# A Surprisingly Simple yet Effective Multi-Query Rewriting Method for Conversational Passage Retrieval

Ivica Kostric
University of Stavanger
Stavanger, Norway
ivica.kostric@uis.no

Krisztian Balog
University of Stavanger
Stavanger, Norway
krisztian.balog@uis.no

## ABSTRACT

Conversational passage retrieval is challenging as it often requires the resolution of references to previous utterances and needs to deal with the complexities of natural language, such as coreference and ellipsis. To address these challenges, pre-trained sequence-to-sequence neural query rewriters are commonly used to generate a single de-contextualized query based on conversation history. Previous research shows that combining multiple query rewrites for the same user utterance has a positive effect on retrieval performance. We propose the use of a neural query rewriter to generate multiple queries and show how to integrate those queries in the passage retrieval pipeline efficiently. The main strength of our approach lies in its simplicity: it leverages how the beam search algorithm works and can produce multiple query rewrites at no additional cost. Our contributions further include devising ways to utilize multi-query rewrites in both sparse and dense first-pass retrieval. We demonstrate that applying our approach on top of a standard passage retrieval pipeline delivers state-of-the-art performance without sacrificing efficiency.

## CCS CONCEPTS

• **Information systems → Query reformulation**.

## KEYWORDS

Conversational search; Conversational passage retrieval; Neural query rewriting

## 1 INTRODUCTION

The main objective of a conversational search system is to effectively retrieve relevant answers to a wide range of information needs expressed in natural language [1]. A major difficulty lies in the conversational nature of the task, namely, that queries are often not standalone and need to be interpreted in the context of the user's previous queries as well as the system's answers to those [35]. Commonly, *query rewriting* (QR) addresses this by employing neural generative models to produce a single de-contextualized query at each conversation turn [9, 18, 33], and then feed that query to a retrieval pipeline. While this works well in many cases, the rewritten query may incorrectly capture the underlying intent, which leads to the retrieval of non-relevant answers. The challenge arises from the discrete generation process, which does not accurately capture the underlying probabilities or importance of terms.

In this paper, we seek to improve retrieval performance by generating multiple queries and modeling the importance of terms based on their presence across the queries. We leverage the beam search algorithm, commonly used in neural QR [10, 19]. Instead of keeping track of only the highest-likelihood sequence in a greedy fashion, the algorithm tracks and considers the best $k$ sequences at each generation step. We utilize the fact that the token probabilities are already computed in order to produce multiple rewrites at no additional cost. Thus, the only modification we need to make to the original beam search algorithm is to return all tracked sequences and their associated probabilities, as opposed to the single most probable sequence. The elegance of this method lies in its simplicity; it is computationally inexpensive yet remarkably effective.

The main research question driving our investigation is: *How can we effectively and efficiently utilize multiple query rewrites in conversational passage retrieval?* To answer this question, we take into consideration that retrieval can be performed using either sparse or dense retrieval methods. Sparse retrieval typically employs pseudo-relevance feedback techniques to expand the query and bridge the vocabulary gap. Our method effectively performs both term-importance estimation and query expansion to represent the underlying information need better and improves MRR by 1.06–6.31 percentage points compared to using a single query rewrite. Dense retrieval, based on contextual neural language models, works better with natural language queries (in contrast to bag-of-words models of sparse retrieval). However, it is computationally expensive and would scale linearly with the number of rewrites, rendering it impractical. Instead, we represent all query rewrites jointly by merging them into a single vector representation in the learned embedding space by weighted average pooling. Our method outperforms a single-query retrieval by 3.52–4.45 percentage points in absolut MRR score.

In summary, the main contribution of this paper is a conversational multi-query rewriting method, CMQR, that can be utilized in conversational passage retrieval and applied on top of any pipeline

that uses generative QR. The novelty of our approach is twofold: (1) it generates multiple query rewrites at no extra cost compared to current neural QR approaches, (2) it effectively utilizes the generated rewrites in both sparse and dense retrieval. Using the QReCC dataset for evaluation, we show that applying our method on top of any pipeline featuring generative QR improves performance, resulting in state-of-the-art results.

All resources developed for this paper (source code, query rewrites, and rankings) can be found at https://github.com/iai-group/sigir2024-multi-query-rewriting.

## 2 RELATED WORK

We focus on the task of *conversational passage retrieval*, where the goal is to retrieve relevant passages to the user query from a large passage collection. Unlike generative approaches, here, hallucinations are ensured not to occur as answers can only come from the collection. While there is some variety between retrieval pipeline architectures for conversational search, the vast majority include QR, followed by a first-pass candidate selection stage and then by one or more re-ranking steps [13, 15, 19, 29, 33]. This setup provides a good balance between efficiency and effectiveness [3]. We demonstrate the benefits of our approach to first-pass retrieval, using both spare and dense retrieval methods.

Automatic QR has a long tradition in IR, predominantly in query expansion, and has been shown effective in a range of tasks [4]. Conversational query rewriting (CQR) aims to generate clear, de-contextualized queries from raw inputs by considering conversation context and addressing coreferences, ellipsis [6], and topic transitions [29]. Crucially, it helps clarify and refine the user's needs in a dialogue setting [25]. Recently, neural rewriting methods leveraging large pre-trained language models, like GPT-2 [13, 28] and T5 [19], have become prevalent. While neural rewriting methods tend to outperform traditional query expansion techniques [7, 8], the best results are achieved by combining the two [13, 19, 28].

Two previous CQR studies are particularly relevant to our work. Lin et al. [19] propose two query reformulation methods: one focused on term importance and another on making human-like queries. They show that fusing ranked lists after separate retrieval stages for both queries increases recall. However, fusing the two lists after re-ranking showed no improvement. The main difference between this work and ours is that we generate multiple natural language query rewrites. Mo et al. [23] presents two neural models: one trained on rewriting queries and another to produce potential answers to the query, the idea being that pre-trained language models can directly answer questions by leveraging their internal knowledge. At inference, these potential answers are used to expand the query. Our approach differs in that we use a single model to estimate term importance and pick expansion terms.

In another line of recent research, deep neural networks are used to generate query embeddings directly from context [21, 22, 34]. These embeddings, used in conjunction with dense retrieval, can handle intricate conversational contexts more effectively. While they integrate seamlessly with advanced neural models for IR, they require systems capable of interpreting them, potentially demanding more computational resources or additional processing steps. In contrast, the traditional approach of QR offers better interpretability and flexibility, translating complex conversational contexts into standalone, understandable queries that can be processed directly and efficiently with existing retrieval pipelines.

## 3 METHOD

This section presents our method for generating multiple query rewrites in Section 3.1. The integration of those rewrites in sparse and dense retrievals is described in Sections 3.2 and 3.3, respectively.

### 3.1 Conversational Multi-Query Rewriting

*3.1.1 Problem Statement.* Conversational query rewriting (CQR) is the task of generating an informative context-independent query from a raw query (i.e., context-dependent user utterance) based on conversation context (i.e., history). Formally, we let $q_i$ be the raw query at conversation turn $i$, and $H = \langle q_1, r_1, q_2, r_2, ..., q_{i-1}, r_{i-1} \rangle$ be the conversation history up to that point, where $r_j$ is a response provided by the system to the $j$th query ($j \in [1..i-1]$). A context-independent query $\hat{q}_i$ is to created from the raw query $q_i$ by considering the conversation history up to that point: $\hat{q}_i = rewrite(\hat{q}_1, r_1, \hat{q}_2, r_2, \ldots, \hat{q}_{i-1}, r_{i_1}, q_i)$. The rewritten query $\hat{q}_i$ is considered self-contained and can be used downstream in various components of the retrieval pipeline.

*3.1.2 Motivation.* The majority of recent approaches employ generative neural models for CQR [15]. However, these models often fail to find omitted information or detect topic shifts in longer conversations [31]. In some cases, query rewrites introduce irrelevant terms, while in other cases, relevant terms are missing (akin to the notion of topic drift in pseudo relevance feedback [20, 27]). We hypothesize that it is often too challenging to accurately capture the user's information need in a single query rewrite. Therefore, instead of returning a single most likely rewrite, we propose to return the top $n$ query rewrites generated with the same model and then utilize these rewrites in all stages of the retrieval pipeline.

More specifically, generative neural approaches to CQR commonly use the beam search algorithm [11]. According to this technique, the probability scores of the $k$ most likely sequences are kept while generating a rewrite. When generation finishes, the sequence with the highest score is returned.

*3.1.3 Conversational Multi-Query Rewriting.* Motivated by the above, we propose *conversational multi-query rewriting* (CMQR), which uses a fine-tuned generative language model to generate the top $n$ query rewrites at each turn $i$, $\hat{q}_i^1, \hat{q}_i^2, \ldots, \hat{q}_i^n$, according to their beam search score. Each query rewrite $\hat{q}_i^j$ has an associated *rewrite score*:

$$RS(\hat{q}_i^j) = P(\hat{q}_i^j|H) = \left( \prod_{l=1}^{|\hat{q}_i^j|} P(t_l|t_{l-1}, \ldots, t_1, H) \right)^{\frac{1}{|\hat{q}_i^j|}},$$

where $t_1, \ldots, t_l$ are the predicted tokens, $|\hat{q}_i^j|$ is the sequence length of the $j$th query rewrite, and $H$ is the conversation history (cf. Section 3.1.1). Considering that RS is the product of the probabilities of all terms in a sequence, length normalization is applied to avoid the query rewriters' tendency to generate very short rewrites.

Due to the quadratic complexity with respect to the input size, a common practice is to limit the input to a maximum of 512 tokens [18]. To accommodate this restriction, we limit the context to the previously rewritten utterances, $\langle \hat{q}_1, \ldots, \hat{q}_{i-1} \rangle$, and the last system response, $r_{i-1}$. We do not rewrite the very first user utterance of a conversation under the assumption that it is already self-contained and states the necessary context (i.e., $\hat{q}_1 = q_1$).

Next, we discuss how to utilize multiple query rewrites in various components of a retrieval pipeline.

## 3.2 Sparse Retrieval

Sparse retrieval relies on a bag-of-words text representation, where each query term contributes to the document relevance estimate according to some scoring function, which is generally of the form $score(q, d) = \sum_{t \in q} w_{t,q} \times w_{t,d}$, where $w_{t,q}$ and $w_{t,d}$ are the term weights associated with query $q$ and document $d$, respectively. Our interest is in setting the term query weights, $w_{t,q}$. In the most commonly used retrieval scoring functions (e.g., BM25), this weight is taken to be the frequency of the term in the query, i.e., $w_{t,q} = c(t, q)$. In our approach, we construct a weighted bag-of-words query from all $n$ rewrites, where we set the weights for each term as the beam search score, i.e., $RS(\hat{q}_i^j)$. For each unique term in such obtained collection of terms, the term weights from all rewrites are summed up and normalized.

Effectively, the method performs both query expansion and a re-estimation of term importance based on multiple query rewrites. A similarity can be drawn to relevance feedback algorithms like RM3 [16], where two weighted queries are interpolated: the original query and the relevance language model query. The difference is, here, we interpolate different queries extracted from conversational context instead of retrieved documents and do not assign a pre-determined portion of the total weight mass to the original query terms. Importantly, our method is seen as complementary to relevance feedback and can be combined with it.

## 3.3 Dense Retrieval

Dense retrieval differs from sparse retrieval in that it aims to compute a relevance score based on the similarity between queries and documents represented in a continuous embedding space instead of matching on exact terms. In the simplest form, this score can be a dot product of the query and the document embedding vectors: $score(q, d) = h_q \cdot h_d$, where $h_q$ and $h_d$ are the learned query and document embedding vectors, respectively. We note that the learned embedding vectors can be pre-computed and stored for all documents in the collection, requiring only the computation of the query embedding vector at retrieval time.

Given $n$ rewrites with associated weights, we first generate embeddings for all rewrites separately and scale them according to the associated weights. Then, the scaled embeddings are summed up into a single vector ($h_{q_i}$) that can be used in a regular dense retrieval system. Formally, the query representation at turn $i$ is obtained by:

$$h_{q_i} = \sum_{j=1}^{n} encode_q(\hat{q}_i^j) RS(\hat{q}_i^j) .$$

Essentially, the final query is a weighted centroid of the query rewrites. This adds robustness to dense retrieval as the center of mass of multiple query rewrites will likely correspond better to the user's information need than a single rewrite would.

## 4 EXPERIMENTAL SETUP

We present the datasets we use in our experimental evaluation, introduce our baselines, and provide implementation details.

## 4.1 Dataset & Evaluation Metrics

Following previous work [23, 30, 32], we use the QReCC [2] dataset, which contains 14k conversations with 80k question-answer pairs, split into training and test sets (63.5k and 16.4k, respectively). The dialogues are based on questions from QuAC [5], TREC CAsT 2019 [6], and Google Natural Questions (NQ) [14], with TREC CAsT appearing only in the test set. Following Ye et al. [32], test instances lacking valid gold passage labels are excluded from our analysis. Consequently, our dataset comprises 8,209 test instances, distributed as 6,396 for QuAC, 1,442 for NQ, and 371 for TREC-CAsT. For a comprehensive evaluation, we present experimental results not only on the overall dataset but also on each subset.

We use mean reciprocal rank (MRR), mean average precision (MAP), and Recall@10 (R@10) as our evaluation metrics.

## 4.2 Baselines

We consider the following baselines for comparison: (1) **Manual rewrite**: Manually rewritten queries provided by the dataset. (2) **T5QR**$_{Manual}$ [18]: A strong T5-based [26] QR model. (3) **ConQRR** [30]: A T5-based model, optimized for retrieval performance using reinforcement learning. (4) **ConvGQR** [23]: An approach employing two T5-based models: one creates a de-contextualized query rewrite, the other predicts an answer to the query. The outputs are merged into a single query used for retrieval. (5) **LLM**$_{adhoc}$ [32]: An LLM query rewrite followed by an LLM query editor in an ad-hoc retrieval pipeline. The authors use ChatGPT 3.5 as their LLM. (6) **T5QR**$_{LLM}$ [32]: A sample of 10k datapoint is taken from the training set and run through the same approach as (4). The outputs are used to train a smaller, distilled model.

For each variant of the retrieval pipeline, we use the same T5-based QR approach, with a beam width of $k = 10$, but consider only the top rewrite. Wherever possible, i.e., code/model is made publicly available, we reproduce results on our system using V100 GPUs. Otherwise, we report the numbers from the original papers (indicated by $\dagger$).

## 4.3 Implementation Details

For QR, we fine-tune a T5 model [26] starting from a *t5-base*[1] checkpoint. We set the beam width to $k = 10$ for both single-query and multi-query approaches, as this was found to produce high-quality rewrites in [19].

For sparse retrieval, following [2], we employ the *Pyserini* [17] toolkit and use BM25 for retreival with hyparameters $k1 = 0.82$ and $b = 0.68$. We generate dense embeddings using a GTR [24]

---

[1]https://huggingface.co/t5-base

**Table 1: Performance of sparse and dense retrieval with QR methods. Bold and underlined indicate the best and second-best results, respectively. * denotes significant improvements with a t-test at p < 0.05 of CMQR over its single-query counterpart.**

| | Method | QReCC | | | QuAC | | | NQ | | | TREC-CAsT | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | MRR | MAP | R@10 | MRR | MAP | R@10 | MRR | MAP | R@10 | MRR | MAP | R@10 |
| Sparse (BM25) | Manual rewrite | 39.81 | 38.45 | 62.65 | 40.32 | 38.98 | 62.90 | 40.78 | 39.05 | 63.80 | **27.34** | **27.04** | **53.77** |
| | $T5QR_{Manual}$ | 31.03 | 29.86 | 50.17 | 30.75 | 29.60 | 49.77 | 34.06 | 32.51 | 52.79 | 24.15 | 23.91 | 46.90 |
| | $ConQRR^{\dagger}$ | 38.30 | – | 60.10 | 39.50 | – | 61.60 | 37.80 | – | 58.00 | 19.80 | – | 43.50 |
| | ConvGQR | 49.18 | 47.66 | 68.01 | 51.34 | 49.82 | 70.10 | 45.57 | _43.85_ | _64.06_ | 25.91 | 25.20 | 47.26 |
| | $LLM_{adhoc}^{\dagger}$ | 49.39 | 47.89 | 67.01 | _53.01_ | _51.52_ | 70.46 | 41.57 | 39.69 | 59.63 | 17.43 | 17.08 | 36.25 |
| | $T5QR_{LLM}$ | 46.72 | 45.19 | 64.00 | 50.13 | 48.64 | 67.50 | 39.26 | 37.14 | 55.99 | 16.98 | 16.97 | 34.64 |
| | $CMQR(T5QR_{Manual})$ | 37.34* | 35.99* | 58.42* | 37.92* | 36.57* | 59.31* | 38.05* | 36.46* | 57.31* | 24.43 | 24.07 | 47.44 |
| | CMQR(ConvGQR) | _50.24_* | _48.79_* | **69.87*** | 52.41* | 50.98* | _72.01_* | **46.83*** | **45.02*** | **65.34** | _26.14_ | _25.78_ | _50.49_ |
| | $CMQR(T5QR_{LLM})$ | **50.73*** | **49.2*** | _69.25_* | **53.96*** | **52.50*** | **72.27*** | 44.11* | 41.95* | 62.40* | 20.78* | 20.46* | 43.80* |
| Dense (GTR) | Manual rewrite | 43.15 | 41.27 | 66.12 | 40.67 | 38.92 | 64.59 | **54.01** | **51.25** | **73.13** | **43.74** | **42.98** | **65.23** |
| | $T5QR_{Manual}$ | 36.08 | 34.41 | 56.95 | 33.70 | 32.16 | 55.36 | 46.11 | 43.65 | 63.50 | 38.11 | 37.30 | 58.76 |
| | $ConQRR^{\dagger}$ | 41.80 | – | 65.10 | 41.60 | – | 65.90 | 45.30 | – | 64.10 | 32.70 | – | 55.20 |
| | ConvGQR | 42.18 | 40.43 | 63.39 | 41.21 | 39.55 | 63.04 | 49.20 | 46.80 | 67.63 | 31.51 | 30.88 | 52.96 |
| | $LLM_{adhoc}^{\dagger}$ | 44.99 | 43.19 | 67.34 | 45.21 | 43.48 | 68.30 | 47.64 | 45.20 | 67.27 | 30.91 | 30.48 | 51.03 |
| | $T5QR_{LLM}$ | 42.46 | 40.67 | 64.47 | 42.78 | 41.06 | 65.61 | 44.78 | 42.29 | 63.48 | 28.02 | 27.55 | 48.65 |
| | $CMQR(T5QR_{Manual})$ | 40.53* | 38.73* | 63.15* | 38.72* | 37.02* | 62.12* | 48.50* | 46.01* | 67.67* | _40.68_* | _39.91_* | _63.21_* |
| | CMQR(ConvGQR) | _45.82_* | _43.96_* | **69.75*** | 45.00* | _43.20_* | _70.00_* | _51.95_* | _49.50_* | _71.17_* | 36.11* | 35.50* | 59.97* |
| | $CMQR(T5QR_{LLM})$ | **45.98*** | **44.17*** | _69.31_* | **45.82*** | **44.08*** | **70.02*** | 49.58* | 47.14* | 69.18* | 34.69* | 34.09* | 57.64* |

model from a publicly available checkpoint.[2] The implementation is based on *Faiss* [12].

## 5 RESULTS

The evaluation results on the QReCC test set along with a breakdown of specific subsets are reported in Table 1. The table is split into two groups: Sparse (Top) and dense retrieval (Bottom), with the same query rewriting methods in each group. When multiple queries are considered for a given method, it is indicated by CMQR(*); in all our experiments, we consider 10 query rewrites.

Our main findings are as follows. First, the CMQR method consistently outperforms both its sparse and dense retrieval counterparts across all datasets. This trend highlights CMQR's effectiveness in improving retrieval metrics. We show significant improvements in the range from 1.06 to 6.31 in sparse retrieval and 3.52 to 4.45 in dense retrieval in terms of MRR (absolute percentage points). The biggest improvement is observed when adding CMQR to the weakest model ($T5QR_{Manual}$), suggesting that the term importance and query expansion methods have a major effect. This observation is further supported by the smallest gain observed with ConvGQR, which has integrated query expansion. Second, for both the sparse and dense retrieval groups, our CMQR method consistently outperforms all other methods on the overall QReCC dataset, achieving the best and second-best results, thereby setting a new state of the art. The strongest performance of CMQR with $T5QR_{LLM}$, a single *t5-base* model, is notable for its compact

size compared to previous best-performing methods, specifically the dual-*t5-base* model (ConvGQR) and the two-step LLM model ($LLM_{adhoc}$). Ye et al. [32] show $LLM_{adhoc}$ is 6 times slower than $T5QR_{LLM}$ making it impractical in real-world conversational applications. Finally, manual rewrites, representing human effort in query rewriting, show strong performance, especially on the test-only TREC-CAsT dataset. However, CMQR methods still surpass these human efforts on the overall QReCC dataset, underscoring the potential of automated systems to enhance retrieval tasks. Interestingly, adding CMQR to the model trained on manual rewrites ($T5QR_{Manual}$) almost reaches the performance of manual rewrites, while CMQR applied to $T5QR_{LLM}$ shows an absolute improvement of 10.91 MRR percentage points compared to the manual rewrite.

## 6 CONCLUSION

In this paper, we developed a method for generating multiple query rewrites for conversational search and explored how these can be incorporated into sparse and dense retrieval. This approach differs from the majority of previous work, where only a single query rewrite is used. We showed how multiple queries can be efficiently integrated at virtually no extra cost for both sparse and dense retrieval. Furthermore, we demonstrated that our method can be applied on top of existing query rewriting methods that employ generative query rewriting, yielding consistent improvements across all methods and resulting in state-of-the-art performance.

---

[2]https://huggingface.co/sentence-transformers/gtr-t5-base

In future work, we plan to employ multi-query rewrites also in the re-ranking components of multi-stage retrieval pipelines and determine automatically the number of rewrites to consider.

## ACKNOWLEDGMENTS

## REFERENCES

[1] Avishek Anand, Lawrence Cavedon, Matthias Hagen, Hideo Joho, Mark Sanderson, and Benno Stein. 2021. Dagstuhl seminar 19461 on conversational search: seminar goals and working group outcomes. *ACM SIGIR Forum* 54, 1 (2021), 3:1–3:11.
[2] Raviteja Anantha, Svitlana Vakulenko, Zhucheng Tu, Shayne Longpre, Stephen Pulman, and Srinivas Chappidi. 2021. Open-Domain Question Answering Goes Conversational via Question Rewriting. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL '21)*. 520–534.
[3] Nima Asadi and Jimmy Lin. 2013. Effectiveness/efficiency tradeoffs for candidate generation in multi-stage retrieval architectures. In *Proceedings of the 36th international ACM SIGIR conference on Research and development in information retrieval (SIGIR '13)*. 997–1000.
[4] Claudio Carpineto and Giovanni Romano. 2012. A Survey of Automatic Query Expansion in Information Retrieval. *ACM Comput. Surv.* 44, 1 (2012).
[5] Eunsol Choi, He He, Mohit Iyyer, Mark Yatskar, Wen-tau Yih, Yejin Choi, Percy Liang, and Luke Zettlemoyer. 2018. QuAC: Question Answering in Context. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing (EMNLP '18)*. 2174–2184.
[6] Jeffrey Dalton, Chenyan Xiong, and Jamie Callan. 2019. CAsT 2019: The Conversational Assistance Track Overview. In *In Proceedings of the Twenty-Eighth Text REtrieval Conference (TREC '19)*.
[7] Jeffrey Dalton, Chenyan Xiong, and Jamie Callan. 2020. CAsT 2020: The Conversational Assistance Track Overview. In *In Proceedings of the Twenty-Ninth Text REtrieval Conference (TREC '20)*.
[8] Jeffrey Dalton, Chenyan Xiong, and Jamie Callan. 2021. TREC CAsT 2021: The Conversational Assistance Track Overview. In *In Proceedings of the Thirtieth Text REtrieval Conference (TREC '21)*.
[9] Ahmed Elgohary, Denis Peskov, and Jordan Boyd-Graber. 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) (EMNLP-IJCNLP '19)*. 5918–5924.
[10] Jianfeng Gao, Chenyan Xiong, Paul Bennett, and Nick Craswell. 2022. Neural Approaches to Conversational Information Retrieval. arXiv:arXiv:2201.05176
[11] Alex Graves. 2012. Sequence Transduction with Recurrent Neural Networks. arXiv:1211.3711
[12] Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2021. Billion-Scale Similarity Search with GPUs. *IEEE Transactions on Big Data* 7, 3 (2021), 535–547.
[13] Vaibhav Kumar and Jamie Callan. 2020. Making Information Seeking Easier: An Improved Pipeline for Conversational Search. In *Findings of the Association for Computational Linguistics: EMNLP 2020 (EMNLP '20)*. 3971–3980.
[14] Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. Natural Questions: A Benchmark for Question Answering Research. *Transactions of the Association for Computational Linguistics* 7 (2019), 452–466.
[15] Weronika Lajewska and Krisztian Balog. 2023. From Baseline to Top Performer: A Reproducibility Study of Approaches at the TREC 2021 Conversational Assistance Track. In *Advances in Information Retrieval: 45th European Conference on Information Retrieval, ECIR 2023, Dublin, Ireland, April 2–6, 2023, Proceedings, Part III (ECIR '23)*. 177–191.
[16] Victor Lavrenko and W. Bruce Croft. 2001. Relevance based language models. In *Proceedings of the 24th annual international ACM SIGIR conference on Research and development in information retrieval (SIGIR '01)*. 120–127.
[17] Jimmy Lin, Xueguang Ma, Sheng-Chieh Lin, Jheng-Hong Yang, Ronak Pradeep, and Rodrigo Nogueira. 2021. Pyserini: A Python Toolkit for Reproducible Information Retrieval Research with Sparse and Dense Representations. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '21)*. 2356–2362.
[18] Sheng-Chieh Lin, Jheng-Hong Yang, Rodrigo Nogueira, Ming-Feng Tsai, Chuan-Ju Wang, and Jimmy Lin. 2020. Conversational Question Reformulation via Sequence-to-Sequence Architectures and Pretrained Language Models. arXiv:2004.01909
[19] Sheng-Chieh Lin, Jheng-Hong Yang, Rodrigo Nogueira, Ming-Feng Tsai, Chuan-Ju Wang, and Jimmy Lin. 2021. Multi-Stage Conversational Passage Retrieval: An Approach to Fusing Term Importance Estimation and Neural Query Rewriting. *ACM Transactions on Information Systems* 39, 4 (2021), 1–29.
[20] Craig Macdonald and Iadh Ounis. 2007. Expertise Drift and Query Expansion in Expert Search. In *Proceedings of the Sixteenth ACM Conference on Conference on Information and Knowledge Management (CIKM '07)*. 341–350.
[21] Kelong Mao, Zhicheng Dou, and Hongjin Qian. 2022. Curriculum Contrastive Context Denoising for Few-shot Conversational Dense Retrieval. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '22)*. 176–186.
[22] Kelong Mao, Zhicheng Dou, Hongjin Qian, Fengran Mo, Xiaohua Cheng, and Zhao Cao. 2022. ConvTrans: Transforming Web Search Sessions for Conversational Dense Retrieval. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing (EMNLP '22)*. 2935–2946.
[23] Fengran Mo, Kelong Mao, Yutao Zhu, Yihong Wu, Kaiyu Huang, and Jian-Yun Nie. 2023. ConvGQR: Generative Query Reformulation for Conversational Search. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers) (ACL '23)*. 4998–5012.
[24] Jianmo Ni, Chen Qu, Jing Lu, Zhuyun Dai, Gustavo Hernandez Abrego, Ji Ma, Vincent Zhao, Yi Luan, Keith Hall, Ming-Wei Chang, and Yinfei Yang. 2022. Large Dual Encoders Are Generalizable Retrievers. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing (EMNLP '22)*. 9844–9855.
[25] Gustavo Penha, Alexandru Balan, and Claudia Hauff. 2019. Introducing MANtIS: a novel Multi-Domain Information Seeking Dialogues Dataset. arXiv:1912.04639
[26] 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. *Journal of Machine Learning Research* 21, 140 (2020), 1–67.
[27] Anna Shtok, Oren Kurland, David Carmel, Fiana Raiber, and Gad Markovits. 2012. Predicting Query Performance by Query-Drift Estimation. *ACM Trans. Inf. Syst.* 30, 2 (2012).
[28] Svitlana Vakulenko, Nikos Voskarides, Zhucheng Tu, and Shayne Longpre. 2021. A Comparison of Question Rewriting Methods for Conversational Passage Retrieval. In *Proceedings of the 43rd European Conference on IR Research (ECIR '21)*. 418–424.
[29] Nikos Voskarides, Dan Li, Pengjie Ren, Evangelos Kanoulas, and Maarten de Rijke. 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 (SIGIR '20)*. 921–930.
[30] Zeqiu Wu, Yi Luan, Hannah Rashkin, David Reitter, Hannaneh Hajishirzi, Mari Ostendorf, and Gaurav Singh Tomar. 2022. CONQRR: Conversational Query Rewriting for Retrieval with Reinforcement Learning. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing (EMNLP '22)*. 10000–10014.
[31] Xinyi Yan, Charles L A Clarke, and Negar Arabzadeh. 2021. WaterlooClarke at the TREC 2021 Conversational Assistant Track. In *In Proceedings of the Thirtieth Text REtrieval Conference (TREC '21)*.
[32] Fanghua Ye, Meng Fang, Shenghui Li, and Emine Yilmaz. 2023. Enhancing Conversational Search: Large Language Model-Aided Informative Query Rewriting. In *Findings of the Association for Computational Linguistics: EMNLP 2023 (EMNLP '23)*. 5985–6006.
[33] Shi Yu, Jiahua Liu, Jingqin Yang, Chenyan Xiong, Paul Bennett, Jianfeng Gao, and Zhiyuan Liu. 2020. Few-Shot Generative Conversational Query Rewriting. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '20)*. 1933–1936.
[34] Shi Yu, Zhenghao Liu, Chenyan Xiong, Tao Feng, and Zhiyuan Liu. 2021. Few-Shot Conversational Dense Retrieval. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '21)*. 829–838.
[35] Hamed Zamani, Johanne R. Trippas, Jeff Dalton, and Filip Radlinski. 2023. Conversational Information Seeking. *Found. Trends Inf. Retr.* 17, 3-4 (2023), 244–456.