

# Ranked List Truncation for Large Language Model-based Re-Ranking

Chuan Meng

University of Amsterdam  
Amsterdam, The Netherlands  
c.meng@uva.nl

Negar Arabzadeh

University of Waterloo  
Waterloo, Canada  
narabzad@uwaterloo.ca

Arian Askari

Leiden University  
Leiden, The Netherlands  
a.askari@liacs.leidenuniv.nl

Mohammad Aliannejadi

University of Amsterdam  
Amsterdam, The Netherlands  
m.aliannejadi@uva.nl

Maarten de Rijke

University of Amsterdam  
Amsterdam, The Netherlands  
m.derijke@uva.nl

## ABSTRACT

Ranked list truncation (RLT) is a key task in information retrieval (IR). In this paper, we explore RLT in a novel “*retrieve-then-re-rank*” perspective, where we optimize re-ranking by truncating the retrieved list (i.e., trim re-ranking candidates). RLT is crucial for re-ranking as it can improve re-ranking efficiency by sending variable-length candidate lists to a re-ranker on a per-query basis, and has the potential to improve re-ranking effectiveness. Despite its importance, there is limited research into applying RLT methods to this new perspective. To address this research gap, we reproduce existing RLT methods in the context of re-ranking, especially newly emerged large language model (LLM)-based re-ranking. In particular, we examine to what extent established findings on RLT for retrieval are generalizable to the “*retrieve-then-re-rank*” setup from three perspectives: (i) assessing RLT methods in the context of LLM-based re-ranking with lexical first-stage retrieval, (ii) investigating the impact of different types of first-stage retrievers on RLT methods, and (iii) investigating the impact of different types of re-rankers on RLT methods. We perform experiments on the TREC 2019 and 2020 deep learning tracks, investigating 8 RLT methods for pipelines involving 3 retrievers and 2 re-rankers. We reach new insights into existing RLT methods in the context of re-ranking.

## CCS CONCEPTS

• **Information systems** → **Retrieval models and ranking**; **Language models**; **Retrieval efficiency**; **Retrieval effectiveness**.

## KEYWORDS

Ranked list truncation, Re-ranking, Large language models

### ACM Reference Format:

Chuan Meng, Negar Arabzadeh, Arian Askari, Mohammad Aliannejadi, and Maarten de Rijke. 2024. Ranked List Truncation for Large Language Model-based Re-Ranking. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '24)*, July 14–18, 2024, Washington D.C., USA. ACM, New York, NY, USA, 12 pages.

Permission to make digital or hard copies of part or all 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 third-party components of this work must be honored. For all other uses, contact the owner/author(s).

SIGIR '24, July 14–18, 2024, Washington D.C., USA  
© 2024 Copyright held by the owner/author(s).

## 1 INTRODUCTION

Ranked list truncation (RLT), a.k.a. query cut-off prediction [10, 25], has been studied for over two decades [3, 34] and recently attracted lots of attention in the information retrieval (IR) community [6, 7, 27, 31, 60, 63]. The task of RLT is to determine how many items in a ranked list should be returned such that a user-defined metric is optimized [7]. The user-defined metric typically considers the balance between the utility of search results and the cost of processing search results [60]. RLT is crucial in various IR applications where it is money- or time-consuming to review a returned item [63]. E.g., in patent search [28] or legal search [57, 60], providing a ranked list with an overwhelming number of items is too costly for patent experts or litigation support professionals [3].

**A new angle.** Existing studies mainly focus on RLT for single-stage retrieval, i.e., optimizing a user-defined metric (e.g., F1) of a retrieved list by truncating it at a certain position. In this paper, we focus on RLT for *re-ranking* [5], i.e., RLT in a “*retrieve-then-re-rank*” setup. In this setup, we still truncate a given retrieved list but focus on enhancing trade-offs between effectiveness and efficiency in re-ranking; truncating the retrieved list directly translates to a reduction in re-ranking depth. RLT for re-ranking is of great importance due to the following reasons: (i) RLT can improve re-ranking efficiency by sending variable-length lists of candidates to a re-ranker on a per-query basis; re-rankers are typically computationally expensive [27] and particularly, recently proposed large language model (LLM)-based re-rankers [29, 48–50, 55, 67, 69] with billions of parameters lead to a substantial increase in computational overhead [71], making it hard to apply them in practice; applying a fixed re-ranking cut-off to all queries is a common practice in literature; however, individual queries can be answered effectively by a shorter or a longer list of re-ranking candidates [9], so RLT can avoid unnecessary re-ranking costs by dynamically trimming the retrieved list; and (ii) RLT has the potential to improve re-ranking effectiveness; indeed, Zamani et al. [66] found that feeding a long retrieved list that includes many irrelevant items to a re-ranker instead can result in inferior re-ranking quality than a shorter retrieved list.

Despite its importance, limited research has explored the application of RLT methods in the “*retrieve-then-re-rank*” setup [13, 26, 62, 66]. E.g., Zamani et al. [66] only use one RLT method to truncate retrieved lists from BM25 to improve the performance of BERT-based re-ranking [44]. Put differently, there is a lack of systematic and comprehensive studies into the use of RLT methods that have

originally been introduced to optimize retrieval, in the context of re-ranking, especially newly emerged LLMs-based re-ranking.

**Research goal.** In this reproducibility paper, we examine *to what extent established findings on RLT for retrieval are generalizable to the “retrieve-then-re-rank” setup*. Specifically, we study the following **findings** from the literature on RLT: (i) Supervised RLT methods generally perform better than their unsupervised counterparts (e.g., set a fixed cut-off for all queries) [6, 27, 60, 63]. (ii) Distribution-based supervised RLT methods (i.e., directly predict a distribution among all candidate cut-off points) perform better than their sequential labeling-based counterpart (i.e., predict whether to truncate at each candidate point) [6, 60, 63]. (iii) Jointly learning RLT with other tasks (e.g., predicting the relevance of each item in the retrieved list) results in better RLT quality [60]. (iv) When truncating a retrieved list returned by a neural-based retriever, incorporating its embeddings improves RLT quality [31].

**Reproducibility challenge.** We highlight the main challenges of applying RLT methods from optimizing retrieval to optimizing re-ranking: (i) the new “retrieve-then-re-rank” setup leads to a new optimization goal for RLT methods, i.e., improving the trade-offs between effectiveness and efficiency in the re-ranking process; more importantly, a specific trade-off can be considered as the optimization goal to meet the requirements of a specific scenario, e.g., effectiveness is more important than efficiency in professional search than web-search; and (ii) the re-ranking setup introduces the *type* of re-ranker as a factor that influences RLT quality; also, it is important to investigate the impact of the interaction between retrievers and re-rankers on RLT; thus, it is important to explore RLT performance under different pipelines of widely-used retrievers, e.g., lexical, leaned sparse [17] and dense [29] retrievers, and different re-rankers, e.g., LLM-based [29] or pre-trained language model-based re-rankers [45].

**Scope.** We consider the challenges and examine each established finding from the literature on RLT in three settings: (i) we begin by checking if RLT methods optimizing for different trade-offs between effectiveness and efficiency of a state-of-the-art LLM-based re-ranker, RankLLaMA [29], with a lexical first-stage retriever; next, to study the impact of retriever types on RLT methods, we assess RLT methods for the LLM re-ranker with other types of retrievers, i.e., learned sparse (SPLADE++ [17]), and dense (RepLLaMA [29]) retrievers; and finally, to study the impact of the choice of re-rankers on RLT methods, we assess RLT methods for a widely-used pre-trained language model-based re-ranker, monoT5 [45].

We perform all experiments on the TREC 2019 and 2020 deep learning (TREC-DL) tracks [11, 12] and consider 8 RLT methods and pipelines involving 3 retrievers and 2 re-rankers, leading to various configurations. We aim to reach new insights into existing RLT methods in the “retrieve-then-re-rank” setup, elaborate on their errors in the new setup, and identify directions for future research.

**Lessons.** Our experiments reveal that while one finding on RLT for first-stage retrieval does not generalize to the “retrieve-then-re-rank” setup well, two findings do generalize to the setup under some specific cases. In short, unsupervised RLT methods outperform supervised ones, and methods that optimize for RLT jointly with other tasks can benefit from the joint training.

We share our supplementary materials to facilitate reproducibility at <https://github.com/ChuanMeng/RLT4Reranking>.

**Contributions.** Our main contributions are as follows:

- We reproduce a comprehensive set of RLT methods in a “retrieve-then-re-rank” perspective.
- We conduct an empirical analysis with a state-of-the-art LLM-based re-ranker, revealing that setting fixed re-ranking cut-offs results in unnecessary computational costs and diminishes re-ranking quality.
- We conduct extensive experiments on 2 datasets, 8 RLT methods and pipelines involving 3 retrievers and 2 re-rankers, leading to various configurations that allow a comprehensive understanding of how RLT well methods generalize to the new perspective.

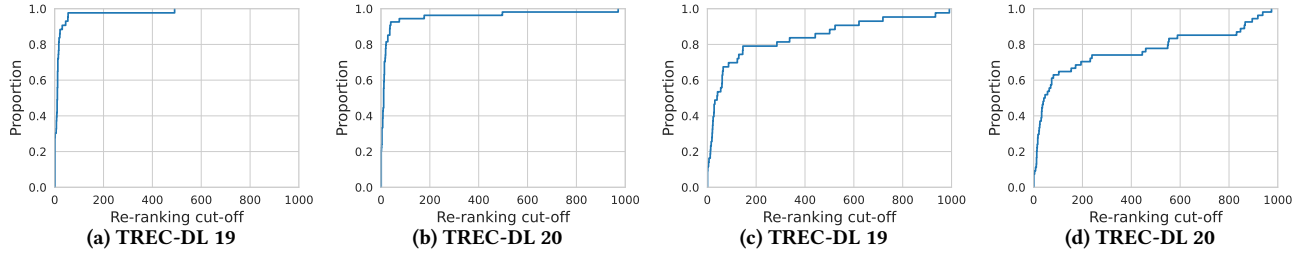
## 2 MOTIVATION

In the literature, applying a fixed re-ranking cut-off to all queries is a common practice [29, 30, 48–50, 55, 67, 69]; however, individual queries may need a shorter or a longer list of re-ranking candidates [9]. We conduct an empirical analysis to demonstrate how RLT holds the potential to enhance both the effectiveness and efficiency in re-ranking compared to fixed cut-offs. To do so, we analyze two “retrieve-then-re-rank” pipelines on the TREC-DL 19 and 20 datasets. We use an LLM-based re-ranker (RankLLaMA [29]) in both pipelines, but for first-stage retrieval, we employ a lexical retriever (BM25) in one pipeline and an LLM-based dense retriever (RepLLaMA [29]) in the other.

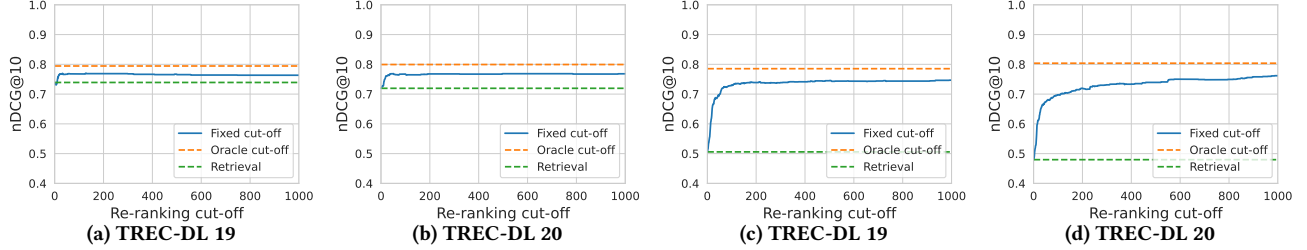
**Query-specific cut-offs improve efficiency.** We study an *oracle* setup in which we define the oracle as the minimum re-ranking cut-offs yielding the highest nDCG@10 values. We find that individual queries have different oracle cut-offs with a wide range. Thus, a fixed cut-off either wastes computational resources or compromises re-ranking quality for queries that need a deeper cut-off. Figure 1 illustrates the cumulative distribution of oracle cut-offs for both pipelines on both datasets. Interestingly, about 30% of queries do not need re-ranking with RepLLaMA as the retriever, and approximately 5% with BM25; thus, calling expensive re-rankers can be omitted for these queries.

**Query-specific cut-offs improve effectiveness.** Figure 2 illustrates the comparison of re-ranking quality between using oracle and fixed cut-offs. We find that oracle cut-offs always perform *statistically significantly* (paired t-test,  $p < 0.05$ ) better than all fixed cut-offs in terms of nDCG@10. Hence, a deeper re-ranking cut-off does not consistently result in improvement and can even be detrimental to re-ranking quality. Our finding is consistent with Zamani et al. [66]. While one might argue that the re-ranking results with a deeper cut-off might be underestimated because of the limited number of judged items within the top 10 ranks, i.e., judged@10 [49], we find that RankLLaMA’s judged@10 values for using a fixed cut-off at 1000 and oracle cut-offs are similar, e.g., 95.35% vs. 96.05% and 97.41% vs. 97.41% when RepLLaMA as the retriever on TREC-DL 19 and 20, respectively.

RLT methods truncate the retrieved list (i.e., trim re-ranking candidates) on a per-query basis, suggesting that effective RLT has the potential to improve re-ranking efficiency and effectiveness.



**Figure 1: Cumulative distribution function of oracle cut-offs for RepLLaMA–RankLLaMA (a, b) and BM25–RankLLaMA (c, d) on TREC-DL 19 and 20. The oracle cut-offs are the minimum re-ranking cut-offs that yield the highest nDCG@10 values.**



**Figure 2: nDCG@10 values for RepLLaMA–RankLLaMA (a, b) and BM25–RankLLaMA (c, d) w.r.t. re-ranking cut-offs on TREC-DL 19 and 20.**

### 3 PRELIMINARIES AND TASK DEFINITION

We recap the task definition of RLT and demonstrate the transition from optimizing retrieval to optimizing re-ranking.

**RLT for retrieval.** Given a user query  $q$ , a collection  $C$  containing  $|C|$  items, and a retrieved list  $L = [d_1, d_2, \dots, d_{|L|}]$  with  $|L|$  ( $|L| \ll |C|$ ) items induced by a first-stage retriever over  $C$  in response to  $q$ . An RLT approach  $f$  aims to predict a truncation point  $k$  that maximizes a target metric that is about the retrieved list  $L$  itself [6, 27, 31, 60, 63], formally:

$$k = f([x_1, x_2, \dots, x_{|L|}]) \in \{1, 2, \dots, |L|\}, \quad (1)$$

where  $[x_1, x_2, \dots, x_{|L|}]$  are item features corresponding one-to-one with the items in the retrieved list  $L = [d_1, d_2, \dots, d_{|L|}]$ . Typically,  $x$  includes the retrieval score [6] and item statistics [27, 60, 63]. As for the target metric,  $F1@k$  and  $DCG@k$  have been widely used in prior studies [6, 7, 27, 31, 60, 63].  $F1@k$  is calculated as:

$$F1@k = \frac{2 \cdot P@k \cdot R@k}{P@k + R@k}, \quad (2)$$

$$P@k = \frac{1}{k} \sum_{i=1}^k \mathbb{I}(y_i = 1), R@K = \frac{1}{N_L} \sum_{i=1}^k \mathbb{I}(y_i = 1),$$

where  $y_i \in \{0, 1\}$  is the relevance label for item  $d_i$  in the truncated retrieved list, and  $N_L$  denotes the number of relevant items in the retrieved list  $L$ . Note that the original discounted cumulative gain (DCG) metric [22] is a monotonic metric since its value always increases with the value of  $k$ ; it cannot evaluate RLT properly because the optimal solution would be not to truncate at all [6]. Therefore, the DCG metric employed in RLT penalizes non-relevant items, rendering it a non-monotonic metric [6, 7, 31, 60, 63]:

$$DCG@k = \sum_{i=1}^k \frac{y_i}{\log_2(i+1)}, \quad (3)$$

where  $y_i \in \{1, -1\}$ ;  $y_i = -1$  if item  $d_i$  is irrelevant to the query.

**RLT for re-ranking.** In the “retrieve-then-re-rank” setup, we no longer focus on optimizing the retrieved list  $L$ , but we aim to test the capability of the RLT in optimizing the trade-offs between effectiveness and efficiency in re-ranking. The truncated retrieved list  $L_{1:k} = [d_1, \dots, d_k]$ , serving as re-ranking candidates, is further forwarded to a re-ranker that returns a re-ranked list  $\hat{L}_{1:k}$ . Specifically, the target metric should assess the re-ranking quality of the re-ranked list  $\hat{L}_{1:k}$  in terms of an IR evaluation metric (e.g., nDCG@10), and also measure the computational cost.

### 4 REPRODUCIBILITY METHODOLOGY

This section elucidates our research questions, the corresponding experiments designed to address them, and our experimental setup.

#### 4.1 Research questions and experimental design

Our work is organized around the following research questions:

**RQ1** Do RLT methods generalize to the context of LLM-based re-ranking with a *lexical first-stage retriever* when optimized for different effectiveness/efficiency trade-offs?

To address **RQ1**, we first quantify the trade-off between re-ranking effectiveness and efficiency, and then optimize RLT methods to model different trade-offs between effectiveness and efficiency, simulating different requirements and scenarios; then, we evaluate their predicted truncation positions in terms of effectiveness and efficiency in LLM-based re-ranking with a lexical retriever.

**RQ2** Do RLT methods generalize to the context of LLM-based re-ranking with *learned sparse or dense first-stage retrievers* when optimized for the different trade-offs, and how does it compare to the case of a lexical retriever?

For answering **RQ2**, we still optimize RLT methods for different trade-offs of the LLM-based re-ranker used in **RQ1**, but study their

performance given learned sparse or dense retrievers, and compare the results with those of using a lexical retriever.

**RQ3** Do RLT methods generalize to the context of *pre-trained language model-based re-ranking*, and how does it compare to the case of an LLM-based re-ranker?

We address **RQ3** by evaluating the truncation points predicted by RLT methods w.r.t. effectiveness and efficiency in the context of a widely-used pre-trained language model-based re-ranker, and compare the results with those of the LLM-based re-ranker.

## 4.2 Experimental setup

**RLT approaches.** We reproduce a variety of unsupervised and supervised RLT methods [6, 7, 27, 31, 60, 63]. Note that although candidate pruning methods [13, 26, 62], specifically designed for pruning the candidate list in a cascading ranking scheme, can be used in our scenario, previous studies study RLT and candidate pruning separately and do not merge these two lines of research [6, 7, 27, 31, 60, 63]. In this work, our focus is on reproducing RLT methods in “retrieve-then-re-rank” setup.

As for unsupervised RLT methods, we experiment with recent RLT studies [6, 7, 31, 60, 63]. We consider the following unsupervised RLT methods.

- **Fixed- $k$**  [27] applies a fixed re-ranking cut-off to all queries; we follow common practice and consider cut-offs that are widely used in the literature about re-ranking, namely 10 [30], 20 [30, 48, 49], 100 [16, 29, 48–50, 55, 67, 69–71], 200 [29], 1000 [52].
- **Greedy- $k$**  [27] uses the fixed truncation position  $k$  that maximizes the target metric value on a training set.
- **Surprise** [7] first calibrates retrieval scores by using generalized Pareto distributions in extreme value theory [47], and truncates a ranked list using a score threshold.

We consider the following supervised RLT methods:

- **BiCut** [27] is a *sequential labeling-based* method; it uses a bidirectional long short-term memory (LSTM) to encode item features over a ranked list, and optimizes the LSTM make a binary prediction (continue or truncate) at each position in a ranked list.
- **Choppy** [6] is a *distribution-based* method, which directly predicts the distribution among all candidate cut-off points, using a transformer encoder [59] to encode item features over a ranked list and predicts the distribution.
- **AttnCut** [63] is also *distribution-based*, encoding item features using a bidirectional LSTM and a transformer encoder.
- **MtCut** [60] is also *distribution-based* and similar to AttnCut, but jointly trains the RLT task with two auxiliary tasks: predicting the relevance of each item in the ranked list and increasing the margin between relevant and irrelevant items. We use the MMoECut variant due to its superior performance.
- **LeCut** [31] is another *distribution-based* and similar to AttnCut, but can only work with a neural-based retriever and incorporates its query-item embeddings as one of the item features. Ma et al. [31] further optimize LeCut with an learning-to-rank (LtR) model jointly. We omit this phase for a fair comparison since other methods are trained without signals from an external LtR model.

We also include **Oracle**, which truncates the retrieved list at the earliest position maximizing re-ranking quality per query.

**Optimizing effectiveness/efficiency trade-offs.** The leading challenge of adapting RLT methods in the context of re-ranking is to optimize RLT methods with a specific trade-off between effectiveness and efficiency. To solve this challenge, we need to score each truncation point (i.e., re-ranking candidate cut-off) under different effectiveness/efficiency trade-offs. To do so, we quantify different trade-offs using the efficiency-effectiveness trade-off (EET) metric [61] and then compute EET scores at a specific trade-off for all re-ranking candidate cut-offs.

EET is defined for as the weighted harmonic mean of effectiveness  $\sigma$  and efficiency  $\gamma$  measures:

$$EET = \frac{(1 + \beta^2) \cdot (\gamma \cdot \sigma)}{\beta^2 \cdot \sigma + \gamma}, \quad (4)$$

where  $\beta$  is a hyperparameter to control the relative importance of effectiveness and efficiency, where  $\beta = 0$  only considers effectiveness and as it increases more attention is paid to efficiency. EET requires instantiation of  $\sigma$  and  $\gamma$  based on the specific use case [61]. We follow [61] to instantiate efficiency  $\gamma$  using “exponential decay”:

$$\gamma = \exp(\alpha \cdot k), \quad (5)$$

where  $k \in \{1, 2, \dots, |L|\}$  is a truncation point (i.e., re-ranking candidate cut-off) in the given retrieved list  $L$ , and  $\alpha < 0$  is a hyperparameter to control how rapidly the efficiency measure decreases; we set  $\alpha$  to -0.001. We instantiate effectiveness  $\sigma$  as the re-ranking gain with a cut-off  $k$ , which is quantified by the difference of re-ranking performance with a cut-off  $k$  minus the performance without re-ranking; the performance is in terms of an IR evaluation metric (e.g., nDCG@10).

Therefore, we can adjust  $\beta$  in Equation 4 and  $\alpha$  in Equation 5 in EET to quantify different effectiveness/efficiency trade-offs in re-ranking, so as to generate EET value distributions across all cut-off candidates under different trade-offs; see Figure 3. All distribution-based RLT methods (Choppy [6], AttnCut [63], MtCut [60] and LeCut [31]) can directly optimize a EET value distribution across all re-ranking cut-off candidates at a certain trade-off.

However, the sequential labeling-based RLT method BiCut [27] cannot optimize a EET value distribution. During training, BiCut optimizes the probability of “continue” and “truncation” at each position in a ranked list via the following loss:

$$\mathcal{L} = \sum_{i=1}^{|L|} (\eta \mathbb{I}(y_i = 0) \frac{p_i}{1-r} + (1-\eta) \frac{1-p_i}{r} \mathbb{I}(y_i = 1)), \quad (6)$$

where  $y_i \in \{0, 1\}$  is the relevance label for an item at a position  $i$ ,  $h$   $p_i$  is the “continue” probability at a position  $i$ , and  $r$  is the proportion of relevant items in the ranked list;  $\eta \in [0, 1]$  is a hyperparameter to control the balance between “continue” and “truncation”. We optimize BiCut for different effectiveness/efficiency trade-offs by adjusting  $\eta$  values, e.g., BiCut trained with a high  $\eta$  value tends to truncate a retrieved list earlier, resulting in more efficiency.

**Datasets.** We experiment with 2 widely-used IR datasets, TREC 2019 and 2020 deep learning (TREC-DL) tracks [11, 12].<sup>1</sup> These datasets offer relevance judgments in multi-graded relevance scales

<sup>1</sup> We also conducted experiments on Robust04 and draw a similar conclusion as TREC-DL 19 and 20; due to the limited space, see our repository for Robust04 results.



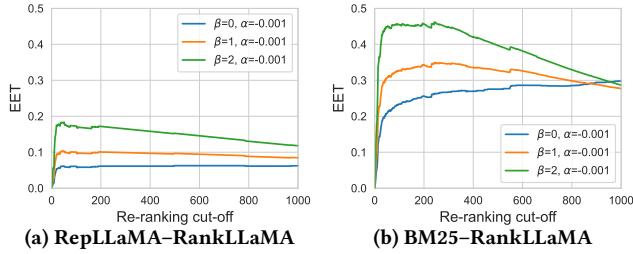


Figure 3: Average EET values across TREC-DL 20 queries w.r.t. re-ranking cut-offs. We use nDCG@10 in effectiveness  $\sigma$ .

per query. TREC-DL 19 and 20 are built upon MS MARCO V1 passage ranking collection encompassing 8.8 million passages.

**Choice of retrievers.** Regarding retrievers, we employ three distinct types: a lexical-based retriever BM25 [51], a learned sparse retriever SPLADE++ (“EnsembleDistil”) [17] and an LLM-based dense retriever RepLLaMA (7B) [29]. To increase the comparability and reproducibility of our paper, we obtain retrieval results of BM25 and SPLADE++ using the publicly available resource from Pyserini<sup>2</sup> and get retrieval results of RepLLaMA from Tevatron,<sup>3</sup> each retriever returns 1000 items per query.

**Choice of re-rankers.** For **RQ1**, **RQ2**, we employ a state-of-the-art LLM-based point-wise reranker, RankLLaMA (7B) [29] and use the resource from Tevatron. For **RQ3**, we employ a widely-used pre-trained language model-based re-ranker, monoT5 (“monot5-base-msmarco”) [45] and use the resource from PyGaggle.<sup>4</sup>

**Evaluation metrics.** We measure re-ranking effectiveness using nDCG@10, the official evaluation metric in TREC deep learning tracks [11, 12], and a widely employed metric in ranking literature [29, 45, 49]. We follow [71] to evaluate re-ranking efficiency by calculating the average re-ranking cut-off across all test set queries, i.e., the number of the average number of re-ranking inferences per query. This consideration is driven by the fact that the re-rankers we employ in this paper are point-wise, and the time spent in point-wise re-ranking is directly proportional to the length of a re-ranking cut-off (i.e., the length of a truncated retrieved list) [33]. We additionally gauge the efficiency by measuring per-query latency; all latency measurements exclude the time to load data and models. We do not consider the latency of first-stage retrieval. Note that RLT methods are lightweight with significantly fewer parameters compared to state-of-the-art re-rankers; the latency introduced by RLT methods can be neglected in the “retrieve-then-re-rank” setup.

**Implementation details.** To implement RLT methods in the context of re-ranking, we refer to the code<sup>5</sup> for RLT in single-stage retrieval, released by Ma et al. [31].

We pass the top-1000 retrieved items to all RLT methods per query because 1000 is typically the deepest re-ranking depth in the literature [29, 46, 52]. For unsupervised RLT methods, Surprise [7] only depends on retrieval scores and uses a score threshold to truncate a ranked list; note that Surprise cannot be tuned for different effectiveness/efficiency trade-offs because the score threshold is set based on Cramer-von-Mises statistic testings [14] and the threshold is not a tunable hyperparameter.

For all supervised RLT methods, we use identical item features to eliminate confounding factors from the input; each item is represented by its retrieval score, length, unique token count, and the cosine similarity between its tf-idf/doc2vec [24] vector and the vectors of its adjacent items. We follow [60, 63] to use gensim<sup>6</sup> for computing tf-idf and doc2vec vectors for each item; The dimension of the tf-idf vectors on the MS MARCO V1 collection is 846,221; we follow [60] to set the dimension of doc2vec vectors as 128. LeCut relies on query-item embeddings from a neural retriever as extra item features; thus, we provide LeCut with the concatenation of query and item embeddings from RepLLaMA. Note that we follow all original work to set hyperparameters: for BiCut, we set # Bi-LSTM layers to 2 and the LSTM hidden size to 128; we train BiCut via Adam [23] with a learning rate of  $1 \times 10^{-4}$ ; for Choppy, we set # transformer layers to 3, # transformer heads to 8, and transformer hidden size to 128; we train Choppy via Adam with a learning rate of  $1 \times 10^{-3}$ ; for AttnCut, we set # Bi-LSTM layers to 2, the LSTM hidden size to 128, # transformer heads to 4, and the transformer hidden size to 128; MtCut and LeCut share almost the same hyperparameters with AttnCut; we train AttnCut, MtCut and LeCut via Adam with a learning rate of  $3 \times 10^{-5}$ . We train all supervised RLT methods for 100 epochs using a batch size of 64 on TREC-DL 19, then infer them on TREC-DL 20, and vice versa. All methods are trained/inferred on an NVIDIA A100 GPU (40GB).

## 5 RESULTS AND DISCUSSIONS

### 5.1 RLT for LLM-based re-ranking

To answer **RQ1**, we focus on optimizing RLT methods for different effectiveness/efficiency trade-offs in the context of an LLM-based re-ranker (RankLLaMA [29]) with a lexical retriever (BM25). We optimize Greedy- $k$ , Choppy, AttnCut and MtCut for three trade-offs between effectiveness and efficiency by adjusting  $\beta$  in Equation 4; we consider  $\beta=0$  (emphasizing effectiveness), 1 (weight effectiveness and efficiency equally), and 2 (prioritizing efficiency). We optimize BiCut by tuning  $\eta$  in Equation 6; we consider  $\eta=0.4$  (emphasizing effectiveness), 0.5 (weight effectiveness and efficiency equally), and 0.6 (prioritizing efficiency).<sup>7</sup> We present the results on TREC-DL 19 and 20 in Figure 4 and 5; we also report results in Table 1. We have four observations.

First, we show the benefits of using a dynamic cutoff approach with RLT methods in uncertain re-ranking depth scenarios to balance effectiveness and efficiency. For example, in a “retrieve-then-rerank” setup using BM25 and RankLLaMA on the TREC-DL 19 dataset shown in Figure 4, we found that a dynamic cutoff method like BiCut ( $\beta = 1$ ,  $\eta = 0.5$ ), with an average predicted cutoff around 200, boosts efficiency fivefold over a fixed cutoff of 1000, with minimal effectiveness loss under 0.2%, demonstrating comparable results with significantly better efficiency.

Second, supervised RLT methods (BiCut, Choppy, AttnCut and MtCut) do not have a clear advantage over their unsupervised counterparts (Fixed- $k$ , Greedy- $k$ , surprise) in terms of nDCG@10 values and re-ranking cost across all the three trade-offs.

<sup>2</sup> <https://github.com/castorini/pyserini>

<sup>3</sup> <https://github.com/texttron/tevatron>

<sup>4</sup> <https://github.com/castorini/pygaggle>

<sup>5</sup> <https://github.com/myx666/LeCut>

<sup>6</sup> <https://radimrehurek.com/gensim> <sup>7</sup> We also explore  $\eta$  values of 0.3 and 0.7. BiCut trained with the former tends not to truncate the retrieved list at all, while BiCut trained with the latter tends to truncate the entire retrieved list.

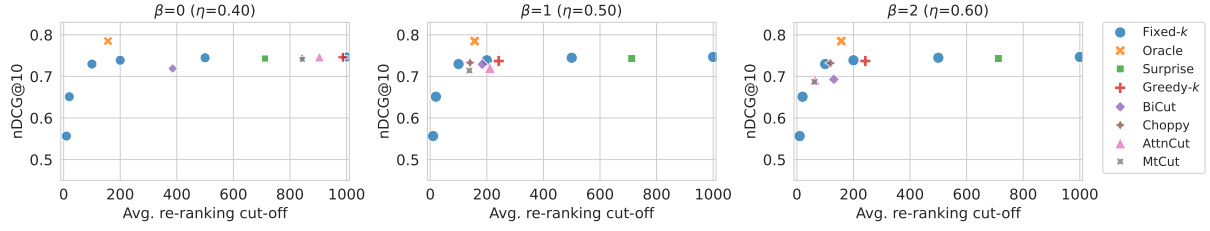


Figure 4: BM25-RankLLaMA on TREC-DL 19.

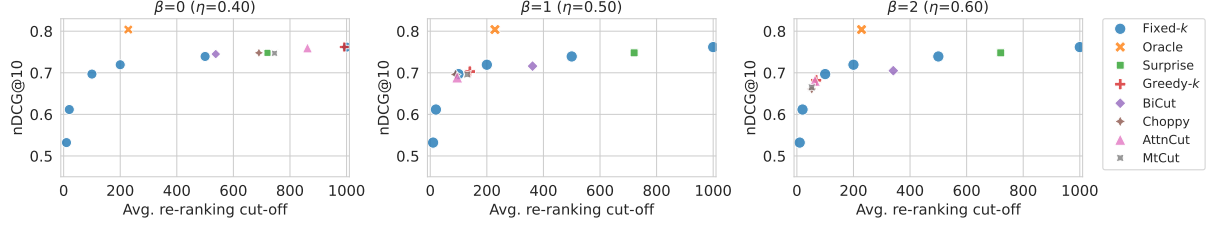


Figure 5: BM25-RankLLaMA on TREC-DL 20.

**Table 1: A comparison of RLT methods, optimized for re-ranking effectiveness/efficiency tradeoffs, in predicting re-ranking cut-off points for the BM25-RankLLaMA pipeline on TREC-DL 19 and 20. Query latency is measured in seconds.**

| Method                 | TREC-DL 19 |              |       | TREC-DL 20 |              |       |
|------------------------|------------|--------------|-------|------------|--------------|-------|
|                        | Avg. k     | nDCG@10      | Lat.  | Avg. k     | nDCG@10      | Lat.  |
| w/o re-ranking         | -          | 0.506        | -     | -          | 0.480        | -     |
| Fixed-k (10)           | 10         | 0.557        | 0.30  | 10         | 0.532        | 0.30  |
| Fixed-k (20)           | 20         | 0.651        | 0.60  | 20         | 0.612        | 0.60  |
| Fixed-k (100)          | 100        | 0.730        | 2.98  | 100        | 0.697        | 2.98  |
| Fixed-k (200)          | 200        | 0.739        | 5.95  | 200        | 0.719        | 5.96  |
| Fixed-k (1000)         | 1000       | <b>0.747</b> | 29.77 | 1000       | <b>0.762</b> | 29.78 |
| Surprise               | 712.12     | 0.743        | 21.20 | 720.56     | 0.748        | 21.46 |
| Greedy-k ( $\beta=0$ ) | 987        | 0.746        | 29.39 | 992        | <b>0.762</b> | 29.54 |
| BiCut ( $\eta=0.40$ )  | 385.51     | 0.719        | 11.48 | 537.57     | 0.745        | 16.01 |
| Choppy ( $\beta=0$ )   | 842.95     | <u>0.744</u> | 25.10 | 690.24     | 0.748        | 20.56 |
| AttnCut ( $\beta=0$ )  | 904.16     | <b>0.747</b> | 26.92 | 861.87     | <u>0.760</u> | 25.67 |
| MtCut ( $\beta=0$ )    | 843.67     | 0.741        | 25.12 | 745.44     | 0.747        | 22.20 |
| Greedy-k ( $\beta=1$ ) | 242        | 0.737        | 7.21  | 140        | 0.703        | 4.17  |
| BiCut ( $\eta=0.50$ )  | 184.02     | 0.729        | 5.48  | 361.59     | 0.716        | 10.77 |
| Choppy ( $\beta=1$ )   | 140.63     | 0.733        | 4.19  | 89.65      | 0.696        | 2.67  |
| AttnCut ( $\beta=1$ )  | 211.21     | 0.720        | 6.29  | 95.11      | 0.689        | 2.83  |
| MtCut ( $\beta=1$ )    | 138.26     | 0.714        | 4.12  | 130.98     | 0.696        | 3.90  |
| Greedy-k ( $\beta=2$ ) | 242        | 0.737        | 7.21  | 68         | 0.682        | 2.03  |
| BiCut ( $\eta=0.60$ )  | 130.74     | 0.693        | 3.89  | 340.81     | 0.705        | 10.15 |
| Choppy ( $\beta=2$ )   | 118.84     | 0.732        | 3.54  | 52.67      | 0.661        | 1.57  |
| AttnCut ( $\beta=2$ )  | 64.16      | 0.692        | 1.91  | 63.89      | 0.681        | 1.90  |
| MtCut ( $\beta=2$ )    | 62.02      | 0.687        | 1.85  | 52.48      | 0.665        | 1.56  |
| Oracle                 | 156.88     | 0.785        | 4.67  | 228.56     | 0.804        | 6.80  |

Third, distribution-based supervised RLT methods (Choppy, AttnCut and MtCut) save more re-ranking cost while maintaining competitive nDCG@10 values compared to the sequential labeling-based method (BiCut) across the trade-offs of weighting effectiveness and efficiency equally or prioritizing efficiency. E.g., as illustrated in Table 1, at  $\beta = 1/\beta = 2$ , Choppy attains a nDCG@10 value of 0.733 (0.732) with an average cut-off of 140.63 (118.84), while

**Table 2: Comparison of RLT methods, optimized for re-ranking effectiveness/efficiency tradeoffs, in predicting re-ranking cut-off points for the SPLADE++-RankLLaMA pipeline on TREC-DL 19 and 20. Query latency measured in seconds.**

| Method                 | TREC-DL 19 |              |       | TREC-DL 20 |              |       |
|------------------------|------------|--------------|-------|------------|--------------|-------|
|                        | Avg. k     | nDCG@10      | Lat.  | Avg. k     | nDCG@10      | Lat.  |
| w/o re-ranking         | -          | 0.731        | -     | -          | 0.720        | -     |
| Fixed-k (10)           | 10         | 0.740        | 0.30  | 10         | 0.741        | 0.30  |
| Fixed-k (20)           | 20         | <u>0.773</u> | 0.60  | 20         | <u>0.778</u> | 0.60  |
| Fixed-k (100)          | 100        | 0.769        | 2.98  | 100        | 0.771        | 2.98  |
| Fixed-k (200)          | 200        | 0.769        | 5.95  | 200        | 0.769        | 5.96  |
| Fixed-k (1000)         | 1000       | 0.768        | 29.77 | 1000       | 0.768        | 29.79 |
| Surprise               | 701.93     | 0.766        | 20.90 | 684.13     | 0.768        | 20.38 |
| Greedy-k ( $\beta=0$ ) | 65         | 0.770        | 1.94  | 42         | 0.772        | 1.25  |
| BiCut ( $\eta=0.40$ )  | 1000       | 0.768        | 29.77 | 1000.0     | 0.768        | 29.79 |
| Choppy ( $\beta=0$ )   | 904.12     | 0.768        | 26.92 | 999.96     | 0.768        | 29.79 |
| AttnCut ( $\beta=0$ )  | 999        | 0.768        | 29.74 | 999        | 0.768        | 29.76 |
| MtCut ( $\beta=0$ )    | 999        | 0.768        | 29.74 | 1000.0     | 0.768        | 29.79 |
| Greedy-k ( $\beta=1$ ) | 65         | 0.770        | 1.94  | 21         | 0.775        | 0.63  |
| BiCut ( $\eta=0.50$ )  | 2          | 0.731        | 0.06  | 1          | 0.720        | 0.00  |
| Choppy ( $\beta=1$ )   | 68.09      | 0.771        | 2.03  | 101.91     | 0.766        | 3.03  |
| AttnCut ( $\beta=1$ )  | 46.37      | 0.771        | 1.38  | 53.20      | 0.771        | 1.58  |
| MtCut ( $\beta=1$ )    | 119.67     | <b>0.775</b> | 3.56  | 996        | 0.768        | 29.67 |
| Greedy-k ( $\beta=2$ ) | 18         | 0.769        | 0.54  | 21         | 0.775        | 0.63  |
| BiCut ( $\eta=0.60$ )  | 2          | 0.731        | 0.06  | 1          | 0.720        | 0.00  |
| Choppy ( $\beta=2$ )   | 1          | 0.731        | 0.00  | 20.57      | 0.764        | 0.61  |
| AttnCut ( $\beta=2$ )  | 23.65      | 0.758        | 0.70  | 52.09      | 0.772        | 1.55  |
| MtCut ( $\beta=2$ )    | 19.74      | 0.756        | 0.59  | 25.93      | <b>0.784</b> | 0.77  |
| Oracle                 | 17.44      | 0.810        | 0.52  | 27.69      | 0.820        | 0.82  |

BiCut achieves the lower nDCG@10 value 0.729 (0.693) with the higher average cut-off 184.02 (130.74).

Fourth, the supervised method (MtCut) optimizing RLT in a multi-task manner does not show a clear advantage over other supervised methods across all three trade-offs.

**Table 3: Comparison of RLT methods optimized for re-ranking effectiveness/efficiency tradeoffs, in predicting cut-off points for the RepLLaMA–RankLLaMA pipeline on TREC-DL 19 and 20. Query latency measured in seconds.**

| Method                    | TREC-DL 19 |              |       | TREC-DL 20 |              |       |
|---------------------------|------------|--------------|-------|------------|--------------|-------|
|                           | Avg. k     | nDCG@10      | Lat.  | Avg. k     | nDCG@10      | Lat.  |
| w/o re-ranking            | -          | 0.738        | -     | -          | 0.720        | -     |
| Fixed- $k$ (10)           | 10         | 0.742        | 0.30  | 10         | 0.729        | 0.30  |
| Fixed- $k$ (20)           | 20         | 0.765        | 0.60  | 20         | 0.761        | 0.60  |
| Fixed- $k$ (100)          | 100        | <b>0.769</b> | 2.98  | 100        | 0.767        | 2.99  |
| Fixed- $k$ (200)          | 200        | <u>0.768</u> | 5.96  | 200        | 0.768        | 5.97  |
| Fixed- $k$ (1000)         | 1000       | 0.763        | 29.81 | 1000       | 0.768        | 29.86 |
| Surprise                  | 458.30     | 0.765        | 13.66 | 459.98     | 0.767        | 13.73 |
| Greedy- $k$ ( $\beta=0$ ) | 770        | 0.763        | 22.95 | 31         | 0.764        | 0.93  |
| BiCut ( $\eta=0.40$ )     | 1000       | 0.763        | 29.81 | 1000       | 0.768        | 29.86 |
| Choppy ( $\beta=0$ )      | 76.72      | 0.766        | 2.29  | 187.48     | 0.764        | 5.60  |
| AttnCut ( $\beta=0$ )     | 928.95     | 0.763        | 27.69 | 572.74     | <b>0.770</b> | 17.10 |
| MtCut ( $\beta=0$ )       | 509.70     | 0.758        | 15.19 | 130.04     | 0.727        | 3.88  |
| LeCut ( $\beta=0$ )       | 417.77     | 0.766        | 12.45 | 441.22     | <b>0.770</b> | 13.17 |
| Greedy- $k$ ( $\beta=1$ ) | 50         | 0.767        | 1.49  | 55         | 0.766        | 1.64  |
| BiCut ( $\eta=0.50$ )     | 166.81     | 0.766        | 4.97  | 322.50     | 0.768        | 9.63  |
| Choppy ( $\beta=1$ )      | 106.53     | 0.766        | 3.18  | 60.52      | 0.766        | 1.81  |
| AttnCut ( $\beta=1$ )     | 458.47     | 0.765        | 13.67 | 341.28     | <u>0.769</u> | 10.19 |
| MtCut ( $\beta=1$ )       | 582.58     | 0.761        | 17.37 | 291.85     | 0.762        | 8.71  |
| LeCut ( $\beta=1$ )       | 315.37     | <b>0.769</b> | 9.40  | 315.56     | <u>0.769</u> | 9.42  |
| Greedy- $k$ ( $\beta=2$ ) | 50         | 0.767        | 1.49  | 55         | 0.766        | 1.64  |
| BiCut ( $\eta=0.60$ )     | 153.70     | 0.766        | 4.58  | 122.41     | 0.764        | 3.65  |
| Choppy ( $\beta=2$ )      | 72.07      | 0.766        | 2.15  | 41.43      | 0.763        | 1.24  |
| AttnCut ( $\beta=2$ )     | 210        | 0.767        | 6.26  | 259.39     | <u>0.769</u> | 7.74  |
| MtCut ( $\beta=2$ )       | 514.53     | 0.763        | 15.34 | 213.74     | 0.764        | 6.38  |
| LeCut ( $\beta=2$ )       | 214.49     | <b>0.769</b> | 6.39  | 261.20     | 0.768        | 7.80  |
| Oracle                    | 23.33      | 0.794        | 0.70  | 42.46      | 0.799        | 1.27  |

## 5.2 The impact of retriever types on RLT

To answer **RQ2**, we optimize RLT methods for different trade-offs in the context of an LLM-based re-ranker (RankLLaMA [29]) with other types of novel retrievers; here we explore learned sparse (SPLADE++ [17]) and dense (RepLLaMA [29]) retrievers. Similarly to the Section 5.1, we optimize RLT methods for three kinds of trade-offs. We present the results for SPLADE++–RankLLaMA on TREC-DL 19 and 20 in Figure 6 and 7, as well as Table 2; we show the results for RepLLaMA–RankLLaMA on TREC-DL 19 and 20 in Figure ?? and 8, as well as Table 3. According to these results, we have four observations.

First, unsupervised methods exhibit a marked advantage over supervised methods for both pipelines. E.g., Fixed- $k$  (10) achieves the best nDCG@10 values while maintaining lowest re-ranking cost for SPLADE++–RankLLaMA on TREC-DL 19 and TREC-DL 20 (except for  $\beta=2$ ); Fixed- $k$  (10) and Greedy- $k$  show the same case for RepLLaMA–RankLLaMA on TREC-DL 19 and 20. We speculate that this is because SPLADE++ and RepLLaMA have much stronger retrieval performance than BM25; many relevant items are at top-ranks in the retrieved lists returned by SPLADE++ and RepLLaMA; therefore, a fixed shallow re-ranking cut-off could work well.

Second, similar to our observation in Section 5.1, the most distribution-based supervised RLT methods outperform the sequential

labeling-based method (BiCut) in terms of effectiveness and efficiency.

Third, the supervised method (MtCut) optimizing RLT in a multi-task manner shows the best nDCG@10 value (except for Oracle), while maintaining the lowest re-ranking cost for SPLADE++–RankLLaMA on TREC-DL 20 ( $\beta = 2$ ); however, MtCut does not keep this trend for RepLLaMA–RankLLaMA.

Fourth, for RepLLaMA–RankLLaMA, LeCut incorporated the query–item embeddings from RepLLaMA, shows better nDCG@10 values when optimized for prioritizing effectiveness ( $\beta = 0$ ) E.g., LeCut ( $\beta = 0$ ) on TREC-DL 20 achieves the highest nDCG@10 value (0.770).

## 5.3 RLT for pre-trained LM-based re-ranking

To answer **RQ3**, we focus on optimizing RLT methods for different effectiveness/efficiency trade-offs in the context of a widely-used pre-trained language model-based re-ranker (monoT5 [45]) with a lexical retriever (BM25); note that using monoT5 [45] to re-rank the retrieved list returned by RepLLaMA [29] and Splade++ [17] yields worse results; hence we only consider the pipeline of BM25–monoT5. We present the results on TREC-DL 19 and 20 in Figure 9 and 10; we also report the raw result numbers in Table 9. Generally, the observations are similar to the ones from Section 5.1: first, same from Section 5.1, supervised RLT methods do not show a clear advantage in terms of effectiveness and efficiency compared to unsupervised ones; second, the most distribution-based supervised RLT methods outperform the sequential labeling-based method (BiCut) in terms of effectiveness and efficiency; and third, optimizing RLT in a multi-task manner does not lead to consistent improvement in effectiveness and efficiency.

## 6 RELATED WORK

### 6.1 Ranked list truncation

Ranked list truncation (RLT) is also known as query cut-off prediction [10, 25]. For a query and a ranked list of documents, the RLT task is to predict the number of items in the ranked list that should be returned, to optimize a user-defined metric [7]. In other words, this task aims to strike a balance between the overall utility of search results and the user cost of processing search results [60]. The task can potentially benefit IR applications where it is money- and time-consuming to review a returned item, e.g., in patent search [28] and legal search [57, 60]. Early work is mainly assumption-based (hence non-neural-based). This kind of research focuses on modeling score distributions by fitting prior distributions to them [3, 34], which helps identify the best cut-off. However, prior assumptions on score distributions do not always hold as retrieval settings change [27, 60]; hence we do not study this line of studies in our work.

Assumption-free methods, on the other hand, learn to predict the truncation position during training and do not rely on a prior assumption. Lien et al. [27] use a bidirectional LSTM to make a binary prediction (continue or truncate) at each position over a ranked list with a loss function that serves as a proxy to an IR metric as a target metric (e.g., F1). Bahri et al. [6] use a transformer encoder [59] to predict a probability distribution over all candidate cut-off positions, training the model by a loss function that optimizes the expected value of a target metric. Wu et al. [63] improve

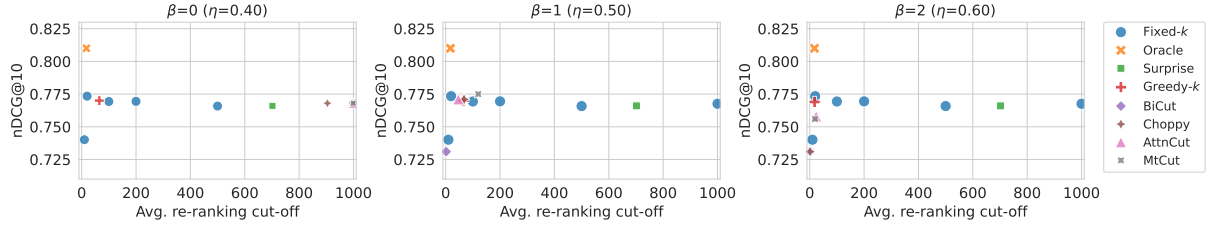


Figure 6: Splade++-RankLLaMA on TREC-DL 19.

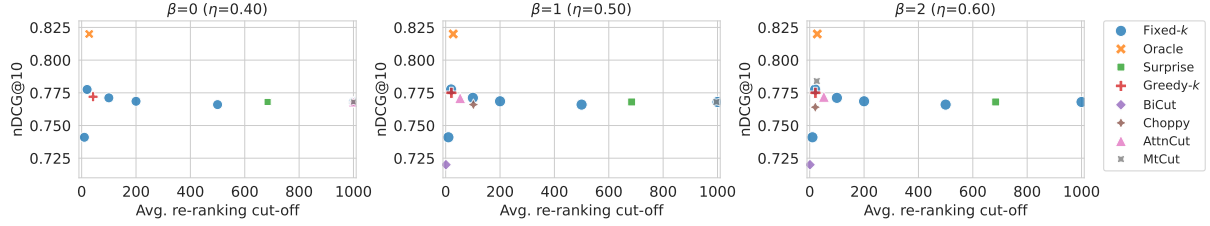


Figure 7: Splade++-RankLLaMA on TREC-DL 20.

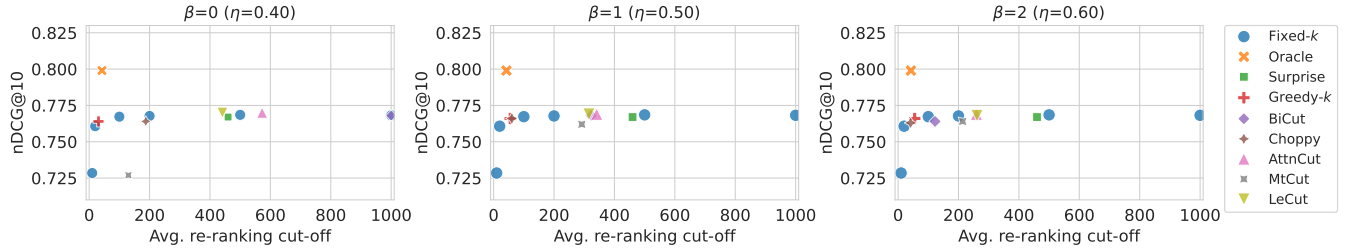


Figure 8: ReplLaMA-RankLLaMA on TREC-DL 20.

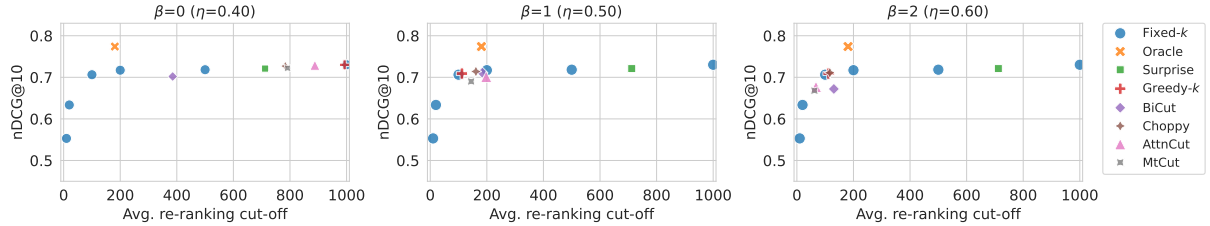


Figure 9: BM25-MonoT5 on TREC-DL 19.

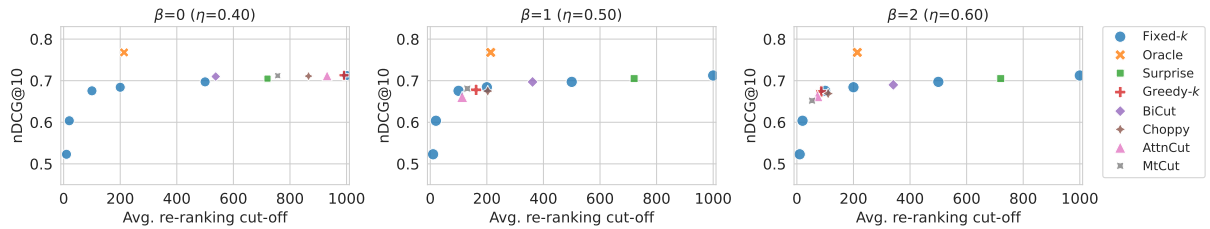


Figure 10: BM25-MonoT5 on TREC-DL 20.

the optimization by employing reward augmented maximum likelihood (RAML) to optimize a target metric directly. Wang et al. [60] identify the negative effect of retrieval bias on RLT and improve the optimization through multi-task learning. Ma et al. [31] introduce semantic features and context information from the retrieval model and jointly optimize RLT and an external Ltr model. Recently, Bahri et al. [7] propose a statistical scoring method based on extreme value theory [47] to calibrate scores in a ranked list, leading to improved RLT quality.

Zamani et al. [66] apply the RLT method from [6] to truncate retrieval lists returned by BM25 for BERT-based re-ranking [44], finding that truncating the retrieval result list to avoid including a large number of non-relevant items in the lower ranks, achieves better re-ranking performance than using fixed cut-offs for all queries.

We differ from Zamani et al. [66] as we provide a systematic and comprehensive study into the use of RLT methods in the context of re-ranking, especially newly emerged LLMs-based re-ranking.



**Table 4: A comparison of RLT methods, optimized for re-ranking effectiveness/efficiency tradeoffs, in predicting re-ranking cut-off points for the BM25-monoT5 pipeline on TREC-DL 19 and 20. Query latency measured in seconds.**

| Method                    | TREC-DL 19 |              |       | TREC-DL 20 |              |       |
|---------------------------|------------|--------------|-------|------------|--------------|-------|
|                           | Avg. k     | nDCG@10      | Lat.  | Avg. k     | nDCG@10      | Lat.  |
| w/o re-ranking            | -          | 0.506        | -     | -          | 0.480        | -     |
| Fixed- $k$ (10)           | 10         | 0.553        | 0.14  | 10         | 0.523        | 0.14  |
| Fixed- $k$ (20)           | 20         | 0.634        | 0.27  | 20         | 0.604        | 0.27  |
| Fixed- $k$ (100)          | 100        | 0.706        | 1.37  | 100        | 0.676        | 1.37  |
| Fixed- $k$ (200)          | 200        | 0.717        | 2.73  | 200        | 0.684        | 2.73  |
| Fixed- $k$ (1000)         | 1000       | <b>0.730</b> | 13.66 | 1000       | <b>0.713</b> | 13.66 |
| Surprise                  | 712.12     | 0.721        | 9.73  | 720.56     | 0.705        | 9.84  |
| Greedy- $k$ ( $\beta=0$ ) | 993        | <b>0.730</b> | 13.57 | 991        | <b>0.713</b> | 13.54 |
| BiCut ( $\eta=0.40$ )     | 385.51     | 0.702        | 5.27  | 537.57     | 0.710        | 7.34  |
| Choppy ( $\beta=0$ )      | 783.56     | 0.727        | 10.70 | 865.56     | 0.711        | 11.82 |
| AttnCut ( $\beta=0$ )     | 888.02     | <u>0.729</u> | 12.13 | 931.26     | <u>0.712</u> | 12.72 |
| MtCut ( $\beta=0$ )       | 791.42     | 0.722        | 10.81 | 756.74     | <u>0.712</u> | 10.34 |
| Greedy- $k$ ( $\beta=1$ ) | 112        | 0.709        | 1.53  | 162        | 0.678        | 2.21  |
| BiCut ( $\eta=0.50$ )     | 184.02     | 0.711        | 2.51  | 361.59     | 0.697        | 4.94  |
| Choppy ( $\beta=1$ )      | 161.37     | 0.714        | 2.20  | 203.17     | 0.675        | 2.78  |
| AttnCut ( $\beta=1$ )     | 197.91     | 0.701        | 2.70  | 112.83     | 0.661        | 1.54  |
| MtCut ( $\beta=1$ )       | 144.51     | 0.690        | 1.97  | 131.06     | 0.681        | 1.79  |
| Greedy- $k$ ( $\beta=2$ ) | 112        | 0.709        | 1.53  | 86         | 0.674        | 1.17  |
| BiCut ( $\eta=0.60$ )     | 130.74     | 0.672        | 1.79  | 340.81     | 0.690        | 4.66  |
| Choppy ( $\beta=2$ )      | 116.72     | 0.711        | 1.59  | 110.56     | 0.669        | 1.51  |
| AttnCut ( $\beta=2$ )     | 67.53      | 0.677        | 0.92  | 72.94      | 0.663        | 1.00  |
| MtCut ( $\beta=2$ )       | 62.02      | 0.668        | 0.85  | 53.56      | 0.652        | 0.73  |
| Oracle                    | 180.91     | 0.774        | 2.47  | 213.67     | 0.768        | 2.92  |

## 6.2 Improving neural re-ranking efficiency

Improving neural re-ranking efficiency has been extensively studied. There are two ideas to improve the efficiency [20]: (i) speed up inference of a neural re-ranker, and (ii) reducing the number of inferences of a neural re-ranker. Approaches to (i) include a simpler re-ranker model [21], distilling knowledge in BERT [15] into a smaller re-ranker [19], pre-computing item representations at indexing time [32], and early-exiting [53, 64]. Studies in (ii) is more related to our work. It includes *multi-stage re-ranking* [35, 46, 62, 68], *candidate pruning* [13, 26, 46] and *early stopping* [65].

*Multi-stage re-ranking* first exploits faster and less effective re-rankers to discard likely non-relevant items and sends fewer candidate items to more expensive re-rankers in later stages. E.g., Zhang et al. [68] first use a feature-based Ltr model to reorder the items returned by BM25 and then send the top- $k$  (applied to all queries) items returned by the faster re-ranker to a BERT re-ranker.

*Candidate pruning* trims the candidate list in the first (or earlier) stage and then forwards the pruned ranked list to the next stage re-ranking. Wang et al. [62] propose a boosting algorithm for jointly learning pruning and ranker stages. Culpepper et al. [13] use a cascade of binary classifiers based on random forests; each classifier is used to predict whether to truncate the given ranked list at a specific cut-off value. Li et al. [26] propose a score-thresholding method, which makes sure the trimmed candidate list produces re-ranking outcomes that satisfy the user-specified error tolerance of an IR evaluation metric.

*Early stopping* is used for an on-the-fly re-ranker to stop scoring the rest of the items in a ranked list, to avoid traversing the entire ranked list. Ying and Huo [65] determine when to stop an on-the-fly re-ranker by comparing the re-ranker scores and a score threshold.

Although candidate pruning methods [13, 26, 62] can be used in our scenario, previous studies study RLT and candidate pruning separately and do not merge these two lines of research [6, 7, 27, 31, 60, 63]. Our focus in this reproducibility study is on specifically reproducing RLT methods in a “*retrieve-then-re-rank*” setup. Exploring all approaches that have the potential to trim the retrieved list and determining the optimal one is beyond the scope of our work.

We also differ from Asadi and Lin [4] and Tonello et al. [58], who investigate improving the efficiency of candidate generation, i.e., first-stage retrieval. In particular, Tonello et al. [58] predict the number of candidate items that should be retrieved by the candidate generation algorithm WAND [8] on a per-query basis. In contrast, we focus on the setting where a first-stage retriever always returns the same number of items, and we improve re-ranking efficiency by truncating retrieved lists.

Our work also differs from Ganguly and Yilmaz [18], who propose variable-depth pooling (VDP) to reduce relevance judgment costs in collection construction; VDP uses query performance prediction (QPP) [2, 36, 37, 39, 40] to predict variable cut-off depths for ranked lists when building a pool of items for a query. In contrast, we focus on using RLT methods in the re-ranking scenario.

## 7 CONCLUSIONS & FUTURE WORK

In this reproducibility study, we examined whether four key findings for ranked list truncation (RLT) for first-stage retrieval hold in the “*retrieve-then-re-rank*” setup. We experimented with RLT methods in three perspectives: (i) assessing RLT methods that optimized to model different trade-offs between effectiveness and efficiency in re-ranking, (ii) investigating the impact of different types of first-stage retrievers on RLT methods, and (iii) investigating the impact of different types of re-rankers on RLT methods.

We showed that findings on RLT for first-stage retrieval do not generalize to the “*retrieve-then-re-rank*” setup well. We found that (i) unsupervised RLT methods generally perform better than supervised ones in terms of the joint consideration of re-ranking effectiveness across various retriever and re-ranker pipeline configurations; (ii) distribution-based supervised RLT methods perform better than their sequential labeling-based counterparts in most cases; (iii) jointly learning RLT with other tasks [60] leads to a marked improvement in terms of re-ranking effectiveness/efficiency for the pipeline of SPLADE++-RankLLaMA, but does not benefit other combinations of other types of retrievers and re-rankers; and (iv) the RLT method incorporating embeddings [31] from a neural retriever leads to better re-ranking effectiveness when optimized for prioritizing effectiveness compared to other supervised RLT methods, but does not improve efficiency markedly.

Suggestions:

We point to two limitations of our study, namely, (i) we only consider point-wise re-rankers, and (ii) we only explore RLT methods in the scenario of two-stage ranking. In future work, we plan to (i) explore RLT for pair-wise and list-wise LLM-based re-rankers [30, 48–

50, 55, 56, 67], (ii) explore RLT for multi-stage ranking, and (iii) explore RLT in the context of conversational search (CS) [1, 38, 41–43, 54].

## ACKNOWLEDGMENTS

We thank Yiding Liu (Baidu Inc.) for insightful discussions that contributed to the development of our idea. This research was partially supported by the China Scholarship Council (CSC) under grant number 202106220041.

## REFERENCES

- [1] Zahra Abbasiantaeb, Chuan Meng, David Rau, Antonis Krasakis, Hossein A Rahmani, and Mohammad Aliannejadi. 2023. LLM-based Retrieval and Generation Pipelines for TREC Interactive Knowledge Assistance Track (iKAT) 2023. In *TREC*.
- [2] Negar Arabzadeh, Chuan Meng, Mohammad Aliannejadi, and Ebrahim Bagheri. 2024. Query Performance Prediction: From Fundamentals to Advanced Techniques. In *ECIR*. Springer, 381–388.
- [3] Avi Arampatzis, Jaap Kamps, and Stephen Robertson. 2009. Where to Stop Reading a Ranked List? Threshold Optimization using Truncated Score Distributions. In *SIGIR*. 524–531.
- [4] Nima Asadi and Jimmy Lin. 2013. Effectiveness/Efficiency Tradeoffs for Candidate Generation in Multi-Stage Retrieval Architectures. In *SIGIR*. 997–1000.
- [5] Arian Askari, Roxana Petcu, Chuan Meng, Mohammad Aliannejadi, Amin Abolghasemi, Evangelos Kanoulas, and Suzan Verberne. 2024. Self-seeding and Multi-intent Self-instructing LLMs for Generating Intent-aware Information-Seeking dialogs. *arXiv preprint arXiv:2402.11633* (2024).
- [6] Dara Bahri, Yi Tay, Che Zheng, Donald Metzler, and Andrew Tomkins. 2020. Choppy: Cut Transformer for Ranked List Truncation. In *SIGIR*. 1513–1516.
- [7] Dara Bahri, Che Zheng, Yi Tay, Donald Metzler, and Andrew Tomkins. 2023. Surprise: Result List Truncation via Extreme Value Theory. In *SIGIR*. 2404–2408.
- [8] Andrei Z Broder, David Carmel, Michael Herscovici, Aya Soffer, and Jason Zien. 2003. Efficient Query Evaluation using a Two-Level Retrieval Process. In *CIKM*. 426–434.
- [9] Sebastian Bruch, Claudio Lucchese, and Franco Maria Nardini. 2023. Efficient and Effective Tree-based and Neural Learning to Rank. *Foundations and Trends in Information Retrieval* 17, 1 (2023), 1–123.
- [10] Daniel Cohen, Bhaskar Mitra, Oleg Lesota, Navid Rekabsaz, and Carsten Eickhoff. 2021. Not All Relevance Scores are Equal: Efficient Uncertainty and Calibration Modeling for Deep Retrieval Models. In *SIGIR*. 654–664.
- [11] Nick Craswell, Bhaskar Mitra, Emine Yilmaz, and Daniel Campos. 2020. Overview of the TREC 2020 Deep Learning Track. In *TREC*.
- [12] Nick Craswell, Bhaskar Mitra, Emine Yilmaz, Daniel Campos, and Ellen M Voorhees. 2019. Overview of the TREC 2019 Deep Learning Track. In *TREC*.
- [13] J Shane Culpepper, Charles LA Clarke, and Jimmy Lin. 2016. Dynamic Cutoff Prediction in Multi-Stage Retrieval Systems. In *Proceedings of the 21st Australasian Document Computing Symposium*. 17–24.
- [14] Donald A Darling. 1957. The Kolmogorov-Smirnov, Cramer-Von Mises Tests. *The Annals of Mathematical Statistics* 28, 4 (1957), 823–838.
- [15] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *NAACL*. 4171–4186.
- [16] Andrew Drozdov, Honglei Zhuang, Zhuyun Dai, Zhen Qin, Razieh Rahimi, Xuanhui Wang, Dana Alon, Mohit Iyyer, Andrew McCallum, Donald Metzler, et al. 2023. PaRaDe: Passage Ranking using Demonstrations with LLMs. In *Findings of EMNLP*. 14242–14252.
- [17] Thibault Formal, Carlos Lassance, Benjamin Piwowarski, and Stéphane Clinchant. 2022. From Distillation to Hard Negative Sampling: Making Sparse Neural IR Models More Effective. In *SIGIR*. 2353–2359.
- [18] Debasis Ganguly and Emine Yilmaz. 2023. Query-specific Variable Depth Pooling via Query Performance Prediction. In *SIGIR*. 2303–2307.
- [19] Luyu Gao, Zhuyun Dai, and Jamie Callan. 2020. Understanding BERT Rankers Under Distillation. In *SIGIR*. 149–152.
- [20] Lukas Gienapp, Maik Fröbe, Matthias Hagen, and Martin Potthast. 2022. Sparse Pairwise Re-ranking with Pre-trained Transformers. In *ICTIR*. 72–80.
- [21] Sebastian Hofstätter, Markus Zlabinger, and Allan Hanbury. 2020. Interpretable & Time-Budget-Constrained Contextualization for Re-Ranking. In *ECAI 2020 24th European Conference on Artificial Intelligence, 29 August-8 September 2020, Santiago de Compostela, Spain-Including 10th Conference on Prestigious Applications of Artificial Intelligence (PAIS 2020)*. IOS Press, 1–8.
- [22] Kalervo Järvelin and Jaana Kekäläinen. 2002. Cumulated Gain-Based Evaluation of IR Techniques. *ACM Transactions on Information Systems* 20, 4 (2002), 422–446.
- [23] Diederik P Kingma and Jimmy Ba. 2015. Adam: A Method for Stochastic Optimization. In *ICLR*.
- [24] Quoc Le and Tomas Mikolov. 2014. Distributed Representations of Sentences and Documents. In *ICML*. PMLR, 1188–1196.
- [25] Oleg Lesota, Navid Rekabsaz, Daniel Cohen, Klaus Antonius Grasserbauer, Carsten Eickhoff, and Markus Schedl. 2021. A Modern Perspective on Query Likelihood with Deep Generative Retrieval Models. In *ICTIR*. 185–195.
- [26] Minghan Li, Xinyu Zhang, Ji Xin, Hongyang Zhang, and Jimmy Lin. 2022. Certified Error Control of Candidate Set Pruning for Two-Stage Relevance Ranking. In *EMNLP*. 333–345.
- [27] Yen-Chieh Lien, Daniel Cohen, and W Bruce Croft. 2019. An Assumption-Free Approach to the Dynamic Truncation of Ranked Lists. In *ICTIR*. 79–82.
- [28] Mihai Lupu and Allan Hanbury. 2013. Patent Retrieval. *Foundations and Trends in Information Retrieval* 7, 1 (2013), 1–97.
- [29] Xueguang Ma, Liang Wang, Nan Yang, Furu Wei, and Jimmy Lin. 2023. Fine-Tuning LLaMA for Multi-Stage Text Retrieval. *arXiv preprint arXiv:2310.08319* (2023).
- [30] Xueguang Ma, Xinyu Zhang, Ronak Pradeep, and Jimmy Lin. 2023. Zero-Shot Listwise Document Reranking with a Large Language Model. *arXiv preprint arXiv:2305.02156* (2023).
- [31] Yixiao Ma, Qingyao Ai, Yueyue Wu, Yunqiu Shao, Yiqun Liu, Min Zhang, and Shaoping Ma. 2022. Incorporating Retrieval Information into the Truncation of Ranking Lists for Better Legal Search. In *SIGIR*. 438–448.
- [32] Sean MacAvaney, Franco Maria Nardini, Raffaele Perego, Nicola Tonello, Nazli Goharian, and Ophir Frieder. 2020. Efficient Document Re-Ranking for Transformers by Precomputing Term Representations. In *SIGIR*. 49–58.
- [33] Sean MacAvaney, Nicola Tonello, and Craig Macdonald. 2022. Adaptive Re-Ranking with a Corpus Graph. In *CIKM*. 1491–1500.
- [34] Raghavan Manmatha, Toni Rath, and Fangfang Feng. 2001. Modeling Score Distributions for Combining the Outputs of Search Engines. In *SIGIR*. 267–275.
- [35] Yoshitomo Matsubara, Thuy Vu, and Alessandro Moschitti. 2020. Reranking for Efficient Transformer-based Answer Selection. In *SIGIR*. 1577–1580.
- [36] Chuan Meng. 2024. Query Performance Prediction for Conversational Search and Beyond. In *SIGIR*.
- [37] Chuan Meng, Mohammad Aliannejadi, and Maarten de Rijke. 2023. Performance Prediction for Conversational Search Using Perplexities of Query Rewrites. In *QPP++2023*. 25–28.
- [38] Chuan Meng, Mohammad Aliannejadi, and Maarten de Rijke. 2023. System Initiative Prediction for Multi-turn Conversational Information Seeking. In *CIKM*. 1807–1817.
- [39] Chuan Meng, Negar Arabzadeh, Mohammad Aliannejadi, and Maarten de Rijke. 2023. Query Performance Prediction: From Ad-hoc to Conversational Search. In *SIGIR*. 2583–2593.
- [40] Chuan Meng, Negar Arabzadeh, Arian Askari, Mohammad Aliannejadi, and Maarten de Rijke. 2024. Query Performance Prediction using Relevance Judgments Generated by Large Language Models. *arXiv preprint arXiv:2404.01012* (2024).
- [41] Chuan Meng, Pengjie Ren, Zhumin Chen, Christof Monz, Jun Ma, and Maarten de Rijke. 2020. RefNet: A Reference-aware Network for Background Based Conversation. In *AAAI*.
- [42] Chuan Meng, Pengjie Ren, Zhumin Chen, Zhaochun Ren, Tengxiao Xi, and Maarten de Rijke. 2021. Initiative-Aware Self-Supervised Learning for Knowledge-Grounded Conversations. In *SIGIR*. 522–532.
- [43] Chuan Meng, Pengjie Ren, Zhumin Chen, Weiwei Sun, Zhaochun Ren, Zhaopeng Tu, and Maarten de Rijke. 2020. DukeNet: A Dual Knowledge Interaction Network for Knowledge-Grounded Conversation. In *SIGIR*. 1151–1160.
- [44] Rodrigo Nogueira and Kyunghyun Cho. 2019. Passage Re-ranking with BERT. *arXiv preprint arXiv:1901.04085* (2019).
- [45] Rodrigo Nogueira, Zhiying Jiang, Ronak Pradeep, and Jimmy Lin. 2020. Document Ranking with a Pretrained Sequence-to-Sequence Model. In *EMNLP*. 708–718.
- [46] Rodrigo Nogueira, Wei Yang, Kyunghyun Cho, and Jimmy Lin. 2019. Multi-Stage Document Ranking with BERT. *arXiv preprint arXiv:1910.14424* (2019).
- [47] James Pickands III. 1975. Statistical Inference Using Extreme Order Statistics. *the Annals of Statistics* (1975), 119–131.
- [48] Ronak Pradeep, Sahel Sharifmoghadam, and Jimmy Lin. 2023. RankVicuna: Zero-Shot Listwise Document Reranking with Open-Source Large Language Models. *arXiv preprint arXiv:2309.15088* (2023).
- [49] Ronak Pradeep, Sahel Sharifmoghadam, and Jimmy Lin. 2023. RankZephyr: Effective and Robust Zero-Shot Listwise Reranking is a Breeze! *arXiv preprint arXiv:2312.02724* (2023).
- [50] Zhen Qin, Rolf Jagerman, Kai Hui, Honglei Zhuang, Junru Wu, Jiaming Shen, Tianqi Liu, Jialu Liu, Donald Metzler, Xuanhui Wang, et al. 2023. Large Language Models are Effective Text Rankers with Pairwise Ranking Prompting. *arXiv preprint arXiv:2306.17563* (2023).
- [51] Stephen Robertson and Hugo Zaragoza. 2009. The Probabilistic Relevance Framework: BM25 and Beyond. *Foundations and Trends in Information Retrieval* 3, 4 (2009), 333–389.
- [52] Devendra Sachan, Mike Lewis, Mandar Joshi, Armen Aghajanyan, Wen-tau Yih, Joelle Pineau, and Luke Zettlemoyer. 2022. Improving Passage Retrieval with Zero-Shot Question Generation. In *EMNLP*. 3781–3797.
- [53] Luca Soldaini and Alessandro Moschitti. 2020. The Cascade Transformer: an Application for Efficient Answer Sentence Selection. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. 5697–5708.
- [54] Weiwei Sun, Chuan Meng, Qi Meng, Zhaochun Ren, Pengjie Ren, Zhumin Chen, and Maarten de Rijke. 2021. Conversations Powered by Cross-Lingual Knowledge. In *SIGIR*. 1442–1451.
- [55] Weiwei Sun, Lingyong Yan, Xinyu Ma, Pengjie Ren, Dawei Yin, and Zhaochun Ren. 2023. Is ChatGPT Good at Search? Investigating Large Language Models as Re-Ranking Agent. In *EMNLP*. 14918–14937.
- [56] 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. *arXiv preprint arXiv:2310.07712* (2023).
- [57] Stephen Tomlinson, Douglas W Oard, Jason R Baron, and Paul Thompson. 2007. Overview of the TREC 2007 Legal Track. In *TREC*.

- [58] Nicola Tonellotto, Craig Macdonald, and Iadh Ounis. 2013. Efficient and Effective Retrieval using Selective Pruning. In *WSDM*. 63–72.
- [59] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention Is All You Need. In *NeurIPS*. 5998–6008.
- [60] Dong Wang, Jianxin Li, Tianchen Zhu, Haoyi Zhou, Qishan Zhu, Yuxin Wen, and Hongming Piao. 2022. MtCut: A Multi-Task Framework for Ranked List Truncation. In *WSDM*. 1054–1062.
- [61] Lidan Wang, Jimmy Lin, and Donald Metzler. 2010. Learning to Efficiently Rank. In *SIGIR*. 138–145.
- [62] Lidan Wang, Jimmy Lin, and Donald Metzler. 2011. A Cascade Ranking Model for Efficient Ranked Retrieval. In *SIGIR*. 105–114.
- [63] Chen Wu, Ruqing Zhang, Jiafeng Guo, Yixing Fan, Yanyan Lan, and Xueqi Cheng. 2021. Learning to Truncate Ranked Lists for Information Retrieval. In *AAAI*, Vol. 35. 4453–4461.
- [64] Ji Xin, Rodrigo Nogueira, Yaoliang Yu, and Jimmy Lin. 2020. Early Exiting BERT for Efficient Document Ranking. In *Proceedings of SustaiNLP: Workshop on Simple and Efficient Natural Language Processing*. 83–88.
- [65] Chengxuan Ying and Chen Huo. 2020. An Adaptive Early Stopping Strategy for Query-based Passage Re-ranking. In *WSDM Cup Reports*.
- [66] Hamed Zamani, Michael Bendersky, Donald Metzler, Honglei Zhuang, and Xuanhui Wang. 2022. Stochastic Retrieval-Conditioned Reranking. In *ICTIR*. 81–91.
- [67] Xinyu Zhang, Sebastian Hofstätter, Patrick Lewis, Raphael Tang, and Jimmy Lin. 2023. Rank-without-GPT: Building GPT-Independent Listwise Rerankers on Open-Source Large Language Models. *arXiv preprint arXiv:2312.02969* (2023).
- [68] Yue Zhang, ChengCheng Hu, Yuqi Liu, Hui Fang, and Jimmy Lin. 2021. Learning to Rank in the Age of Muppets: Effectiveness–Efficiency Tradeoffs in Multi-Stage Ranking. In *Proceedings of the Second Workshop on Simple and Efficient Natural Language Processing*. 64–73.
- [69] Honglei Zhuang, Zhen Qin, Kai Hui, Junru Wu, Le Yan, Xuanhui Wang, and Michael Bendersky. 2023. Beyond Yes and No: Improving Zero-Shot LLM Rankers via Scoring Fine-Grained Relevance Labels. *arXiv preprint arXiv:2310.14122* (2023).
- [70] Shengyao Zhuang, Bing Liu, Bevan Koopman, and Guido Zuccon. 2023. Open-source Large Language Models are Strong Zero-shot Query Likelihood Models for Document Ranking. *arXiv preprint arXiv:2310.13243* (2023).
- [71] Shengyao Zhuang, Honglei Zhuang, Bevan Koopman, and Guido Zuccon. 2023. A Setwise Approach for Effective and Highly Efficient Zero-shot Ranking with Large Language Models. *arXiv preprint arXiv:2310.09497* (2023).