

A Passage-Based Approach to Learning to Rank Documents

Eilon Sheerit

seilon@campus.technion.ac.il
Technion

Anna Shtok

annie.shtok@gmail.com

Oren Kurland

kurland@ie.technion.ac.il
Technion

ABSTRACT

According to common relevance-judgments regimes, such as TREC's, a document can be deemed relevant to a query even if it contains a very short passage of text with pertinent information. This fact has motivated work on passage-based document retrieval: document ranking methods that induce information from the document's passages. However, the main source of passage-based information utilized was passage-query similarities. We address the challenge of utilizing richer sources of passage-based information to improve document retrieval effectiveness. Specifically, we devise a suite of learning-to-rank-based document retrieval methods that utilize an effective ranking of passages produced in response to the query; the passage ranking is also induced using a learning-to-rank approach. Some of the methods quantify the ranking of the passages of a document. Others utilize the feature-based representation of passages used for learning a passage ranker. Empirical evaluation attests to the clear merits of our methods with respect to highly effective baselines. Our best performing method is based on learning a document ranking function using document-query features and passage-query features of the document's passage most highly ranked.

1 INTRODUCTION

The ad hoc retrieval task is ranking documents in a corpus in response to a query by presumed relevance to the information need the query represents. Often, documents are deemed relevant even if they contain only a short passage with pertaining information; e.g., by TREC's relevance judgment regime [52].

As a result, there has been a large body of work on *passage-based document retrieval*: utilizing information induced from document passages to rank the documents; e.g., [4, 7, 25, 32, 55]. The most commonly used passage-based document retrieval methods rank a document by the highest query similarity exhibited by any of its passages [4, 7, 25, 32, 55] and by integrating this similarity with the document-query similarity [4, 7, 55].

The passage-query (surface level) similarity is one out of many possible estimates for passage relevance. Indeed, various passage-relevance estimates were devised for the task of passage retrieval, a.k.a focused retrieval; e.g., [5, 8–10, 15, 16, 22, 27, 39, 42, 43, 48,

56, 58]. That is, passages are ranked in response to a query using passage-relevance estimates. The merits of integrating the estimates using learning-to-rank (LTR) approaches were also demonstrated [5, 9, 10, 39, 56, 58].

Motivated by the (recent) progress in devising effective passage retrieval methods, specifically, using LTR methods, and the fact that the main passage-based information used by most passage-based document retrieval methods is confined to passage-query similarities, we address the following challenge: devising LTR methods for document retrieval that utilize an effective query-based passage ranking. Some of the methods we present are not based on any assumptions regarding the passage retrieval approach used to rank passages. Others are based on the premise that passages were ranked in response to the query using an LTR method that utilizes passage-based features. A case in point, the most effective LTR-based document retrieval method that we present uses both document-based and passage-based features; the latter are those of the document's passage which is the most highly ranked by an LTR method used to rank passages.

Each of the methods we present can be viewed as a conceptual analog, or generalization, of previously proposed approaches for either (i) passage-based document retrieval, where these approaches do not utilize learning-to-rank or feature-based representations, or (ii) cluster-based document retrieval.

In addition to presenting novel passage-based document retrieval methods, we also propose new features for learning-to-rank passages. These features are query-independent passage relevance priors adapted from work on document retrieval over the Web [2].

Extensive empirical evaluation shows that our passage-based document retrieval approaches significantly outperform strong baselines. Further analysis demonstrates the importance of (i) utilizing an effective passage ranking, and (ii) using information induced from the document's passage that is the most highly ranked. In addition, we demonstrate the merits of using the query-independent passage features we propose for the task of passage retrieval. Specifically, integrating these features with previously proposed ones in a learning-to-rank approach results in passage retrieval performance that transcends the state-of-the-art.

Our contributions can be summarized as follows:

- A study of different methods that utilize passage-based features in a learning-to-rank approach for ranking documents.
- The utilization of an effective passage ranking for inducing document ranking, or in other words, addressing the question of how passage ranking can be transformed to document ranking.
- Some of our methods conceptually generalize previously proposed passage-based document retrieval methods which do not use learning-to-rank or feature-based representation.
- Some of our methods are conceptual reminiscent of cluster-based document retrieval approaches. This is the first work,

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

© 2018 Association for Computing Machinery.

to the best of our knowledge, to make the connection between passage-based document retrieval and cluster-based document retrieval.

- Attaining state-of-the-art retrieval performance across different collections and different feature sets.
- Demonstrating the effectiveness for passage retrieval of using passage-relevance priors adopted from work on document-relevance priors in Web retrieval.

2 RELATED WORK

The line of work most related to ours is on passage-based document retrieval [3, 4, 7, 13, 21, 24, 25, 28, 30, 32, 41, 44, 53–55]. As already noted, the most commonly used passage-based document retrieval methods are ranking a document by the maximum query-similarity of its passages [4, 7, 24, 25, 32, 44, 55] and by interpolating this similarity with the document-query similarity [4, 7, 44, 55]. We show that our best-performing methods substantially outperform a highly effective method that integrates document-query and passage-query similarities [4].

In [54], features based on passage-query similarities were used to learn a document ranker [54]. The induced ranking was fused with a query-similarity-based document ranking. Our FPD method, described in Section 3, generalizes this approach by using many more passage features, integrating the resultant passage-based document ranking with that produced by learning to rank documents, and applying state-of-the-art learning-to-rank approaches. While FPD is highly effective, it is outperformed by our best performing method.

Recently, Yulianti et al. [57] presented a method that selects (or generates) a passage from a document in response to a query using information induced from a community question answering system. Then, features of the passage (not necessarily those used for selecting the passage) along with document features are used to represent the document. This approach is reminiscent of our JPDs method which uses passage features and document features to represent a document. There are, however, major differences between the two. Our method is not based on an external resource. Furthermore, we utilize passage ranking that is induced using a learning-to-rank approach with passage features while in [57] this is not the case. In addition, the passage features used in our method are the same as those used for ranking passages which is not the case in Yulianti et al. [57]. We demonstrate the merits of using the passage features that are used for (effective) passage ranking to represent a document. We also show the merits of using passage-relevance prior estimates adopted from work on Web retrieval to rank passages. Some of these estimates were used by Yulianti+al et al. [57] to rank documents but not passages.

Recently, a neural-network approach was presented for passage-based document retrieval [14]. Passage-query relevance signals (scores) are estimated using neural-network matching models and then aggregated to yield a document score. A difference with several of our models, in addition to using neural networks rather than a feature-based approach, is that ranking induced over passages from different documents is not utilized. A feature-based learning-to-rank baseline used in this work [14] represents a document using its

features and the average, maximum and minimum values of query-similarities of its constituent passages. Therefore, this baseline is conceptually reminiscent of our JPDm method which uses various aggregates of the feature values of document’s passages together with the document features to represent documents. We show that there are passage-based features much more effective than passage-query similarities for estimating passage relevance, and accordingly, use aggregates of these features’ values to represent documents.

Some passage-based document retrieval methods use query expansion [30, 32] or inter-passage similarities [28, 53, 54]. Integrating query expansion and information induced from inter-passage similarities in our approaches is an interesting future venue.

Passage-based document retrieval approaches utilize term proximity information by the virtue of using passages. There are many other approaches for utilizing term proximities [30, 33, 34, 36, 37, 40, 51, 60]. We show that our best performing method outperforms a state-of-the-art term proximity model: the sequential dependence model from the Markov Random Field framework [36].

The vast majority of previous work on passage-based document retrieval has focused on using passages marked prior to retrieval time. There are some methods that simultaneously mark passages and use them for retrieval [13, 25, 41]. Hence, our methods are not committed to a specific approach of passage markup.

To implement and evaluate our passage-based document retrieval methods, we use a passage ranking method that is based on learning-to-rank. Some of the features we use for passage retrieval are adopted from work on retrieving sentences to create snippets [39] and retrieving sentences (and more generally passages) as answers to non-factoid questions [9, 27, 56]. We show that passage retrieval performance can be significantly improved if we also use query-independent passage relevance priors adapted from work on devising document relevance priors for Web retrieval [2]. Query-independent sentence priors different than ours, mainly based on opinion/sentiment analysis, were used in past work on sentence retrieval [15]. More generally, there is a big body of work on retrieving passages; e.g., [8, 10, 22, 26, 27, 41, 48]. Our focus is different: we devise methods that utilize passage retrieval to improve document retrieval. Yet, we empirically show that the passage retrieval method we use in our document retrieval methods outperforms state-of-the-art passage retrieval approaches. Still, as already noted, our document retrieval methods are not committed to a specific passage retrieval approach.

3 RETRIEVAL FRAMEWORK

Our goal is to rank documents in corpus D with respect to query q . We devise document retrieval methods that utilize information induced from document passages. We assume that passages were marked in documents using some approach; $g \in d$ indicates that passage g is part of document d . The retrieval methods we present are not dependent on the type of passages used. If S is a document set, $G(S)$ denotes the ranked list of all passages of documents in S , where ranking was performed using some passage retrieval method.

Let D_{init} be an initially retrieved document list produced in response to q by using some retrieval method; e.g., in the experiments reported in Section 4 we use standard language-model-based retrieval. Then, a learning-to-rank (LTR) method [31] is used to re-rank D_{init} ; the resultant ranked list is denoted D_{LTR} . The only assumption we make about the LTR method is that it uses a feature-based vector representation, $\vec{v}_{(d,q)}$, for every pair of a document d and the query q .

We devise document ranking methods that re-rank D_{LTR} using information induced from the ranked list $G(D_{LTR})$ of all passages in documents in D_{LTR} .¹ Some of the approaches we present do not depend on the passage ranking method used to produce $G(D_{LTR})$. Others are based on the assumption that the ranking is induced using an LTR approach applied to passages; a pair of passage g and query q is represented using the feature vector $\vec{v}_{(g,q)}$. The basic premise is that effective passage ranking can be utilized to improve document ranking.

3.1 Passage-Based Document Ranking

We now present five passage-based document retrieval approaches that can be used to re-rank D_{LTR} . They mainly differ by the way they utilize information about the ranking of passages in $G(D_{LTR})$. These methods are either inspired by, or bear important connections to, existing passage-based and cluster ranking approaches.

3.1.1 A fusion-based approach. The first method we consider is conceptually reminiscent of a commonly used passage-based document retrieval approach. The approach linearly interpolates the document-query similarity score with the highest query similarity score of a passage in the document [4, 7, 55].

Here, instead of relying on query similarities, we use the ranking of documents in D_{LTR} and that of the passages in $G(D_{LTR})$ to induce document and passage retrieval scores, respectively. Specifically, we apply the rank-to-score transformation used in the highly effective reciprocal rank fusion method [11]. That is, the score assigned to item x , passage or document, with respect to the list L it is in, $G(D_{LTR})$ or D_{LTR} , is:

$$Score_L(x) \stackrel{def}{=} \frac{1}{v + r_L(x)};$$

$r_L(x)$ is x 's rank in L ; the top item is at rank 1; v is a free parameter.

The final retrieval score of document d ($\in D_{LTR}$) is:

$$Score(d; q) \stackrel{def}{=} \alpha Score_{D_{LTR}}(d) + (1 - \alpha) \max_{g \in d} Score_{G(D_{LTR})}(g); \quad (1)$$

α is a free parameter. Thus, d is ranked high if it was originally ranked high in D_{LTR} and at least one of its passages was ranked high in $G(D_{LTR})$.

The method just presented essentially applies the reciprocal rank fusion approach to fuse two rankings of the documents in D_{LTR} and is therefore denoted **RRF**. The first is the LTR-based ranking of D_{LTR} . That is, documents are ranked using a ranking function learned based on document-only features. The second ranking is based on the highest rank in $G(D_{LTR})$ of a document's passage. In other words, the retrieval score of a document with respect to this ranking is based on the reciprocal rank of its passage that is the

highest ranked. Note that the method is agnostic to the retrieval methods that were used to produce D_{LTR} and $G(D_{LTR})$; e.g., these need not even be LTR methods. All the method relies on is the ranking of documents and the ranking of passages of these documents.

3.1.2 Utilizing various passage-ranking statistics. The RRF method utilizes only the highest ranked passage of a document to assign its final retrieval score in Equation 1. The next method, **SMPD** ("statistics about multiple passages per document"), ranks a document by utilizing various *statistics* regarding the ranking of the document's passages in $G(D_{LTR})$.

The feature vector used to represent a query-document pair is:

$$\vec{v}_{(d,q)}^{SMPD} \stackrel{def}{=} \vec{v}_{(d,q)} \oplus \vec{v}'_{(g \in d, q)}.$$

$\vec{v}_{(d,q)}^{SMPD}$ is the concatenation of $\vec{v}_{(d,q)}$: the original feature vector used to learn and apply the ranking function that served to induce D_{LTR} and $\vec{v}'_{(g \in d, q)}$: a vector composed of passage-based estimates. The estimates are the (i) maximum (max), (ii) minimum (min), (iii) average (avg), and (iv) standard deviation (std) of $Score_{G(D_{LTR})}(g)$ for $g \in d$; (v) the fraction of passages in d that are among the 50 (top50) and (vi) 100 (top100) highest ranked passages in $G(D_{LTR})$; and, (vii) the number of passages in d (numPsg).

The rationale behind the SMPD method is to augment the original document-query representation with "statistics" about the potential relevance of its passages. The premise is that the relative ranking of passages in $G(D_{LTR})$ can attest to their relevance to some extent. While SMPD is based on the fact that D_{LTR} was indeed produced using an LTR approach, it is not committed to a specific passage ranking method used to produce $G(D_{LTR})$.

We note an interesting conceptual connection between SMPD and a cluster-based document retrieval method [29]. The method ranks clusters of similar documents using measures that quantify the ranking of their constituent documents in a document ranking. In SMPD, we rank a document using measures that quantify the ranking of its constituent passages.

3.1.3 Joint passage-document representation using a single passage.

The next method, **JPDs** ("joint passage document with a single passage"), similarly to the RRF method, uses d 's passage g_{max} that is the highest ranked in $G(D_{LTR})$. However, JPDs does not rely on g_{max} 's absolute rank in $G(D_{LTR})$, but only on the fact that it is the highest ranked among d 's passages. JPDs is based on the premise that both D_{LTR} and $G(D_{LTR})$ were produced using LTR methods with feature vectors $\vec{v}_{(d,q)}$ and $\vec{v}_{(g_{max},q)}$, respectively. These two feature vectors are concatenated, and the resultant feature vector

$$\vec{v}_{(d,q)}^{JPDs} \stackrel{def}{=} \vec{v}_{(d,q)} \oplus \vec{v}_{(g_{max},q)},$$

is used for learning a ranker.

An important principle underlying JPDs is to avoid *metric divergence* [36]. That is, the features used to estimate the relevance of the document's passage that is presumably the most relevant — according to $G(D_{LTR})$'s ranking — are used directly, along with document-based features, to learn a document ranking function.

JPDs could be viewed as a conceptual generalization of the approach of smoothing a document language model with that induced from its passage which is the most similar to the query [3]. That is, both approaches augment the document representation

¹Note that these passages are also the passages of documents in D_{init} since D_{LTR} is a re-rank of D_{init} .

with information about its passage which is either the most query similar [3] or the most highly ranked using a learning-to-rank approach (JPDs). The difference is unsupervised method [3] vs. a supervised method (JPDs), and in accordance, representations (language models vs. feature vectors) and their integration (linear interpolation vs. concatenation).

3.1.4 Joint passage-document representation using multiple passages.

The JPDs method uses information induced from a single passage of d to augment the document-query feature-vector representation. We next consider an alternative, “joint passage document with multiple passages” — **JPDm** in short. The document-query representation in JPDm utilizes information induced, potentially, from multiple passages. Specifically, we define a feature vector, $agg_{g \in d}(\vec{v}_{(g,q)})$, based on the same passage features used to represent passages in the LTR method that produced $G(D_{LTR})$. Each feature value in $agg_{g \in d}(\vec{v}_{(g,q)})$ is the aggregate of the corresponding feature values of all d ’s passages. The feature vector is then concatenated with the original document-query feature vector

$$\vec{v}_{(d,q)}^{JPDm} \stackrel{def}{=} \vec{v}_{(d,q)} \oplus agg_{g \in d}(\vec{v}_{(g,q)});$$

$\vec{v}_{(d,q)}^{JPDm}$ is used for learning a document ranking function. The resultant methods are termed **JPDm-avg**, **JPDm-max** and **JPDm-min** when using the average, maximum and minimum aggregate functions, respectively. We note that JPDm is the only approach we consider which does not use the ranking of passages in $G(D_{LTR})$.

It is important to highlight an additional difference between the JPDm and SMPD methods, as both augment the document-query feature vector for learning a document ranking function with information induced from multiple passages in the document. While SMPD uses statistics mainly about the ranking of the document’s passages, JPDm utilizes passage-level features which were used to learn a passage ranker.

There is an interesting conceptual connection between JPDm and the ClustMRF method that ranks clusters of similar documents by the presumed percentage of relevant documents they contain [46]. In ClustMRF, clusters are represented using aggregates of feature values of their constituent documents — e.g., aggregates of document-query similarity scores, document-relevance prior estimates and more. JPDm represents documents using aggregates of feature values of their constituent passages.

Finally, we note the important difference between JPDs and JPDm. In JPDs, the passage-based features that are added to the document features represent a single passage; this is the document’s most highly ranked passage. In contrast, in JPDm, the passage-based features used to augment the document features do not represent a single passage: these are aggregates, over the document’s passages, of feature values used in the passages’ feature-vector representations. For example, in JPDm-avg, a single passage-based feature value would be the average feature value — where average is computed over the document’s passages — for some feature in the feature-vector representation of the documents’ passages.

3.1.5 Two-stage retrieval. To further study the merits of simultaneously using document and passage features to learn a document ranking function as in the JPDs and JPDm methods presented above, we next explore the **FPD** method (“first passage then document”).

A *document* ranking function is learned by representing the document-query pair with $\vec{v}_{(g_{max},q)}$ — the feature vector for the document’s passage g_{max} that is the most highly ranked in $G(D_{LTR})$. That is, the learned document ranker utilizes only passage-based features. The ranker is then used to re-rank D_{LTR} . The resultant ranking is fused with D_{LTR} ’s original ranking using the reciprocal rank approach as in RRF. See Section 3.1.1 for further details².

It is important to contrast the FPD and RRF methods. Both fuse the original ranking of D_{LTR} with a ranking based on utilizing passage-based information. The difference is the type of passage-based information used. While RRF utilizes the rank in $G(D_{LTR})$ of the document’s most highly ranked passage to directly induce document ranking, FPD utilizes the passage-query feature vector of this passage to learn and apply a document ranker.

We further note that FPD depends on the fact that $G(D_{LTR})$ was induced using an LTR approach. In contrast, FPD is not committed to a specific retrieval method used to induce D_{LTR} .

4 EXPERIMENTAL SETTING

The datasets used for experiments are specified in Table 1. ROBUST, WT10G, GOV2 and ClueWeb are TREC datasets. ROBUST mostly contains newswire documents. WT10G is a small Web corpus. GOV2 is a crawl of the .gov domain. ClueWeb is a large-scale (noisy) Web collection. For ClueWeb we removed from the initial document rankings, described below, documents with a Waterloo’s spam classifier score below 50 [12].

The TREC datasets do not have passage-level relevance judgments that are needed for learning a passage-ranking method. Thus, to learn a passage ranker we used the INEX dataset. The learned ranker was utilized by our passage-based document retrieval methods over all datasets. The INEX dataset was used for the focused (passage) retrieval tracks in 2009 and 2010 [1, 20]. It includes relevance judgments for virtually every character in a relevant document; that is, annotators marked the pieces of relevant text in relevant documents. The dataset contains English Wikipedia documents from which we removed all XML tags; i.e., we treated the documents as plaintext. We use this dataset not only for learning a passage ranker, but also for evaluating the effectiveness of the learned ranker, as well as evaluating the effectiveness of our passage-based document retrieval methods in addition to the evaluation performed over the TREC datasets.

The passage features we propose are also used for learning and evaluating a passage ranker over the AQUAINT collection which was used for the novelty tracks in TREC 2003 and 2004 [49, 50]. In these tracks, relevant documents have sentence-level relevance judgments. To perform sentence (passage) retrieval using the queries in both tracks, we follow the experimental setting in the 2003 track and rank the sentences in the set of relevant documents that were provided to participants.

Titles of topics served for queries. (Queries with no relevant documents in the qrels were removed.) The Indri toolkit was used for

²Experiments — actual numbers are omitted as they convey no additional insight — showed that simply using the passage-based document ranking without the additional fusion often yields performance (substantially) inferior to that of FPD.

Table 1: Datasets used for experiments.

Corpus	Data	# of docs	Avg doc. length	Queries
ROBUST	Disks 4&5-CR	528,155	479	301-450, 601-700
WT10G	WT10g	1,692,096	607	451-550
GOV2	GOV2	25,205,179	930	701-850
ClueWeb	ClueWeb09 (Category B)	50,220,423	807	1-200
INEX	2009&2010	2,666,190	552	2009001-2009115, 2010001-2010107
AQUAINT	AQUAINT	1,033,461	436	N1-N100

all experiments³. We applied Krovetz stemming to queries, documents (and their passages) and removed stopwords on the INQUERY list only from queries. We used non-overlapping fixed-length windows of 300 terms for passages in our document retrieval methods. Such passages were shown to be effective for passage-based document retrieval [25]. In Section 5 we study the effect of passage length on passage retrieval performance.

Our main experiments are conducted with two learning-to-rank (LTR) methods for ranking documents and passages: LambdaMART [6] (LMart in short)⁴ or a linear RankSVM [23]⁵ (SVM in short). LambdaMART was trained for NDCG@10. In Section 5.1.7 we present experimental results for two additional learning-to-rank methods.

We measure the similarity between texts x and y (e.g., a query, a document or a passage) using the minus cross entropy between the unigram language models induced from them:

$$Sim(x, y) \stackrel{def}{=} \exp(-CE(\theta_x^{MLE} || \theta_y^{Dir})); \quad (2)$$

θ_x^{MLE} is the unsmoothed maximum likelihood estimate induced from x and θ_y^{Dir} is a Dirichlet smoothed language model induced from y [59].

The two-tailed paired t-test with a 95% confidence level was used to determine statistically significant retrieval performance differences. We applied Bonferroni correction for multiple hypothesis testing; i.e., when comparing a method with multiple baselines.

4.1 Document Retrieval

We use a standard (unigram) language model approach (**LM**) to retrieve an initial document list D_{init} of 1000 documents for q : document d is scored by $Sim(q, d)$. We then (re-)rank D_{init} using an LTR method to obtain D_{LTR} ; **init-LTR** denotes this ranking. Since some of the datasets used for evaluation do not have hyperlink and hypertext information, we only use highly effective content-based features. Specifically, the first three features in the document-query feature vector $\vec{v}_{(d, q)}$ are those of the sequential dependence model (SDM) from the Markov Random Field (MRF) framework [36]: unigrams, ordered bigrams and unordered bigrams (biterns). SDM is a state-of-the-art term-proximity model. The next three features are the most effective *document relevance priors* reported in [2]: (i) **SW1** and (ii) **SW2** are the fraction of terms in d that are stopwords on the INQUERY list, and the fraction of stopwords on the INQUERY list that appear in d respectively, and (iii) the entropy,

Ent, of the term distribution in d . High presence of stopwords, and high entropy, presumably attests to rich use of language and therefore to content breadth [2]. In Section 5.1.8 we also present experimental results when using the MSLR⁶ features used in the LETOR datasets.

The set of all passages in documents in D_{LTR} is ranked to yield $G(D_{LTR})$. The same LTR method used to produce D_{LTR} is used to produce $G(D_{LTR})$ with the passage-based features described in Section 4.2. Then, D_{LTR} is re-ranked using the document retrieval methods from Section 3 that utilize $G(D_{LTR})$. We use MAP and p@10 to evaluate document retrieval performance.

Baselines. Recall that D_{LTR} was attained by re-ranking D_{init} using an LTR approach; i.e., the set of documents in these two lists is the same. All the baselines we describe and our passage-based document retrieval methods from Section 3 are used to rank this document set.

The initial language-model-based ranking of D_{init} , denoted **LM**, is the first baseline. The second is the initial LTR-based ranking of D_{LTR} , **init-LTR**. MRF’s **SDM** with its three features [36] also serves as a reference comparison. SDM is a special case of the LTR method used to induce D_{LTR} where document relevance priors are not used. Another reference comparison is **DocPsg** [4] where document d is scored with $\lambda Sim(q, d) + (1 - \lambda) \max_{g \in d} Sim(q, g)$; the value of λ is negatively correlated with d ’s length which serves as a document homogeneity measure [4]. DocPsg is an effective representative of the approach of interpolating document-query and passage-query similarity estimates [4, 7, 55].

4.2 Features for Learning to Rank Passages

All our passage-based document ranking approaches (except for JPDm) utilize a ranking of the documents’ passages; i.e., the ranked list $G(D_{LTR})$. We now turn to describe the features used for learning a passage ranker. Some of these are novel to this study. The features are estimates of passage g ’s relevance to the query q . Let d_g denote g ’s ambient document which we assume is part of a document set S_{doc} retrieved for q . S_{psg} denotes the set of passages of documents in S_{doc} . If S_{doc} is the set of documents in D_{LTR} , the list we aim to re-rank, then S_{psg} is $G(D_{LTR})$.

The **PsgQuerySim** feature is the (normalized) passage-query similarity: $\frac{Sim(q, g)}{\sum_{g' \in S_{psg}} Sim(q, g')}$. Since passages are relatively short, the ambient document can provide context in estimating query

³www.lemurproject.org

⁴https://code.google.com/p/ijforests/.

⁵https://www.cs.cornell.edu/people/tj/svm_light/svm_rank.html.

⁶www.research.microsoft.com/en-us/projects/mslr

similarities (cf. [43]): **DocQuerySim** is $\frac{Sim(q, d_g)}{\sum_{d' \in S_{doc}} Sim(q, d')}$. Additional document-based features are the maximum, average, and standard deviation of $Sim(q, g')$ for $g' \in d_g$: **MaxPDSim**, **AvgPDSim** and **StdPDSim**, respectively. The longer g is with respect to d_g , the less reliance on document-based query-similarity information is called for [4]. Therefore, the ratio between g 's and d_g 's lengths serves as a query independent feature: **LengthRatio**.

Passages (if exist) that precede (g_{pre}) and follow (g_{follow}) g in d_g provide focused context for g [16]. Hence, we use **QuerySimPre** and **QuerySimFollow**: $Sim(q, g_{pre})$ and $Sim(q, g_{follow})$, respectively. If g is the first or the last passage in the document, we use $Sim(q, g)$ for QuerySimPre and QuerySimFollow, respectively.

The next features – the use of which for passage retrieval is novel to this study – are query-independent passage relevance priors. These are adopted from work on document relevance priors in Web search [2]. Specifically, we use the entropy (Ent) and stopwords (SW1, SW2) features described above, but now for passages rather than documents.

The passage independent feature **QueryLength** is the number of unique terms in the query. This feature can potentially help to improve the performance of non-linear rankers (cf., [35]).

The next features are adopted from work on selecting sentences for results' snippets [39]. These were also used to retrieve sentences (passages) for questions [9, 56]. **ExactMatch** is true if q is a substring of g and false otherwise. **TermOverlap** and **SynonymsOverlap** are the fraction of query terms and their synonyms (determined using Wordnet) in g . **PsgLength** is the number of terms in g after removing stopwords, and **PsgLocation** is g 's position (in terms of passages) in d_g over the number of d_g 's passages.

We also compare g with q using the following three semantic-similarity measures utilized for sentence-answer retrieval [56]. (The first two were also used in [9].) The **ESA** similarity [19] is computed by using, separately, q and the 20 terms in g with the highest TF.IDF values for query likelihood retrieval over the INEX Wikipedia collection. The cosine measure is used to compare the lists of min-max normalized retrieval scores of the top-100 documents.

W2V is the cosine similarity between the centroid of the Word2Vec vectors representing q 's terms and the centroid of the Word2Vec vectors representing g 's terms. We used the 300 dimensional newswire-based Word2Vec vectors from <https://code.google.com/p/word2vec/>.

Entity is the Jaccard coefficient between the set-based entity representations of q and g . Wikipedia entities (i.e., titles) marked with a confidence level ≥ 0.1 by TagMe [17] were used.

4.2.1 Evaluating passage retrieval. Most of our passage-based document ranking methods rely on the ranking of document passages. Hence, we also evaluate the effectiveness of the learned passage ranker using the INEX and AQUAINT datasets – this is a focused (passage) retrieval task. For INEX, the set S_{psg}^{init} , of all passages of documents in the language-model-based initially retrieved document list D_{init} , is ranked; the top-1500 passages are evaluated using MAiP and iP[x]: precision at recall level $x \in \{.01, .1\}$ [1, 20]. These evaluation measures were devised for the focused retrieval task where the percentage of relevant information in a passage is accounted for. For AQUAINT, following the novelty track in 2013 [50], we set D_{init} to be the provided set of relevant documents, and S_{psg}^{init} is the set of all sentences in these documents which are ranked

using our passage ranker. The top 1500 ranked sentences are evaluated using MAP and p@10. (The tracks provided sentence-level binary relevance judgments.)

We use the following baselines for passage ranking. The first method, **QSF** (“query-similarity fusion”) [7, 8], scores g by $(1 - \lambda) \frac{Sim(q, g)}{\sum_{g' \in S_{psg}^{init}} Sim(q, g')} + \lambda \frac{Sim(q, d_g)}{\sum_{d' \in D_{init}} Sim(q, d')}$; λ is a free parameter. The two components of this interpolation are among the features used above for learning a passage ranker.

A tf.idf-based positional model was used for passage retrieval [8]. We use a language-model-based positional approach [33], **PLM**, with a Gaussian kernel, as other methods also utilize language models: g is scored by $\lambda \frac{Sim(q, i_{max}(g))}{\sum_{g' \in S_{psg}^{init}} Sim(q, i_{max}(g'))} + \beta \frac{Sim(q, g)}{\sum_{g' \in S_{psg}^{init}} Sim(q, g')} + (1 - \lambda - \beta) \frac{Sim(q, d_g)}{\sum_{d' \in D_{init}} Sim(q, d')}$; $i_{max}(g)$ is the position in g whose Dirichlet induced language model yields the highest query similarity among all positions i in g ; λ and β are free parameters. Using PLM as a feature in our passage ranking approach showed no merit.

We adapt the **owpc** method [5], originally used to rank structured XML elements, as an additional baseline. For compliance with our setting, all features except for those which rely on XML structure are used in the two LTR methods used for all experiments. Most features rely on the query-similarity of the passage and its ambient document; most of the features described above, which we use for learning a passage ranker, were not utilized.

The state-of-the-art LTR-based baseline, **MKS**, utilizes all the features proposed in [56] for retrieving answer sentences to non-factoid questions. Our passage ranker utilizes some of these features.

The LTR-based approaches, owpc, MKS and our methods, are used to re-rank the top 1500 passages retrieved by QSF which is considered an effective method. Applying LTR methods on an initially retrieved list is common practice [31]; specifically, the list size, for document retrieval, is often the same as that of the number of documents to be retrieved (e.g., 1000); hence, LTR methods often operate as re-ranking approaches. Similarly, the 1500 threshold used here for passage retrieval corresponds to the standard passage list size used in the focused retrieval track of INEX [1, 20].

4.3 Additional Experimental Details

As already noted, we use the INEX dataset to train a passage ranker with the features described in Section 4.2. The ranker is also used for passage-based document retrieval over the TREC corpora which lack focused (passage) relevance judgments. To learn a ranker, all passages of documents in the initial language-model-based document list retrieved from INEX, D_{init} , are ranked using the QSF method described in Section 4.2.1; thus, D_{init} serves for the set S_{doc} in Section 4.2. The top 1500 passages serve for training. We explored a few binary/graded passage relevance-grade definitions for learning a passage ranker. These use the fraction of relevant characters in a passage, denoted $RFrac$. A bucket-based approach which produces five relevance grades resulted in effective performance of our passage ranker and the owpc and MKS baselines (see Section 4.2.1 for details): 0: $RFrac < .10$; 1: $.10 \leq RFrac < .25$; 2: $.25 \leq RFrac < .50$; 3: $.50 \leq RFrac < .75$; 4: $.75 \leq RFrac$.

To learn a passage ranking function for the sentence retrieval (ranking) task over AQUAINT, we use the sentences’ binary relevance judgments as relevance grades.

For the JPDs passage-based document retrieval approach, the DocQuerySim passage feature is not used, as it is the unigram feature of SDM that is used as a document-based feature. For the JPDm-avg and JPDm-max passage-based document retrieval methods, we do not use the passage-query similarity feature PsgQuerySim (see Section 4.2) in $agg_{g \in d}(\vec{v}_{(g,q)})$ since aggregating this feature value across the passages in the document amounts to the AvgPDSim and MaxPDSim features, respectively, which are already used in $\vec{v}_{(g,q)}$.

We used leave-one-out cross validation over queries for training and testing; i.e., each query was used once for test wherein all other queries were used for training. For the LTR methods we randomly split the train set to train (80%) and validation (20%);⁷ the latter was used to set the hyper parameters of the LTR methods. For consistency, we use the same train set to set the free-parameter values of the non-LTR baselines (i.e., the validation set is not used for these methods). MAP and MAiP served as the optimization criteria for values of (hyper-) parameters in document and passage retrieval, respectively. We min-max normalized the feature values used in the learning-to-rank methods on a per-query basis.

The Dirichlet smoothing parameter was set to 1000 [59] for the initial language-model-based document retrieval, and to values in {500, 1500, 2500} in all other cases. The three parameters of MRF’s SDM are set to values in {0, 0.1, . . . , 1}. The value of λ in QSF is in {0.1, 0.2, . . . , 0.9}. RankSVM’s regularization parameter is set to {0.0001, 0.01, 0.1}; all other hyper parameters of RankSVM, and those of LambdaMART, are set to default values of the implementations.

For PLM, the value of the steepness parameter of the Gaussian kernel is in {50, 100, . . . , 300}; λ and β were set to values in {0, 0.2, . . . , 1} [33]. α (in the RRF and FPD methods from Section 3) and ν (in the RRF, SMPD and FPD methods from Section 3) are in {0, 0.1, . . . , 1} and {0, 30, 60, 90, 100}, respectively.

5 EXPERIMENTAL RESULTS

In Section 5.1 we analyze the performance of our passage-based document retrieval methods described in Section 3. As these methods rely on passage ranking, in Section 5.2 we analyze the performance of our learning-to-rank-based passage retrieval method.

5.1 Passage-Based Document Retrieval

5.1.1 Main Result. Table 2 presents our main result. We see that in all relevant comparisons (5 datasets \times 2 evaluation measures), JPDs, which is shown below to be our best performing approach, substantially outperforms all baselines: LM (unigram language-model-based retrieval), DocPsg (a representative passage-based document retrieval approach), SDM (a state-of-the-art term proximity method) and init-LTR (a learning-to-rank approach that utilizes document-query features). Most improvements are statistically significant. (We applied Bonferroni correction for multiple comparisons.) Refer back to Section 4.1 for more details about the baselines.

⁷The only exception was that the passage LTR method applied on TREC corpora was learned using all queries in the INEX dataset.

Recall that JPDs learns a document ranker by utilizing the document-query features used to induce init-LTR and the passage-query features of the document’s passage most highly ranked in response to the query. Its clear superiority with respect to the init-LTR methods attest to the merits of the way JPDs leverages passage-based information.

Given the performance superiority in most relevant comparisons of init-SVM and init-LMart to the other baselines, below we use them as reference comparisons. We note that their effectiveness attests to the effectiveness of the document features we use⁸ (See Section 4.1 for details regarding the features.)

Since our methods utilize init-SVM and init-LMart (i.e., the initial list D_{LTR} or features used to induce it), and using each of the two entails a different experimental setting, we compare X-SVM and X-LMart methods separately.

5.1.2 Comparing All Our Methods. Table 3 presents the performance comparison of all our proposed passage-based document retrieval methods from Section 3. The init-LTR methods serve for reference comparison.

We see in Table 3 that all the proposed methods outperform the init-LTR baselines — often statistically significantly — in the vast majority of relevant comparisons and are never outperformed in a statistically significant manner by a baseline.

JPDs is the most effective approach among those we proposed: its block in the table has the highest number of boldfaced numbers, it outperforms any other approach in most relevance comparisons, and it is never statistically significantly outperformed by other approaches while the reverse often holds. These findings attest to the merits of using the passage-query features of the document’s passage most highly ranked together with the document-query features to learn a document ranker.

The JPDm-max approach is the second-best performing. This finding is not entirely surprising: JPDs, which is our best performing method, uses the features of the document’s passage most highly ranked while JPDm-max uses per each passage-based feature the maximum value over the document’s passages. As could be expected, both JPDm-max and JPDm-avg outperform JPDm-min. That is, using the average or the maximum of a feature value across the document’s passages yields better performance than using the minimal value.

Table 3 also shows that RRF outperforms SMPD in most relevant comparisons when using SVM and the reverse holds when using LMart. However, only the MAP differences between RRF-SVM and SMPD-SVM for ROBUST and INEX are statistically significant. We thus conclude that the most important passage-rank-based information is the rank of a document’s most highly ranked passage. (Recall that SMPD uses additional statistics about the ranking of passages of a document.) We attribute these findings to the fact that a document can be deemed relevant even if it contains only a single short relevant passage.

Another observation that we make based on Table 3 is that FPD and JPDs outperform RRF in most relevant comparisons; i.e., using the query-passage features of the passage most highly ranked

⁸The finding that init-LMart underperforms init-SVM can be attributed to the fact that LMart is a non-linear ranker while SVM is, and the number of queries used for training is not very large.

Table 2: Main result. Comparison between document retrieval baselines and JPDs-LTR which is shown below to be our best performing method. 'l', 'd', 's' and 'i' mark statistically significant differences with LM, DocPsg, SDM and init-LTR respectively. Comparisons between LTR-based methods are performed between two methods utilizing the same LTR approach. Boldface: best result per column.

	ROBUST		WT10G		GOV2		ClueWeb		INEX	
	MAP	p@10								
LM	.254	.433	.195	.290	.292	.534	.187	.339	.367	.554
DocPsg	.254	.424	.209	.292	.298	.523	.168	.306	.368	.538
SDM	.261	.440	.202	.293	.304	.576	.192	.338	.385	.568
init-SVM	.261	.439	.213	.334	.336	.643	.222	.406	.392	.577
init-LMart	.245	.427	.198	.311	.326	.651	.224	.394	.378	.584
JPDs-SVM	.290^{ld}	.480^{ld}	.235^{ld}	.381^{ld}	.350^{ld}	.656^{ld}	.246^{ld}	.452^{ld}	.417^{ld}	.589^{ld}
JPDs-LMart	.290^{ld}	.471^{ld}	.229^{ld}	.378^{ld}	.345^{ld}	.655^{ld}	.234^{ld}	.423^{ld}	.410^{ld}	.593^{ld}

Table 3: Comparison of all our passage-based document retrieval methods. 'i' and 'j' mark statistically significant differences with init-LTR and JPDs-LTR, respectively. Comparisons between LTR-based methods are performed between two methods utilizing the same LTR approach. Boldface: best result per column.

	ROBUST		WT10G		GOV2		ClueWeb		INEX	
	MAP	p@10	MAP	p@10	MAP	p@10	MAP	p@10	MAP	p@10
init-SVM	.261	.439	.213	.334	.336	.643	.222	.406	.392	.577
init-LMart	.245	.427	.198	.311	.326	.651	.224	.394	.378	.584
JPDs-SVM	.290	.480	.235	.381	.350	.656	.246	.452	.417	.589
JPDs-LMart	.290	.471	.229	.378	.345	.655	.234	.423	.412	.593
RRF-SVM	.275 ^{ij}	.462 ^{ij}	.231 ⁱ	.376 ⁱ	.346 ⁱ	.639	.234 ^{ij}	.425 ^{ij}	.408 ^{ij}	.601 ⁱ
RRF-LMart	.281 ^{ij}	.462 ⁱ	.230 ⁱ	.367 ⁱ	.339 ^{ij}	.638	.232 ⁱ	.427 ⁱ	.410 ⁱ	.603
SMPD-SVM	.271 ^{ij}	.455 ^{ij}	.223 ^{ij}	.363 ⁱ	.344 ^{ij}	.647	.233 ^{ij}	.418 ^j	.401 ^{ij}	.598 ⁱ
SMPD-LMart	.280 ^{ij}	.460 ⁱ	.236 ⁱ	.370 ⁱ	.341 ⁱ	.641	.239 ⁱ	.433 ⁱ	.412 ⁱ	.600
JPDm-avg-SVM	.285 ^{ij}	.465 ^{ij}	.228 ⁱ	.363 ⁱ	.343 ^j	.639	.244 ⁱ	.434 ^{ij}	.415 ⁱ	.598 ⁱ
JPDm-avg-LMart	.288 ⁱ	.471 ⁱ	.223 ⁱ	.355 ^{ij}	.342 ⁱ	.663	.237 ⁱ	.422 ⁱ	.417 ⁱ	.595
JPDm-max-SVM	.293 ⁱ	.476 ⁱ	.235 ⁱ	.374 ⁱ	.350 ⁱ	.643	.242 ⁱ	.429 ^j	.420 ⁱ	.601 ⁱ
JPDm-max-LMart	.289 ⁱ	.468 ⁱ	.228 ⁱ	.363 ⁱ	.349 ⁱ	.654	.230	.416	.416 ⁱ	.602
JPDm-min-SVM	.270 ^{ij}	.451 ^j	.233 ⁱ	.342 ^j	.334 ^j	.630 ^j	.236 ⁱ	.430 ^{ij}	.404 ^{ij}	.583
JPDm-min-LMart	.271 ^{ij}	.454 ^{ij}	.220 ⁱ	.338 ^j	.337 ^{ij}	.640	.230	.403 ^j	.394 ^{ij}	.578
FPD-SVM	.288 ^{ij}	.474 ⁱ	.228 ^{ij}	.372 ⁱ	.348 ⁱ	.643	.238 ^{ij}	.434 ^{ij}	.411 ^{ij}	.588
FPD-LMart	.291 ⁱ	.468 ⁱ	.228 ⁱ	.362 ⁱ	.349 ⁱ	.655	.236 ⁱ	.423 ⁱ	.414 ⁱ	.609ⁱ

of a document is more effective than using its rank. Using these features together with document features (JPDs) is more effective than using them separately (FPD) to induce document ranking.

5.1.3 Further Analysis of JPDs. We saw above that JPDs is the most effective passage-based document retrieval approach among those we proposed. JPDs uses together the document-query features and the passage-query features of the document's most highly ranked passage so as to learn a document ranking function. In Table 4 we contrast the performance of JPDs with that of its variants that use the passage-query features of the document's second (JPDs-second), third (JPDs-third) and lowest (JPDs-lowest) ranked passages in $G(D_{LTR})$.

Table 4 shows that the original version, JPDs, outperforms in most relevant comparisons its variants (JPDs-second, JPDs-third and JPDs-lowest). More generally, we see that for almost all datasets, the lower the document's passage, whose passage-query features are used, is ranked, the lower the retrieval performance of the JPDs approach that uses these features.⁹ These findings attest to the merits of using the features of the document's most highly ranked passage, and to the fact that the relative ranking of the document's passages with respect to the query can imply to the benefits of using information induced from them to rank the document.

⁹We note that the use of the lowest ranked passage did not result in substantial performance decrease due to the length of passages used here: 300; that is, such passages can incorporate a descent amount of information from the entire document, especially in cases of relatively short documents.

Table 4: Comparing variants of JPDs. Boldface: the best result in a column for each LTR method (SVM or LMart). 'j' marks statistically significant differences with JPDs-LTR.

	ROBUST		WT10G		GOV2		ClueWeb		INEX	
	MAP	p@10	MAP	p@10	MAP	p@10	MAP	p@10	MAP	p@10
JPDs-SVM	.290	.480	.235	.381	.350	.656	.246	.452	.417	.589
JPDs-second-SVM	.277 ^j	.464 ^j	.236	.363	.341 ^j	.646	.238 ^j	.430 ^j	.414	.594
JPDs-third-SVM	.273 ^j	.455 ^j	.232	.363	.338 ^j	.633 ^j	.238 ^j	.434 ^j	.412	.598
JPDs-lowest-SVM	.271 ^j	.452 ^j	.231	.347 ^j	.335 ^j	.629 ^j	.226 ^j	.410 ^j	.402 ^j	.577
JPDs-LMart	.290	.471	.229	.378	.345	.655	.234	.423	.412	.593
JPDs-second-LMart	.280 ^j	.458	.226	.358	.341	.655	.240	.429	.410	.588
JPDs-third-LMart	.273 ^j	.455 ^j	.218	.361	.341	.650	.235	.422	.401 ^j	.587
JPDs-lowest-LMart	.270 ^j	.448 ^j	.219	.336 ^j	.337 ^j	.649	.232	.411	.400 ^j	.581

Table 5: Comparing JPDs with JPD-2 where the features of the document's two most highly ranked passages are used in addition to those of the document. Boldface: the best result in a column for each LTR method (SVM or LMart). 'j' marks statistically significant differences with JPDs-LTR.

	ROBUST		WT10G		GOV2		ClueWeb		INEX	
	MAP	p@10	MAP	p@10	MAP	p@10	MAP	p@10	MAP	p@10
JPDs-SVM	.290	.480	.235	.381	.350	.656	.246	.452	.417	.589
JPD-2-SVM	.291	.473	.235	.373	.351	.655	.250^j	.452	.421	.601
JPDs-LMart	.290	.471	.229	.378	.345	.655	.234	.423	.412	.593
JPD-2-LMart	.291	.473	.236	.375	.349	.649	.235	.430	.418	.605

Table 6: The effect on document ranking effectiveness of the passage ranker: LTR-based (PsgLTR) vs. integrating the passage-query similarity with the query-similarity of the passage's ambient document (QSF). '**' marks statistically significant differences between PsgLTR and QSF. Boldface: the best result for evaluation measure in a block.

	ROBUST		WT10G		GOV2		ClueWeb		INEX	
	MAP	p@10								
RRF-SVM PsgLTR	.275*	.462*	.231*	.376*	.346*	.639	.234*	.425*	.408*	.601*
RRF-SVM QSF	.261	.442	.215	.324	.336	.643	.223	.406	.390	.574
RRF-LMart PsgLTR	.281*	.462*	.230*	.367*	.339*	.638	.232*	.427*	.410*	.603*
RRF-LMart QSF	.257	.442	.204	.318	.326	.645	.224	.396	.382	.581
SMPD-SVM PsgLTR	.271*	.455*	.223*	.363*	.344	.647	.233	.418	.401*	.598*
SMPD-SVM QSF	.259	.439	.213	.337	.337	.642	.227	.409	.386	.564
SMPD-LMart PsgLTR	.280*	.460*	.236*	.370*	.341	.641	.239*	.433*	.412*	.600
SMPD-LMart QSF	.258	.442	.211	.327	.336	.651	.223	.407	.389	.579
JPDs-SVM PsgLTR	.290	.480	.235	.381	.350	.656	.246	.452	.417	.589
JPDs-SVM QSF	.288	.474	.233	.373	.347	.647	.245	.441	.414	.595
JPDs-LMart PsgLTR	.290	.471	.229	.378	.345	.655	.234	.423	.412	.593
JPDs-LMart QSF	.289	.473	.230	.365	.343	.641	.228	.407	.410	.597
FPD-SVM PsgLTR	.288	.474	.228	.372	.348*	.643	.238*	.434	.411*	.588
FPD-SVM QSF	.286	.475	.230	.365	.345	.632	.230	.422	.405	.591
FPD-LMart PsgLTR	.291	.468	.228	.362	.349*	.655*	.236	.423	.414	.609
FPD-LMart QSF	.287	.468	.225	.361	.344	.631	.233	.417	.411	.605

5.1.4 Utilizing Two Passages. Our JPDs method utilizes the features of the document's most highly ranked passage in addition

to the document's features. We now consider a variant of JPDs,

denoted **JPD-2**, which uses in addition the features of the document’s passage which is the second ranked¹⁰. The feature vectors of the two passages are concatenated with that of the document for learning a document ranker. Table 5 presents the results.

We see in Table 5 that using the two passages (JPD-2-LTR) yields performance that is very similar in most relevant comparisons to that of using a single passage (JPDs-LTR). In only a single case, the performance difference is statistically significant.

5.1.5 The Effect of the Passage Ranker. Our passage-based document retrieval approaches (except for JPDm) utilize information induced from the ranking of passages in the initially retrieved document list, D_{init} . In Table 6 we compare the performance of the approaches when using two different passage ranking methods. The first is the QSF method described in Section 4.2.1 which integrates the passage-query similarity value with the query-similarity value of the passage’s ambient document. The second passage ranking method, **PsgLTR**, was used insofar: SVM or LMart applied with our proposed passage-based features from Section 4.2¹¹. In Section 5.2 we show that the passage-ranking effectiveness of PsgLTR is substantially better than that of QSF.

The message rising from Table 6 is clear: our passage-based document retrieval methods post better performance when using the LTR-based passage ranker than when using the QSF method to rank passages. While most improvements are statistically significant, those for JPDs are not. This finding attests to the robustness of JPDs with respect to the passage ranker used.

5.1.6 Feature Analysis for Document Retrieval. We now present feature analysis for our best performing approach, JPDs. We start by analyzing JPDs-SVM which outperforms JPDs-LMart (see Table 2).

First, we average, per dataset, the weights assigned to features in JPDs-SVM using the different training folds. (Recall that we use leave-one-out cross validation.) Then, the features are ordered in descending order of these averages. Each feature is assigned a score which is the reciprocal of its rank position in the ordered list. Finally, features are ordered by averaging their scores across datasets. The top 10 features¹² according to this analysis are (p and d indicate that the feature is of the passage or the document, respectively): SDM unigrams (d), ESA (p), Entity (p), Ent (d), AvgPDSim (p), MaxPDSim (p), SW2 (d), SDM biterns (d), SynonymsOverlap (p), W2V (p). Thus, both document-based and passage-based features are among the top-5 and top-10. This finding attests to the merits of using both types of features to learn a document ranking function.

We also performed ablation tests for JPDs where we removed one feature at a time. Actual numbers are omitted as they convey no additional insight. We order the features in descending order of the number of cases where their removal resulted in statistically significant performance drop. A case is defined by a dataset and

evaluation measure. (We include JPDs-SVM and JPDs-LMart together in this analysis.) We mark the features with (d/p,x): whether the feature is document-based or passage-based (d/p) and the number of cases (x) its removal caused statistically significant performance drop. The ordered list of features is: ESA (p,15), SDM unigrams (d,4), SDM biterns (d,2), SW1 (d,2), Ent (d,1), SW2 (d,1), SDM bigrams (d,1), MaxPDSim (p,1), LengthRatio (p,1), SynonymsOverlap (p,1), pLocation (p,1), Entity (p,1). Thus, as was the case for the SVM-based feature weight analysis from above, ESA which is a passage feature and SDM unigrams which is a document feature are the most important. More generally, the list contains both document and passage features. We note that while the removal of each of the document features resulted in at least one case of statistically significant drop, for quite a few passage features this was not the case; i.e., there is redundancy between the passage features.

We next turn to present feature analysis for the SMPD approach¹³. SMPD uses the same document features as JPDs, but different passage-based features: mainly those which quantify the rank positions of the document’s passages in the passage ranking. The results of an ablation test, as that performed above, are: max (p,5), SW2 (d,4), SDM unigrams (d,3), SDM biterns (d,2), avg (p,2), numPsg (p,2), Ent (d,1), SW1 (d,1), SDM bigrams (d,1), min (p,1), std (p,1), top50 (p,1). We observe again a mix of document and passage features. The max feature, which quantifies the rank position of the document’s most highly ranked passage, is more important than the min and avg features. This finding provides further support to the merits of using information about the highest ranked passage of the document.

5.1.7 LTR Methods. Heretofore, we applied our methods using two LTR approaches: RankSVM and LambdaMART. In Table 7, we study the performance of our JPDs method with two additional LTR approaches: MART [18] and coordinate ascent [38]. MART, known as gradient boosted regression trees, is a non-linear pairwise ranker which combines the outputs obtained by different regression trees. On the other hand, coordinate ascent (CAscent in short) is a linear listwise approach. We used the RankLib implementations of the MART and CAscent algorithms¹⁴. CAscent was trained for NDCG@10.

Table 7 shows that the JPDs method improves over the initial LTR ranking in all relevant comparisons (5 datasets \times 2 evaluation measures \times 4 LTR methods). Most of the improvements for SVM and LMart are statistically significant while some of the improvements for MART and CAscent are statistically significant.

We also see in Table 7 that in most relevant comparisons, using JPDs with SVM and LMart results in performance that transcends that of its implementations that use MART and CAscent. This finding can be attributed to some extent to the effectiveness of the passage ranking utilized by JPDs. The MAiP effectiveness of the passage ranking induced using MART and CAscent is lower than that attained by using SVM and LMart when using the INEX dataset for passage retrieval evaluation. Specifically, the MAiP performance of SVM, LMart, MART and CAscent is .267, .275, .250 and .259, respectively.

¹⁰To avoid having the same features used for the two passages, the following features were removed from the feature vector of the second ranked passage: DocQuerySim, MaxPDSim, AvgPDSim, StdPDSim and QueryLength.

¹¹We do not present the comparison for the JPDm approach as it is independent of the passage ranking.

¹²JPDs-SVM uses 24 features and JPDs-LMart uses 25 features — the additional feature is the query length which is not useful for a linear ranker.

¹³In this analysis we set v , the free parameter of SMPD, to a value which is effective across the train folds.

¹⁴<https://sourceforge.net/p/lemur/wiki/RankLib/>.

Table 7: Varying the LTR method used in JPDs and in init-LTR. 'i' marks statistically significant difference with init-LTR. Boldface: the best result in a column for each LTR method (SVM, LMart, MART or CAscent)

	ROBUST		WT10G		GOV2		ClueWeb		INEX	
	MAP	p@10	MAP	p@10	MAP	p@10	MAP	p@10	MAP	p@10
init-SVM	.261	.439	.213	.334	.336	.643	.222	.406	.392	.577
JPDs-SVM	.290ⁱ	.480ⁱ	.235ⁱ	.381ⁱ	.350ⁱ	.656	.246ⁱ	.452ⁱ	.417ⁱ	.589
init-LMart	.245	.427	.198	.311	.326	.651	.224	.394	.378	.584
JPDs-LMart	.290ⁱ	.471ⁱ	.229ⁱ	.378ⁱ	.345ⁱ	.655	.234ⁱ	.423ⁱ	.412ⁱ	.593
init-MART	.258	.439	.203	.305	.332	.640	.216	.403	.381	.565
JPDs-MART	.285ⁱ	.462ⁱ	.211	.345ⁱ	.343ⁱ	.659	.223	.415	.407ⁱ	.577
init-CAscent	.257	.443	.211	.324	.329	.649	.212	.406	.377	.586
JPDs-CAscent	.273ⁱ	.471ⁱ	.226ⁱ	.372ⁱ	.339ⁱ	.647	.215	.420	.382	.602

Table 8: Using the MSLR (LETOR) document features in comparison to using the features used thusfar for the initial document ranking and in our JPDs method. 'i' and 'm' mark statistically significant differences with init-LTR and init-MSLR-LTR, respectively. Boldface: the best result in a column, per block of either the original features (first block) or the MSLR features (second block), for each LTR method (SVM or LMart).

	GOV2		ClueWeb	
	MAP	p@10	MAP	p@10
init-SVM	.336	.643	.222	.406
init-LMart	.326	.651	.224	.394
JPDs-SVM	.350ⁱ	.656	.246ⁱ	.452ⁱ
JPDs-LMart	.345ⁱ	.655	.234ⁱ	.423ⁱ
init-MSLR-SVM	.323	.595	.251	.437
init-MSLR-LMart	.315	.599	.241	.428
JPDs-MSLR-SVM	.353^m	.634^m	.264^m	.452^m
JPDs-MSLR-LMart	.342^m	.633^m	.244	.437

5.1.8 *Using LETOR Features.* Insofar, we have used the document features described in Section 4.1. This practice resulted in highly effective document ranking performance as exhibited by the init-LTR baselines as well as our methods. We now turn to explore the performance of our methods with a much larger set of document(-query) features. Specifically, we use the MSLR¹⁵ features from the LETOR datasets for retrieval over the GOV2 and ClueWeb collections with the queries specified in Table 1. We used all MSLR features except for the Outlink number, SiteRank, QualityScore, QualityScore2, Query-url click count, url click count, and url dwell time. In addition to the MSLR features, we also use here the highly effective query-independent document quality measures used above: the fraction of terms in the document that are stopwords, the fraction of stopwords that appear in the document, and the entropy of the term distribution in the document. The stopword list used for the two stopword features is composed of the collection’s 100 most frequent alphanumeric terms [45, 47]. For ClueWeb we also used the spam score assigned to a document by the Waterloo spam classifier and the PageRank score. All together, we used, at the document level, 149 features for GOV2 and 151 features for ClueWeb.

The results are presented in Table 8. We first see that in terms of the initial ranking, the MSLR features are more effective than those we used above for ClueWeb, but the reverse holds for GOV2. (This could potentially be attributed to the fact that for GOV2 there are fewer queries than for ClueWeb.) We further see in Table 8 that our JPDs method is also effective with the MSLR features. It always outperforms the initial ranking; in most relevant comparisons, the improvements are statistically significant.

5.2 Passage Retrieval

Heretofore, we have focused on the document retrieval task. Our passage-based document retrieval methods utilize a ranking of passages induced using our proposed passage retrieval approach. (See Section 4.2 for details.) We now turn to compare the performance of our passage ranker with that of the passage retrieval baselines described in Section 4.2.1.

Table 9 presents the performance numbers of the passage retrieval methods for the INEX collection. We see that our LTR methods, PsgLTR-SVM and PsgLTR-LMart, outperform all other passage retrieval methods in most relevant comparisons (3 passage lengths \times 3 evaluation measures) with many of the improvements being statistically significant. We note that the MKS baseline [56]

¹⁵www.research.microsoft.com/en-us/projects/mslr

Table 9: Passage retrieval over INEX with passages of length 300, 150 and 50. LM is standard language-model-based document retrieval (i.e., documents serve for passages). Boldface: the best result in a column. Statistically significant differences with LM, QSF and PLM are marked with 'l', 'f' and 'm', respectively. 'o' and 'k' mark a statistically significant difference between PsgLTR-X and owpc-X and between PsgLTR-X and MKS-X, respectively.

	INEX								
	Psg300			Psg150			Psg50		
	MAiP	iP[.01]	iP[.1]	MAiP	iP[.01]	iP[.1]	MAiP	iP[.01]	iP[.1]
LM	.256	.523	.449	.256	.523	.449	.256	.523	.449
QSF	.248	.577	.453	.234	.575	.455	.209	.581	.449
PLM	.253	.586	.472	.240	.596	.471	.215	.605	.469
owpc-SVM	.242	.577	.440	.229	.569	.438	.202	.570	.431
owpc-LMart	.255	.578	.460	.240	.566	.450	.208	.577	.443
MKS-SVM	.247	.593	.468	.235	.602	.459	.199	.626	.457
MKS-LMart	.262	.620	.479	.241	.629	.479	.200	.644	.459
PsgLTR-SVM	.267 _{ok}	.637 _o ^{lf}	.487 _o	.253_o	.662_{ok} ^{lfm}	.492 _o	.213 ^l	.647_o ^l	.467
PsgLTR-LMart	.275 _o ^{fm}	.644_o ^{fm}	.496	.253	.650 _o ^{fm}	.494_o ^l	.209 ^l	.634 ^l	.454

Table 10: Sentence retrieval over AQUAINT. Boldface: the best result in a column. Statistically significant differences with QSF and PLM are marked with 'f' and 'm', respectively. 'o' and 'k' mark a statistically significant difference between PsgLTR-X and owpc-X and between PsgLTR-X and MKS-X, respectively.

	AQUAINT	
	MAP	p@10
QSF	.471	.624
PLM	.518	.669
owpc-SVM	.579	.701
owpc-LMart	.589	.716
MKS-SVM	.569	.664
MKS-LMart	.585	.701
PsgLTR-SVM	.602 _{ok} ^{fm}	.713 _k ^f
PsgLTR-LMart	.606_{ok} ^{fm}	.710 ^f

was recently shown to yield state-of-the-art passage retrieval performance.

Table 10 presents the effectiveness of our passage retrieval approach, PsgLTR, in ranking sentences in the AQUAINT collection. We see that PsgLTR-SVM and PsgLTR-LMart statistically significantly outperform all other passage retrieval methods in terms of MAP. In the single case where our methods are outperformed by another method (MKS-LMart) in terms of p@10, the performance differences are not statistically significant.

The findings presented above for focused (passage) retrieval over INEX, and sentence retrieval over AQUAINT, attest to the fact that our passage ranker posts state-of-the-art passage retrieval performance.

5.2.1 Feature Analysis for Passage Retrieval. We first use the SVM-based feature analysis, as was performed above for document retrieval, to analyze the relative importance of features used in our passage retrieval approach (PsgLTR-SVM). For INEX, we consider each of the three passage lengths as a different experimental setting. The top 10 features for INEX are: ESA, SW1, MaxPDSim, Entity, StdPDSim, SW2, Ent, DocQuerySim, AvgPDSim and SynonymsOverlap. For AQUAINT, the top-10 features are: Ent, SW1, ESA, LengthRatio, TermOverlap, AvgPDSim, QuerySimPre, SynonymsOverlap, QuerySimFollow, PsgLength. Recall that using stopwords-based passage priors (SW1 and SW2) to rank passages is novel to this study. We see that SW1 is the second most important feature for both INEX and AQUAINT. Another observation is that, as expected, the relative ordering of passages in this analysis, and the set of features that are among the top-10, are not identical to those presented above when using the passage features for document retrieval.

In addition, we perform ablation tests for PsgLTR. When using passages of 300 terms for INEX, the features whose removal resulted in statistically significant performance drop of MAiP are: ESA, MaxPDSim, AvgPDSim, SW1. The features whose removal resulted in statistically significant performance drop of MAP for AQUAINT are: Ent, SW1, ESA, SW2. The features are ordered in both cases in a descending order of the performance drop. Given that the retrieval tasks over INEX (passage retrieval) and AQUAINT (sentence retrieval) are different, it is not surprising that the feature lists are a bit different. Yet, ESA and SW1 are in both cases among the most important features, which was also the case above in the SVM-based analysis.

6 CONCLUSIONS AND FUTURE WORK

Our focus in this work was on passage-based document retrieval: document ranking methods that utilize information induced from document passages. Previous work on passage-based document retrieval has focused on methods that integrate passage-query and

document-query similarity values. Here, we addressed the challenge of utilizing richer sources of passage-based information for improving document retrieval effectiveness.

We presented a suite of learning-to-rank methods for document retrieval that use passage-based information. Most of the methods rely on ranking passages in response to the query using an effective approach, specifically, utilizing learning-to-rank. Some of the methods use information about the ranking of the passages of a document. Other methods use the passage-based features utilized for passage ranking and integrate them with document-based features so as to learn a document ranking function. We described connections between our methods and past unsupervised approaches for passage-based document retrieval as well as approaches for ranking clusters of similar documents.

To learn a passage-ranking method, we used previously proposed features along with features which were not used before for learning passage ranking functions. These features are query-independent passage-relevance priors adopted from work on using document relevance priors for Web search.

Empirical evaluation performed with a suite of datasets demonstrated the effectiveness of our methods. Our most effective method integrates document-based features with passage-based features of the document's most highly ranked passage. In addition, our best performing method was shown to outperform the use of different sets of document-based features. Further exploration provided support to the merits of using an effective passage ranking method. We also showed that our passage-ranking method yields state-of-the-art passage retrieval performance.

For future work we intend to integrate in our methods additional passage-based features; e.g., those induced from inter-passage similarities. We also plan to explore how our methods can be used for, and with, pseudo-feedback-query expansion. A case in point, we can apply query expansion at the passage-level, document-level, or both, so as to enrich the feature set used.

REFERENCES

- [1] Paavo Arvola, Shlomo Geva, Jaap Kamps, Ralf Schenkel, Andrew Trotman, and Johanna Vainio. 2011. Overview of the INEX 2010 ad hoc track. In *Comparative Evaluation of Focused Retrieval*. 1–32.
- [2] Michael Bendersky, W Bruce Croft, and Yanlei Diao. 2011. Quality-biased ranking of web documents. In *Proc. of WSDM*. 95–104.
- [3] Michael Bendersky and Oren Kurland. 2008. Re-ranking search results using document-passage graphs. In *Proc. of SIGIR*. 853–854.
- [4] Michael Bendersky and Oren Kurland. 2010. Utilizing passage-based language models for ad hoc document retrieval. *Information Retrieval* 13, 2 (2010), 157–187.
- [5] David Buffoni, Nicolas Usunier, and Patrick Gallinari. 2010. Lip6 at INEX: OWPC for ad hoc track. In *Focused Retrieval and Evaluation*. 59–69.
- [6] Christopher J.C. Burges. 2010. *From RankNet to LambdaRank to LambdaMART: An Overview*. Technical Report. Microsoft Research.
- [7] James P Callan. 1994. Passage-level evidence in document retrieval. In *Proc. of SIGIR*. 302–310.
- [8] David Carmel, Anna Shtok, and Oren Kurland. 2013. Position-based contextualization for passage retrieval. In *Proc. of CIKM*. 1241–1244.
- [9] Ruy-Cheng Chen, Damiano Spina, W. Bruce Croft, Mark Sanderson, and Falk Scholer. 2015. Harnessing Semantics for Answer Sentence Retrieval. In *Proc. of ESAIR*. 21–27.
- [10] Ruy-Cheng Chen, Evi Yulianti, Mark Sanderson, and W Bruce Cro. 2017. On the Benefit of Incorporating External Features in a Neural Architecture for Answer Sentence Selection. In *Proc. of SIGIR*. 1017–1020.
- [11] Gordon V Cormack, Charles LA Clarke, and Stefan Buettcher. 2009. Reciprocal rank fusion outperforms condorcet and individual rank learning methods. In *Proc. of SIGIR*. 758–759.
- [12] Gordon V Cormack, Mark D Smucker, and Charles LA Clarke. 2011. Efficient and effective spam filtering and re-ranking for large web datasets. *Information retrieval* 14, 5 (2011), 441–465.
- [13] Ludovic Denoyer, Hugo Zaragoza, and Patrick Gallinari. 2001. HMM-based Passage Models for Document Classification and Ranking. In *Proc. of ECIR*.
- [14] Yixing Fan, Jiafeng Guo, Yanyan Lan, Jun Xu, Chengxiang Zhai, and Xueqi Cheng. 2018. Modeling Diverse Relevance Patterns in Ad-hoc Retrieval. In *Proceedings of SIGIR*. 375–384.
- [15] Ronald T. Fernández and David E. Losada. 2012. Effective sentence retrieval based on query-independent evidence. *Information Processing and Management* 48, 6 (2012), 1203–1229.
- [16] Ronald T Fernández, David E Losada, and Leif A Azzopardi. 2011. Extending the language modeling framework for sentence retrieval to include local context. *Information Retrieval* 14, 4 (2011), 355–389.
- [17] Paolo Ferragina and Ugo Scaiella. 2012. Fast and accurate annotation of short texts with wikipedia pages. *IEEE software* 29, 1 (2012), 70–75.
- [18] Jerome H Friedman. 2001. Greedy function approximation: a gradient boosting machine. *Annals of statistics* (2001), 1189–1232.
- [19] Evgeniy Gabrilovich and Shaul Markovitch. 2007. Computing semantic relatedness using Wikipedia-based explicit semantic analysis. In *Proc. of IJCAI*, Vol. 7. 1606–1611.
- [20] Shlomo Geva, Jaap Kamps, Miro Lethonen, Ralf Schenkel, James A Thom, and Andrew Trotman. 2010. Overview of the INEX 2009 ad hoc track. In *Focused retrieval and evaluation*. 4–25.
- [21] Marti A. Hearst and Christian Plaunt. 1993. Subtopic Structuring for Full-Length Document Access. In *Proc. of SIGIR*. 59–68.
- [22] Jing Jiang and Chengxiang Zhai. 2004. UIUC in HARD 2004–Passage Retrieval Using HMMs. In *TREC*.
- [23] Thorsten Joachims. 2006. Training linear SVMs in linear time. In *Proc. of KDD*. 217–226.
- [24] Marcin Kaszkiel and Justin Zobel. 1997. Passage Retrieval Revisited. In *Proc. of SIGIR*. 178–185.
- [25] Marcin Kaszkiel and Justin Zobel. 2001. Effective ranking with arbitrary passages. *Journal of the American Society for Information Science and Technology* 52, 4 (2001), 344–364.
- [26] Mostafa Keikha, Jae Hyun Park, and W. Bruce Croft. 2014. Evaluating answer passages using summarization measures. In *Proc. of SIGIR*. 963–966.
- [27] Mostafa Keikha, Jae Hyun Park, W. Bruce Croft, and Mark Sanderson. 2014. Retrieving Passages and Finding Answers. In *Proc. of ADCS*. 81.
- [28] Eyal Krikon, Oren Kurland, and Michael Bendersky. 2010. Utilizing inter-passage and inter-document similarities for reranking search results. *ACM Trans. Inf. Syst.* 29, 1 (2010), 3:1–3:28.
- [29] Oren Kurland and Carmel Domshlak. 2008. A rank-aggregation approach to searching for optimal query-specific clusters. In *Proceedings of SIGIR*. 547–554.
- [30] Hao Lang, Donald Metzler, Bin Wang, and Jin-Tao Li. 2010. Improved latent concept expansion using hierarchical markov random fields. In *Proc. of CIKM*. 249–258.
- [31] Tie-Yan Liu. 2009. Learning to rank for information retrieval. *Foundations and Trends in Information Retrieval* 3, 3 (2009), 225–331.
- [32] Xiaoyong Liu and W. Bruce Croft. 2002. Passage retrieval based on language models. In *Proc. of CIKM*. 375–382.
- [33] Yuanhua Lv and Chengxiang Zhai. 2009. Positional language models for information retrieval. In *Proc. of SIGIR*. 299–306.
- [34] Yuanhua Lv and Chengxiang Zhai. 2010. Positional relevance model for pseudo-relevance feedback. In *Proc. of SIGIR*. 579–586.
- [35] Craig Macdonald, Rodrygo LT Santos, and Iadh Ounis. 2012. On the usefulness of query features for learning to rank. In *Proc. of CIKM*. 2559–2562.
- [36] Donald Metzler and W. Bruce Croft. 2005. A Markov Random Field Model for Term Dependencies. In *Proc. of SIGIR*. 472–479.
- [37] Donald Metzler and W. Bruce Croft. 2007. Latent concept expansion using markov random fields. In *Proc. of SIGIR*. 311–318.
- [38] Donald Metzler and W Bruce Croft. 2007. Linear feature-based models for information retrieval. *Information Retrieval* 10, 3 (2007), 257–274.
- [39] Donald Metzler and Tapas Kanungo. 2008. Machine learned sentence selection strategies for query-biased summarization. In *Proc. of SIGIR*. 40–47.
- [40] Jun Miao, Jimmy Xiangji Huang, and Zheng Ye. 2012. Proximity-based rochio's model for pseudo relevance. In *Proc. of SIGIR*. 535–544.
- [41] Elke Mittendorf and Peter Schäuble. 1994. Document and passage retrieval based on hidden Markov models. In *Proc. of SIGIR*. Springer-Verlag New York, Inc., 318–327.
- [42] Vanessa Murdock and W Bruce Croft. 2005. A translation model for sentence retrieval. In *Proc. of HLT/EMNLP*. Association for Computational Linguistics, 684–691.
- [43] Vanessa Graham Murdock. 2006. *Aspects of sentence retrieval*. Ph.D. Dissertation. University of Massachusetts Amherst.
- [44] Seung-Hoon Na, In-Su Kang, Yeha Lee, and Jong-Hyeok Lee. 2008. Completely-Arbitrary Passage Retrieval in Language Modeling Approach. In *Proc. of AIRS*. 22–33.

- [45] Alexandros Ntoulas, Marc Najork, Mark Manasse, and Dennis Fetterly. 2006. Detecting spam web pages through content analysis. In *Proc. of WWW*. 83–92.
- [46] Fiana Raiber and Oren Kurland. 2013. Ranking document clusters using markov random fields. In *Proceedings of SIGIR*. 333–342.
- [47] Nimrod Raifer, Fiana Raiber, Moshe Tennenholtz, and Oren Kurland. 2017. Information Retrieval Meets Game Theory: The Ranking Competition Between Documents’ Authors. In *Proc. of SIGIR*. 465–474.
- [48] Gerard Salton, James Allan, and Chris Buckley. 1993. Approaches to passage retrieval in full text information systems. In *Proc. of SIGIR*. 49–58.
- [49] Ian Soboroff. 2004. Overview of the TREC 2004 Novelty Track. In *Proc. of TREC*.
- [50] Ian Soboroff and Donna Harman. 2003. Overview of the TREC 2003 Novelty Track. In *Proc. of TREC*. 38–53.
- [51] Tao Tao and ChengXiang Zhai. 2007. An exploration of proximity measures in information retrieval. In *Proc. of SIGIR*. 295–302.
- [52] Ellen M. Voorhees and Donna K. Harman. 2005. *TREC: Experiments and evaluation in information retrieval*. The MIT Press.
- [53] Xiaojun Wan, Jianwu Yang, and Jianguo Xiao. 2008. Towards a unified approach to document similarity search using manifold-ranking of blocks. *Information Processing and Management* 44, 3 (2008), 1032–1048.
- [54] Mengqiu Wang and Luo Si. 2008. Discriminative probabilistic models for passage based retrieval. In *Proc. of SIGIR*. 419–426.
- [55] Ross Wilkinson. 1994. Effective Retrieval of Structured Documents. In *Proc. of SIGIR*. 311–317.
- [56] Liu Yang, Qingyao Ai, Damiano Spina, Ruey-Cheng Chen, Liang Pang, W Bruce Croft, Jiafeng Guo, and Falk Scholer. 2016. Beyond factoid QA: Effective methods for non-factoid answer sentence retrieval. In *Proc. of ECIR*. Springer, 115–128.
- [57] Evi Yulianti, Ruey-Cheng Chen, Falk Scholer, W. Bruce Croft, and Mark Sanderson. 2018. Ranking Documents by Answer-Passage Quality. In *Proceedings of SIGIR*. 335–344.
- [58] Evi Yulianti, Ruey-Cheng Chen, Falk Scholer, and Mark Sanderson. 2016. Using Semantic and Context Features for Answer Summary Extraction. In *Proc. of ADCS*. 81–84.
- [59] Chengxiang Zhai and John Lafferty. 2001. A study of smoothing methods for language models applied to ad hoc information retrieval. In *Proc. of SIGIR*. 334–342.
- [60] Jinglei Zhao and Yeogirl Yun. 2009. A proximity language model for information retrieval. In *Proc. of SIGIR*. 291–298.