overall the worst performing among those considered. Instead, methods E and F were designed to place the least important clusters towards the center, since LLMs struggle to utilize the information in the middle of their prompt effectively. However, we can see that both approaches are ineffective: we suspect this is due to the length of

<sup>8</sup>The similarity between a cluster  $C$  and the query is defined as the maximum cosine similarity between the query  $q \in Q$  with any sentence  $s_{i,j} \in C$  belonging to the cluster.

<sup>9</sup>The clusters are placed first, last, second, second-to-last, third, and so on, e.g., [A, B, C, D, E] becomes [A, C, E, D, B].

<sup>10</sup>The clusters are placed last, first, second-to-last, second, third-to-last, and so on, e.g., [A, B, C, D, E] becomes [B, D, E, C, A].

**Table 2: Comparisons between the four approaches proposed for RQ2: “What is the best ordering strategy for sentences within the same cluster?”. In the top half, each row reports three numbers, which are the wins for the approach in the column label, the ties, and the wins for the approach in the row label, respectively. In the bottom half, the overall results are reported.**

	A vs.	B vs.	C vs.	D vs.
<b>A</b>	—	53-3-59	48-8-59	<b>55-4-56</b>
<b>B</b>	59-3-53	—	54-3-58	<b>60-7-48</b>
<b>C</b>	59-8-48	58-3-54	—	<b>57-6-52</b>
<b>D</b>	56-4-55	48-7-60	52-6-57	—
<b>Overall</b>	<b>174-15</b>	159-13	154-17	<b>172-17</b>
<b>Avg. Score</b>	0.6281	0.6143	0.6124	<b>0.6451</b>

the input text being much smaller than the maximum context window of the model. Different results may be observed when varying the amount of input data provided to the LLM for generation.

#### 4.2 RQ2: Order of Sentences within the same Cluster

In this second experiment, we evaluate different sorting schemes for sentences within the same cluster, keeping the cluster’s order fixed at the best strategy determined in RQ1. We test four different strategies for ordering sentences within the same cluster: sentences selected in random order (strategy A); sentences selected in descending order by re-ranker score (strategy B); sentences selected by visiting order<sup>11</sup> (strategy C); sentences selected by aggregation order (strategy D).

As shown in Table 2, the best results are achieved by two different strategies: option D, sorting sentences within the same cluster based on aggregation order, and interestingly, option A, randomly sorting the sentences. Both strategies are preferable to the other two methods considered, performing 8.18% and 11.69% better w.r.t. options B and C, respectively. We note however that the difference in performance of the various strategies are not large as the sentences are grouped in the clusters by their similarity. The LLM response appears to be more impacted by the order of the clusters than by the order of sentences within each cluster.

#### 4.3 RQ3: Comparison with Baselines

Our last experiment investigates whether our proposed approach is beneficial in enhancing the overall effectiveness of the RAG system w.r.t. four simpler baseline methods that may be used in practice by current state-of-the-art RAG systems. We test five different strategies: i) the top-5 retrieved documents (A), ii) the top-40 sentences taken in random order (B), iii) the top-40 sentences taken in descending order by re-ranker score (C), iv) the top-40 sentences selected by visiting order (D), v) the best clusterization-based approach determined from RQ1 and RQ2 (CL).

<sup>11</sup>The sentences are sorted based on the order in which they appear when sequentially scanning through the set of top- $k$  retrieved documents.

The results obtained are shown in Table 3. The clusterization-based approach demonstrate superior performance, resulting as the best strategy in this comparison. The four baselines yield notably lower results: 15.14%, 54.94%, 8.66%, and 15.67%, respectively. Among the methods considered in this work, randomly sorting the top- $h$  sentences is by far the least performing approach. This, in turn, proves our starting intuition about coherent, fluent, and well-structured text being critical factors for LLMs to generate high quality output.

### 5 ADDITIONAL EXPERIMENTS

The clusterization-based ordering strategy proposed in this work is designed to position sentences sharing analogous semantic content close together in the LLM prompt. Given the results obtained in Section 4.3, we have shown its effectiveness in our experimental settings. Nevertheless, we answer two additional research questions in this section to gain additional insights. Specifically,

RQ4 Is there a correlation between the similarity of subsequent sentences in the LLM prompt and the quality of the generated response?

RQ5 Is the proposed clusterization strategy more effective than directly optimising the similarity of subsequent sentences?

**Experimental Settings.** We determine heuristically the two ordering  $ord^+$  and  $ord^-$ , which maximize and minimize the overall similarity between subsequent sentences. Let  $sum^+$  and  $sum^-$  be the sum of similarity between subsequent sentences for  $ord^+$  and  $ord^-$  respectively. The similarity  $sim(p)$  for a sentence permutation  $p$  is given by the following equation, where min-max normalization is used, and  $s_i$  are the embedding representations of the respective sentences:

$$sim(p) = \frac{\left(\sum_{i=2}^h \cos(s_{i-1}, s_i)\right) - sum^-}{sum^+ - sum^-}$$

In our experiments, for each query, we generate one million random permutations, then we determine which is the permutation with similarity closer to each of the following thresholds: 0.125, 0.250, 0.375, 0.500, and 0.625. We decided to stop at 0.625 because

**Table 3: Comparisons between the five approaches considered for RQ3: “Can our proposed strategy enhance the effectiveness of the RAG system w.r.t. baseline methods?”. In the top half, each row reports three numbers, which are the wins for approach in the column label, the ties, and the wins for approach in the row label, respectively. In the bottom half, the overall results are reported.**

	A vs.	B vs.	C vs.	D vs.	CL vs.
<b>A</b>	—	45-4-62	54-1-56	54-0-57	<b>66-2-43</b>
<b>B</b>	62-4-45	—	71-1-39	64-8-39	<b>67-5-39</b>
<b>C</b>	56-1-54	39-1-71	—	50-4-57	<b>59-3-49</b>
<b>D</b>	57-0-54	39-8-64	57-4-50	—	<b>59-3-49</b>
<b>CL</b>	43-2-66	39-5-67	49-3-59	49-3-59	—
<b>Overall</b>	<b>218-7</b>	162-18	231-9	217-15	<b>251-13</b>
<b>Avg. Score</b>	0.5882	0.5533	0.6177	0.6016	<b>0.6392</b>

**Table 4: Comparisons between the seven approaches proposed for RQ4: “Is there a correlation between the similarity of subsequent sentences in the LLM prompt and the quality of the generated response?”. In the top half, each row reports three numbers, which are the wins for the approach in the column label, the ties, and the wins for the approach in the row label, respectively. In the bottom half, the overall results are reported.**

	1.000 vs.	0.625 vs.	0.500 vs.	0.375 vs.	0.250 vs.	0.125 vs.	0.000 vs.
<b>1.000</b>	—	46-2-43	38-2-51	45-0-46	40-2-49	40-1-50	38-1-52
<b>0.625</b>	43-2-46	—	37-2-52	42-2-47	41-1-49	36-0-55	35-1-55
<b>0.500</b>	51-2-38	52-2-37	—	51-2-38	52-0-39	37-0-54	44-2-45
<b>0.375</b>	46-0-45	47-2-42	38-2-51	—	42-2-47	37-3-51	37-1-53
<b>0.250</b>	49-2-40	49-1-41	39-0-52	47-2-42	—	43-3-45	42-1-48
<b>0.125</b>	50-1-40	55-0-36	54-0-37	51-3-37	45-3-43	—	44-1-46
<b>0.000</b>	52-1-38	55-1-35	45-2-44	53-1-37	48-1-42	46-1-44	—
<b>Overall</b>	291-8	<b>304-8</b>	251-8	289-10	268-9	239-8	240-7
<b>Avg. Score</b>	0.5731	<b>0.5866</b>	0.5480	0.5617	0.5516	0.5349	0.5143

higher values are unlikely to be observed given that the average similarity of these permutations is 0.3433 with standard deviation 0.0530.

**Results.** We determine how the quality of the generated response is influenced when varying the similarity between subsequent sentences at various predefined thresholds, as shown in Table 4. It is interesting to note that the highest results are obtained by permutations with 0.625 normalised similarity, rather than 1.000 which is the ordering maximising the similarity between subsequent sentences ( $ord^+$ ). This method achieves 4.47% and 26.67% more pairwise wins w.r.t.  $ord^+$  and  $ord^-$ , respectively. To answer RQ5, we assess the responses generated using the best clustering strategy against the approach defined above. The average scores are 0.7652 and 0.7348 while the pairwise wins and ties are 38 - 46 - 31, respectively.

From these experiments, we can conclude that a positive correlation exists between similarity between subsequent sentences and response quality, while proving that sentence similarity may not be the only factor that should be considered. Moreover, subdividing and explicitly grouping together sentences by subtopic is beneficial w.r.t. considering the sentence similarity only in a pairwise fashion and thus lacking a global vision of the retrieved knowledge.

## 6 CONCLUSIONS AND FUTURE WORK

In this work, we presented a novel pipelined RAG architecture aimed at selecting a set of relevant sentences for each query and arranging them in a specific order to optimize the quality of responses generated by a LLM. For this purpose, sentences are first extracted from the top documents retrieved. Then, they are reranked, and the most relevant sentences are organized in clusters by similarity. We proposed different strategies for ordering clusters and the sentences within clusters in the input given to the LLM for response generation. To the best of our knowledge, this is the first work investigating sentence clustering and re-ordering to improve the quality of the response generated by RAG systems. Our empirical assessment is based on a well-known—public—framework for conversational search. The results of the experiments show that different sequences of sentences in the LLM prompt significantly impact response quality despite all methodologies processing identical

information from the same set of sentences. Random permutations yield the lowest results, whereas our proposed approach based on sentence clusterization yields superior results. Additionally, we examined whether maximizing the similarity between consecutive sentences in the LLM prompt enhances response quality. While a positive correlation between these factors was observed, it is not the exclusive determinant. Consequently, while we infer that sentence similarity constitutes a pivotal aspect, other contributing factors remain unidentified, warranting further investigation. Moreover, although our experimental evaluation employs a well-known conversational collection, the methodology and results shown in this work are general. They could also be applied to other scenarios, such as ad-hoc search.

In future work, we intend to evaluate the impact of the number of clusters selected by our method for generating the response. Our intuition is that the number of clusters identified for a given query is a proxy of the difficulty of the query itself. Fewer clusters or even a single large should characterize simple and close queries. In contrast, difficult—multi-faceted—queries are possibly characterized by more clusters, each addressing a different facet of the query. This intuition paves the way for the extension of the evaluation methodology by adopting diversification-based metrics [3], allowing us to understand how well the generated answers cover the query facets and the topical distribution of the clusters.

## REFERENCES

- [1] Satanjeev Banerjee and Alon Lavie. 2005. METEOR: An Automatic Metric for MT Evaluation with Improved Correlation with Human Judgments. In *Proceedings of the Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization@ACL 2005, Ann Arbor, Michigan, USA, June 29, 2005*, Jade Goldstein, Alon Lavie, Chin-Yew Lin, and Clare R. Voss (Eds.). Association for Computational Linguistics, 65–72. <https://aclanthology.org/W05-0909/>
- [2] Elizabeth Clark, Shruti Rijhwani, Sebastian Gehrmann, Joshua Maynez, Roei Aharoni, Vitaly Nikolaev, Thibault Sellam, Aditya Siddhant, Dipanjan Das, and Ankur P. Parikh. 2023. SEAHORSE: A Multilingual, Multifaceted Dataset for Summarization Evaluation. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6–10, 2023*, Houda Bouamor, Juan Pino, and Kalika Bali (Eds.). Association for Computational Linguistics, 9397–9413. <https://doi.org/10.18653/V1/2023.EMNLP-MAIN.584>
- [3] Charles L.A. Clarke, Maheedhar Kolla, Gordon V. Cormack, Olga Vechtomova, Azin Ashkan, Stefan Büttcher, and Ian MacKinnon. 2008. Novelty and diversity in information retrieval evaluation. In *Proceedings of the 31st Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (Singapore, Singapore) (SIGIR '08)*. Association for Computing Machinery, New York, NY, USA, 659–666. <https://doi.org/10.1145/1390334.1390446>
- [4] Nick Craswell, Bhaskar Mitra, Emine Yilmaz, and Daniel Campos. 2020. Overview of the TREC 2020 Deep Learning Track. In *Proceedings of the Twenty-Ninth Text Retrieval Conference, TREC 2020, Virtual Event [Gaithersburg, Maryland, USA], November 16–20, 2020 (NIST Special Publication, Vol. 1266)*, Ellen M. Voorhees and Angela Ellis (Eds.). National Institute of Standards and Technology (NIST). <https://trec.nist.gov/pubs/trec29/papers/OVERVIEW.DL.pdf>
- [5] Nick Craswell, Bhaskar Mitra, Emine Yilmaz, Daniel Campos, and Ellen M. Voorhees. 2020. Overview of the TREC 2019 deep learning track. *CoRR abs/2003.07820* (2020). [arXiv:2003.07820](https://arxiv.org/abs/2003.07820) <https://arxiv.org/abs/2003.07820>
- [6] Florin Cuconasu, Giovanni Trappolini, Federico Siciliano, Simone Filice, Cesare Campagnano, Yoelle Maarek, Nicola Tonellotto, and Fabrizio Silvestri. 2024. The power of noise: Redefining retrieval for rag systems. *arXiv preprint arXiv:2401.14887* (2024).
- [7] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2–7, 2019, Volume 1 (Long and Short Papers)*, Jill Burstein, Christy Doran, and Thamar Solorio (Eds.). Association for Computational Linguistics, 4171–4186. <https://doi.org/10.18653/V1/N19-1423>
- [8] Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, Meng Wang, and Haofen Wang. 2024. Retrieval-Augmented Generation



- for Large Language Models: A Survey. arXiv:2312.10997 [cs.CL]
- [9] Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, and Ting Liu. 2023. A Survey on Hallucination in Large Language Models: Principles, Taxonomy, Challenges, and Open Questions. *CoRR* abs/2311.05232 (2023). <https://doi.org/10.48550/ARXIV.2311.05232> arXiv:2311.05232
  - [10] Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Yejin Bang, Andrea Madotto, and Pascale Fung. 2023. Survey of Hallucination in Natural Language Generation. *ACM Comput. Surv.* 55, 12 (2023), 248:1–248:38. <https://doi.org/10.1145/3571730>
  - [11] Weronika Lajewska and Krisztian Balog. 2023. Towards Filling the Gap in Conversational Search: From Passage Retrieval to Conversational Response Generation. In *Proceedings of the 32nd ACM International Conference on Information and Knowledge Management, CIKM 2023, Birmingham, United Kingdom, October 21-25, 2023*, Ingo Frommholz, Frank Hopfgartner, Mark Lee, Michael Oakes, Mounia Lalmas, Min Zhang, and Rodrygo L. T. Santos (Eds.). ACM, 5326–5330. <https://doi.org/10.1145/3583780.3615132>
  - [12] Weronika Lajewska and Krisztian Balog. 2024. Towards Reliable and Factual Response Generation: Detecting Unanswerable Questions in Information-Seeking Conversations. In *Advances in Information Retrieval - 46th European Conference on Information Retrieval, ECIR 2024, Glasgow, UK, March 24-28, 2024, Proceedings, Part III (Lecture Notes in Computer Science, Vol. 14610)*, Nazli Goharian, Nicola Tonello, Yulan He, Aldo Lipani, Graham McDonald, Craig Macdonald, and Iadh Ounis (Eds.). Springer, 336–344. [https://doi.org/10.1007/978-3-031-56063-7\\_25](https://doi.org/10.1007/978-3-031-56063-7_25)
  - [13] Chin-Yew Lin. 2004. ROUGE: A Package for Automatic Evaluation of Summaries. In *Text Summarization Branches Out*. Association for Computational Linguistics, Barcelona, Spain, 74–81. <https://aclanthology.org/W04-1013>
  - [14] Sheng-Chieh Lin, Jheng-Hong Yang, and Jimmy Lin. 2021. In-Batch Negatives for Knowledge Distillation with Tightly-Coupled Teachers for Dense Retrieval. In *Proceedings of the 6th Workshop on Representation Learning for NLP, Repl4NLP@ACL-IJCNLP 2021, Online, August 6, 2021*, Anna Rogers, Iacer Calixto, Ivan Vulic, Naomi Saphra, Nora Kassner, Oana-Maria Camburu, Trapit Bansal, and Vered Shwartz (Eds.). Association for Computational Linguistics, 163–173. <https://doi.org/10.18653/V1/2021.REPL4NLP-1.17>
  - [15] Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2023. Lost in the Middle: How Language Models Use Long Contexts. *CoRR* abs/2307.03172 (2023). <https://doi.org/10.48550/ARXIV.2307.03172> arXiv:2307.03172
  - [16] Xueguang Ma, Xinyu Zhang, Ronak Pradeep, and Jimmy Lin. 2023. Zero-Shot Listwise Document Reranking with a Large Language Model. *CoRR* abs/2305.02156 (2023). <https://doi.org/10.48550/ARXIV.2305.02156> arXiv:2305.02156
  - [17] Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das, Daniel Khashabi, and Hananeh Hajishirzi. 2023. When Not to Trust Language Models: Investigating Effectiveness of Parametric and Non-Parametric Memories. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023*, Anna Rogers, Jordan L. Boyd-Graber, and Naoaki Okazaki (Eds.). Association for Computational Linguistics, 9802–9822. <https://doi.org/10.18653/V1/2023.ACL-LONG.546>
  - [18] Tomáš Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Efficient Estimation of Word Representations in Vector Space. In *1st International Conference on Learning Representations, ICLR 2013, Scottsdale, Arizona, USA, May 2-4, 2013, Workshop Track Proceedings*, Yoshua Bengio and Yann LeCun (Eds.). <http://arxiv.org/abs/1301.3781>
  - [19] Tri Nguyen, Mir Rosenberg, Xia Song, Jianfeng Gao, Saurabh Tiwary, Rangan Majumder, and Li Deng. 2016. MS MARCO: A Human Generated Machine Reading Comprehension Dataset. In *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 (CEUR Workshop Proceedings, Vol. 1773)*, Tarek Richard Besold, Antoine Bordes, Artur S. d'Ávila Garcez, and Greg Wayne (Eds.). CEUR-WS.org. [http://ceur-ws.org/Vol-1773/CoCoNIPS\\_2016\\_paper9.pdf](http://ceur-ws.org/Vol-1773/CoCoNIPS_2016_paper9.pdf)
  - [20] OpenAI. 2023. GPT-4 Technical Report. *CoRR* abs/2303.08774 (2023). <https://doi.org/10.48550/ARXIV.2303.08774> arXiv:2303.08774
  - [21] Shon Otmazgin, Arie Cattán, and Yoav Goldberg. 2022. F-coref: Fast, Accurate and Easy to Use Coreference Resolution. In *Proceedings of the 2nd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 12th International Joint Conference on Natural Language Processing, AACL/IJCNLP 2022 - System Demonstrations, Taipei, Taiwan, November 20 - 23, 2022*. Association for Computational Linguistics, 48–56. <https://aclanthology.org/2022.aacl-demo.6>
  - [22] Shon Otmazgin, Arie Cattán, and Yoav Goldberg. 2023. LingMess: Linguistically Informed Multi Expert Scorers for Coreference Resolution. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics, EACL 2023, Dubrovnik, Croatia, May 2-6, 2023*, Andreas Vlachos and Isabelle Augenstein (Eds.). Association for Computational Linguistics, 2744–2752. <https://doi.org/10.18653/V1/2023.EACL-MAIN.202>
  - [23] Paul Owoicho, Jeff Dalton, Mohammad Aliannejadi, Leif Azzopardi, Johanne R. Trippas, and Svitlana Vakulenko. 2022. TREC CAsT 2022: Going Beyond User Ask and System Retrieve with Initiative and Response Generation. In *Proceedings of the Thirty-First Text REtrieval Conference, TREC 2022, online, November 15-19, 2022 (NIST Special Publication, Vol. 500-338)*, Ian Soboroff and Angela Ellis (Eds.). National Institute of Standards and Technology (NIST). [https://trec.nist.gov/pubs/trec31/papers/Overview\\_cast.pdf](https://trec.nist.gov/pubs/trec31/papers/Overview_cast.pdf)
  - [24] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a Method for Automatic Evaluation of Machine Translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics, July 6-12, 2002, Philadelphia, PA, USA*. ACL, 311–318. <https://doi.org/10.3115/1073083.1073135>
  - [25] Jeffrey Pennington, Richard Socher, and Christopher D. Manning. 2014. Glove: Global Vectors for Word Representation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing, EMNLP 2014, October 25-29, 2014, Doha, Qatar, A meeting of SIGDAT, a Special Interest Group of the ACL*, Alessandro Moschitti, Bo Pang, and Walter Daelemans (Eds.). ACL, 1532–1543. <https://doi.org/10.3115/V1/D14-1162>
  - [26] Fabio Petroni, Aleksandra Piktus, Angela Fan, Patrick S. H. Lewis, Majid Yazdani, Nicola De Cao, James Thorne, Yacine Jernite, Vladimir Karpukhin, Jean Maillard, Vassilis Plachouras, Tim Rocktäschel, and Sebastian Riedel. 2021. KILT: a Benchmark for Knowledge Intensive Language Tasks. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2021, Online, June 6-11, 2021*, Kristina Toutanova, Anna Rumshisky, Luke Zettlemoyer, Dilek Hakkani-Tür, Iz Beltagy, Steven Bethard, Ryan Cotterell, Tanmoy Chakraborty, and Yichao Zhou (Eds.). Association for Computational Linguistics, 2523–2544. <https://doi.org/10.18653/V1/2021.NAACL-MAIN.200>
  - [27] Ronak Pradeep, Sahel Sharifmoghaddam, and Jimmy Lin. 2023. RankVicuna: Zero-Shot Listwise Document Reranking with Open-Source Large Language Models. *CoRR* abs/2309.15088 (2023). <https://doi.org/10.48550/ARXIV.2309.15088> arXiv:2309.15088
  - [28] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer. *J. Mach. Learn. Res.* 21 (2020), 140:1–140:67. <http://jmlr.org/papers/v21/20-074.html>
  - [29] Pengjie Ren, Zhumin Chen, Zhaochun Ren, Evangelos Kanoulas, Christof Monz, and Maarten de Rijke. 2021. Conversations with Search Engines: SERP-based Conversational Response Generation. *ACM Trans. Inf. Syst.* 39, 4 (2021), 47:1–47:29. <https://doi.org/10.1145/3432726>
  - [30] Ruiyang Ren, Yuhao Wang, Yingqi Qu, Wayne Xin Zhao, Jing Liu, Hao Tian, Hua Wu, Ji-Rong Wen, and Haifeng Wang. 2023. Investigating the Factual Knowledge Boundary of Large Language Models with Retrieval Augmentation. *CoRR* abs/2307.11019 (2023). <https://doi.org/10.48550/ARXIV.2307.11019> arXiv:2307.11019
  - [31] Weiwei Sun, Zheng Chen, Xinyu Ma, Lingyong Yan, Shuaiqiang Wang, Pengjie Ren, Zhumin Chen, Dawei Yin, and Zhaochun Ren. 2023. Instruction Distillation Makes Large Language Models Efficient Zero-shot Rankers. *CoRR* abs/2311.01555 (2023). <https://doi.org/10.48550/ARXIV.2311.01555> arXiv:2311.01555
  - [32] Weiwei Sun, Lingyong Yan, Xinyu Ma, Shuaiqiang Wang, Pengjie Ren, Zhumin Chen, Dawei Yin, and Zhaochun Ren. 2023. Is ChatGPT Good at Search? Investigating Large Language Models as Re-Ranking Agents. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, Houda Bouamor, Juan Pino, and Kalika Bali (Eds.). Association for Computational Linguistics, 14918–14937. <https://doi.org/10.18653/V1/2023.EMNLP-MAIN.923>
  - [33] Raphael Tang, Xinyu Zhang, Xueguang Ma, Jimmy Lin, and Ferhan Ture. 2023. Found in the Middle: Permutation Self-Consistency Improves Listwise Ranking in Large Language Models. *CoRR* abs/2310.07712 (2023). <https://doi.org/10.48550/ARXIV.2310.07712> arXiv:2310.07712
  - [34] Nandan Thakur, Nils Reimers, Andreas Rücklé, Abhishek Srivastava, and Iryna Gurevych. 2021. BEIR: A Heterogenous Benchmark for Zero-shot Evaluation of Information Retrieval Models. *CoRR* abs/2104.08663 (2021). [arXiv:2104.08663](https://arxiv.org/abs/2104.08663)
  - [35] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. LLaMA: Open and Efficient Foundation Language Models. *CoRR* abs/2302.13971 (2023). <https://doi.org/10.48550/ARXIV.2302.13971> arXiv:2302.13971
  - [36] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shriti Bhoale, Dan Bikel, Lukas Blecher, Cristian Canton-Ferrer, Moya Chen, Guillem Cucurull, David Esiohu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross



- Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurélien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open Foundation and Fine-Tuned Chat Models. *CoRR* abs/2307.09288 (2023). <https://doi.org/10.48550/ARXIV.2307.09288> arXiv:2307.09288
- [37] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is All you Need. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett (Eds.). 5998–6008. <https://proceedings.neurips.cc/paper/2017/hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html>
- [38] Fangyuan Xu, Weijia Shi, and Eunsol Choi. 2023. RECOMP: Improving Retrieval-Augmented LMs with Compression and Selective Augmentation. *CoRR* abs/2310.04408 (2023). <https://doi.org/10.48550/ARXIV.2310.04408> arXiv:2310.04408
- [39] Dayu Yang, Yue Zhang, and Hui Fang. 2022. An Exploration Study of Mixed-initiative Query Reformulation in Conversational Passage Retrieval. In *Proceedings of the Thirty-First Text REtrieval Conference, TREC 2022, online, November 15-19, 2022 (NIST Special Publication, Vol. 500-338)*, Ian Soboroff and Angela Ellis (Eds.). National Institute of Standards and Technology (NIST). [https://trec.nist.gov/pubs/trec31/papers/udel\\_fang.C.pdf](https://trec.nist.gov/pubs/trec31/papers/udel_fang.C.pdf)
- [40] Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 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. <https://openreview.net/forum?id=SkeHuCVFDr>
- [41] Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu, Tingchen Fu, Xinting Huang, Enbo Zhao, Yu Zhang, Yulong Chen, Longyue Wang, Anh Tuan Luu, Wei Bi, Freda Shi, and Shuming Shi. 2023. Siren's Song in the AI Ocean: A Survey on Hallucination in Large Language Models. *CoRR* abs/2309.01219 (2023). <https://doi.org/10.48550/ARXIV.2309.01219> arXiv:2309.01219
- [42] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. Judging LLM-as-a-Judge with MT-Bench and Chatbot Arena. In *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (Eds.). [http://papers.nips.cc/paper\\_files/paper/2023/hash/91f18a1287b398d378ef22505bf41832-Abstract-Datasets\\_and\\_Benchmarks.html](http://papers.nips.cc/paper_files/paper/2023/hash/91f18a1287b398d378ef22505bf41832-Abstract-Datasets_and_Benchmarks.html)