

Re-ranking search results using language models of query-specific clusters

Oren Kurland

Received: 13 October 2007 / Accepted: 16 June 2008 / Published online: 10 July 2008
© Springer Science+Business Media, LLC 2008

Abstract To obtain high precision at top ranks by a search performed in response to a query, researchers have proposed a cluster-based re-ranking paradigm: clustering an initial list of documents that are the most highly ranked by some initial search, and using information induced from these (often called) *query-specific* clusters for re-ranking the list. However, results concerning the effectiveness of various *automatic* cluster-based re-ranking methods have been inconclusive. We show that using query-specific clusters for automatic re-ranking of top-retrieved documents is effective with several methods in which clusters play different roles, among which is the *smoothing* of *document language models*. We do so by adapting previously-proposed cluster-based retrieval approaches, which are based on (static) query-independent clusters for ranking all documents in a corpus, to the re-ranking setting wherein clusters are query-specific. The best performing method that we develop outperforms both the initial document-based ranking and some previously proposed cluster-based re-ranking approaches; furthermore, this algorithm consistently outperforms a state-of-the-art pseudo-feedback-based approach. In further exploration we study the performance of cluster-based smoothing methods for re-ranking with various (soft and hard) clustering algorithms, and demonstrate the importance of clusters in providing context from the initial list through a comparison to using single documents to this end.

Keywords Query-specific clusters · Cluster-based language models · Cluster-based re-ranking · Cluster-based smoothing

1 Introduction

Users of search engines expect to see the documents most pertaining to their queries at the top ranks of the retrieved results (Croft 1995). A paradigm suggested by several

Initial version of this paper appeared as Chapter 5 of Kurland (2006).

O. Kurland (✉)
Faculty of Industrial Engineering and Management, Technion—Israel Institute of Technology,
Technion City, Haifa 32000, Israel
e-mail: kurland@ie.technion.ac.il

researchers for achieving this goal is to perform an initial search over the entire corpus in response to a query, and then to automatically *re-rank* the most highly ranked documents, so as to improve the precision at the very top ranks of the resultant list. (See, for example, Preece (1973); Willett (1985); Kleinberg (1998); Liu and Croft (2004); Diaz (2005); Kurland and Lee (2005, 2006); Liu and Croft (2006a).) The motivating idea behind the *re-ranking* paradigm is that the ratio of relevant to non-relevant documents in the *initial list* to be re-ranked, that is, the most highly ranked documents from the initial search, tends to be much higher than that in the entire corpus. However, since documents in the list were retrieved in response to a query, it is a challenging task to differentiate the relevant documents from the non-relevant ones.

To approach this challenge of (automatic) re-ranking, several researchers (Preece 1973; Willett 1985; Liu and Croft 2004; Kurland and Lee 2006; Liu and Croft 2006a, b) proposed to cluster the documents in the initial list and utilize information induced from the clusters; those are often termed *query-specific* clusters since the documents upon which clustering is performed were retrieved in response to a query. A potential advantage in using query-specific clusters for re-ranking that researchers (Hearst and Pedersen 1996; Tombros et al. 2002; Liu and Croft 2004; Kurland and Lee 2006) have pointed out is that relevant documents in the initial list might be clustered together—a manifestation of van Rijsbergen’s *cluster hypothesis* (van Rijsbergen 1979) in the re-ranking setting. Indeed, there is some empirical evidence that (under different clustering algorithms) there are often some query-specific clusters that contain a high percentage of relevant documents (Hearst and Pedersen 1996; Tombros et al. 2002; Kurland 2006). However, automatically finding these clusters is a very hard challenge (Willett 1985; Liu and Croft 2004). On the other hand, it was shown that users of interactive search systems can use query-specific clusters for quickly detecting relevant documents that they contain (Hearst and Pedersen 1996; Leuski 2001).

A different way by which clusters can be utilized has recently been proposed in the language modeling framework to information retrieval (Ponte and Croft 1998; Croft and Lafferty 2003). Researchers suggested to use information induced from document clusters to *smooth* document language models so as to “enrich” the document representation with corpus-related information (Azzopardi et al. 2004; Kurland and Lee 2004; Liu and Croft 2004; Tao et al. 2006; Wei and Croft 2006). Indeed, cluster-based smoothing has shown promise for ranking an entire corpus when *static query-independent clusters*, which are created offline, were used (Azzopardi et al. 2004; Kurland and Lee 2004; Liu and Croft 2004; Tao et al. 2006; Wei and Croft 2006). However, the results regarding the effectiveness of cluster-based smoothing for the re-ranking setting (using query-specific clusters) have remained inconclusive (Liu and Croft 2004).

We show that using query-specific clusters for automatic re-ranking is in fact effective whether clusters are used for *selecting* documents—specifically, detecting relevant documents by the patterns of membership of documents in clusters—or for *smoothing* document language models. We do so by adapting recently proposed cluster-based retrieval algorithms (Kurland and Lee 2004), which utilize information induced from static query-independent clusters for ranking all documents in a corpus, to the re-ranking setting wherein clusters are query-specific.

We empirically show that the most effective (cluster-based smoothing) re-ranking algorithm that we present not only significantly outperforms the initial document-based ranking over all tested TREC corpora, but also consistently outperforms a state-of-the-art pseudo-feedback-based approach, namely *the relevance model* (Lavrenko and Croft 2001). Moreover, the algorithm also outperforms some previously-proposed cluster-based approaches for re-ranking that utilize information induced from query-specific clusters. In further exploration we study the performance of cluster-based smoothing methods for re-ranking with

various clustering algorithms, and demonstrate the importance of clusters in providing context from the initial list through a comparison to using single documents to this end.

The rest of the paper is organized as follows. In Sect. 2 we present the different re-ranking algorithms that we explore. Section 3 describes the connection of our approach to previously-suggested models for re-ranking and to previous approaches for utilizing cluster-based information. We then present an empirical evaluation of the performance of our algorithms in Sect. 4 and conclude in Sect. 5.

2 Retrieval framework

Since we are focused on the re-ranking setting, the algorithms we present are applied not to the entire corpus \mathcal{C} , but to a subset $\mathcal{D}_{\text{init}}^{N,q}$ (henceforth $\mathcal{D}_{\text{init}}$), defined as the top N documents retrieved in response to the query q by a given initial retrieval engine. The algorithms also take into account a set $Cl(\mathcal{D}_{\text{init}})$ of (query-specific) clusters of the documents in $\mathcal{D}_{\text{init}}$. We assume that documents in $\mathcal{D}_{\text{init}}$ and clusters in $Cl(\mathcal{D}_{\text{init}})$ are assigned with unique IDs.

The algorithms we present utilize statistical language models (Ponte and Croft 1998; Croft and Lafferty 2003). We use $p_x(y)$ to denote the language-model-based similarity between x (a document or a cluster) and y (a document, a cluster, or a query).¹ Our language-model-induction methods are described in Sect. 2.2.

Clustering Previous work on utilizing query-specific clustering has focused on hard-clustering techniques (e.g., Willett 1985; Hearst and Pedersen 1996; Leuski 2001; Tombros et al. 2002; Liu and Croft 2004). In contrast, here we focus on using overlapping nearest-neighbor clusters that were shown to be effective when utilized in a query-independent fashion (Griffiths et al. 1986; Kurland and Lee 2004; Kurland et al. 2005; Kurland 2006), and which were recently used in the re-ranking setting (Kurland and Lee 2006; Liu and Croft 2006a, b).

Formally, for each document $d \in \mathcal{D}_{\text{init}}$ we define a cluster that contains d and the $k - 1$ documents d_i ($d_i \neq d$) from $\mathcal{D}_{\text{init}}$ that yield the highest language-model similarity $p_{d_i}(d)$ (we break ties by document IDs); k is a free parameter. Thus, the set $Cl(\mathcal{D}_{\text{init}})$ contains N (overlapping) clusters. We study the relative merits of this nearest-neighbor clustering approach with respect to hard-clustering techniques in Sect. 4.6.

2.1 Re-ranking algorithms

In what follows we adapt cluster-based retrieval algorithms that were originally designed by Kurland and Lee (2004) for use with *query-independent* (static) clusters, and were shown to be effective for ranking the entire corpus, to the re-ranking setting wherein the clusters are *query-specific*.

The original versions of the algorithms that we consider (Kurland and Lee 2004) operate on clusters that are most similar to the query—i.e., *top-retrieved clusters*—for anchoring the query-independent clustering information to the query at retrieval time. The variants that we present here, on the other hand, consider *all* clusters in $Cl(\mathcal{D}_{\text{init}})$ as these are constructed from $\mathcal{D}_{\text{init}}$ —documents that are the most highly ranked by some initial search.

In the algorithms that we present clusters play two different roles, namely *selection* of documents and *smoothing* of documents' language models (Kurland and Lee 2004).

¹ Some other work uses these language-model-based estimates for forming links between textual items and utilizing them with graph-based methods (Kurland and Lee 2005, 2006). We discuss the relation of our methods to these approaches in Sects. 3 and 4.

2.1.1 Cluster-based document selection

Most cluster-based document-selection algorithms aim to identify a subset of clusters that potentially contain a large number of relevant documents (Croft 1980; Willett 1985; Kurland and Lee 2004; Liu and Croft 2004). However, finding query-specific clusters that contain a high percentage of relevant documents is known to be a very hard task (Hearst and Pedersen 1996; Tombros et al. 2002; Liu and Croft 2004). One of the reasons is that query-specific clusters contain documents that are similar to the query to begin with.

Therefore, we will focus here on a different cluster-based selection approach, which exploits the structure induced by overlapping clusters. Specifically, if we think of clusters as potentially representing aspects manifested in the initial list $\mathcal{D}_{\text{init}}$, one might opt to rank high documents that exhibit as many such aspects as possible, specifically, documents that belong to many clusters. An alternative view for the potential in utilizing the structure induced by clusters might be based on the fact that documents that belong to many of the clusters exhibit (high) similarity to many other documents in the initial list $\mathcal{D}_{\text{init}}$. Thus, such documents could be considered as *central* with respect to the initial list—a notion recently explored via a graph-based framework and which was shown to be connected with relevance (Kurland and Lee 2005, 2006; Kurland 2006).

Utilizing the structure induced by clusters as described above is the idea underlying Kurland and Lee’s (2004) best-performing cluster-based selection method—the **bag-select** algorithm. In its original form, the bag-select algorithm ranks high documents from the (entire) corpus that exhibit high similarity to the query and that belong to many *top-retrieved* query-independent clusters. Dropping the notion of “top-retrieved clusters”, as we deal with query-specific clusters, we focus on the centrality of a document with respect to the initial list $\mathcal{D}_{\text{init}}$ as measured by its membership in clusters from $Cl(\mathcal{D}_{\text{init}})$.

As noted above, the original version of the bag-select algorithm (Kurland and Lee 2004) also takes into account the document-query similarity information. This is done for coping with the fact that the clusters in this work (Kurland and Lee 2004) are query-independent. Case in point, top-retrieved query-independent clusters might contain documents that do not pertain to the query, but which are similar to documents that are based on information not related to the query. While it might seem at a first glance that using document-query similarity information for re-ranking $\mathcal{D}_{\text{init}}$ is redundant, experimental results show that using this information is actually important. This finding is in line with some recent work on graph-based re-ranking (Kurland and Lee 2005). Indeed, some of the query-specific clusters might exhibit “aspects” not pertaining to the query, or more specifically, contain a high percentage of non-relevant documents. Therefore, using document-query similarity information, and hence considering document-specific characteristics, might help to ameliorate the overgeneralization caused by scoring documents based only on cluster-induced information. Further support to the importance of “query-anchoring” is given in work on *score regularization* for re-ranking (Diaz 2005), which shows that documents from the initial list that are highly similar both to the query and to other documents in the list that are similar to the query tend to be relevant.

Given the observations made at the above, we set the re-ranking version of the bag-select algorithm to score document d by

$$Score_{\text{bag-select}}(d) \stackrel{\text{def}}{=} p_d(q) \cdot \#(c \in Cl(\mathcal{D}_{\text{init}}) : d \in c).$$

The bag-select algorithm utilizes two sources of information: the number of clusters to which the document belongs and the document-query similarity. In the next section we

show how these two sources of information, along with additional ones, can be modeled and integrated using a probabilistic approach.

2.1.2 Cluster-based smoothing

In work on language models for ad hoc retrieval, several researchers have proposed to smooth the document language model with that of the cluster(s) with which it is associated (Azzopardi et al. 2004; Kurland and Lee 2004; Liu and Croft 2004; Wei and Croft 2006). The idea is to enrich the document representation with corpus-context information. Such an approach can help, for example, to deal with the *synonymy* problem, and more generally, with the *sparse data problem*. Applying cluster-based smoothing in the re-ranking setting with query-specific clusters means that the context-information is drawn from the initial list $\mathcal{D}_{\text{init}}$ rather than from the entire corpus. Hence, such an approach can be thought of as query-specific (cluster-based) smoothing: the information used for smoothing is drawn from documents that are (relatively) similar to the document in hand and to the query.

To study whether utilizing context from $\mathcal{D}_{\text{init}}$ to enrich a document representation yields effective re-ranking performance, we adapt Kurland and Lee’s (2004) *aspect-based* algorithms, which are named after the aspect models of Hofmann and Puzicha (1998). Aspect models are an approach for modeling a corpus based on the assumption that each document exhibits (or is “generated” by) a mixture of aspects. The algorithm for finding the aspects, in terms of language models, induces clustering of documents as it estimates document-aspect association probabilities, and aspects might be thought of as (soft) clusters.

Kurland and Lee (2004) conceptually adopt the basic formulation underlying the aspect models and use it with static (query-independent) existing clusters for ranking all documents in a corpus. Specifically, the **aspect-t** algorithm (Kurland and Lee 2004) is based on estimating the conditional probability $p(q|d)$ —often termed *query likelihood* (Song and Croft 1999). The idea is to estimate the probability that the query terms are generated by a (probabilistic) model induced from a document. Using simple probability rules, this probability can be written as

$$p(q|d) = \sum_{c \in \mathcal{C}(\mathcal{D}_{\text{init}})} p(q|d, c)p(c|d). \quad (1)$$

The basic conceptual assumption underlying aspect models is that a query is independent of a document given a cluster (Hofmann and Puzicha 1998). That is, the query terms could be viewed as being generated directly from models of clusters (aspects) that generate the terms in the document. Using this assumption we get that the above probability is

$$\sum_{c \in \mathcal{C}(\mathcal{D}_{\text{init}})} p(q|c)p(c|d). \quad (2)$$

Following recent work on cluster-based smoothing (Kurland and Lee 2004; Liu and Croft 2004; Tao et al. 2006), we post the constraint, which we will later relax, that a document can be “represented” (i.e., smoothed) *only* by the clusters to which it belongs.²

² Such a constraint can potentially alleviate the computational cost of estimating the document-cluster association strength for all available clusters and documents; this cost is significant when using, for example, static overlapping clusters (Kurland and Lee 2004). An implicit assumption underlying this constraint is that the best clusters to use for representing a document are those that contain it. We return to this point later on.

Thus, we truncate the summation from Eq. 2 (hence the suffix “-t” for “truncated”)³ and in addition use a language-model-based similarity measure for conditional probabilities to derive the aspect-t algorithm:

$$Score_{aspect-t}(d) \stackrel{def}{=} \sum_{c \in CI(\mathcal{D}_{init}): d \in c} p_c(q)p_d(c).$$

It is important to note that the original scoring function of the aspect-t algorithm (Kurland and Lee 2004) is slightly different than the one presented here, and not only due to the shift from using (top-retrieved) query-independent clusters to using *all* available query-specific clusters from $CI(\mathcal{D}_{init})$. Directly adapting Kurland and Lee’s model to the re-ranking setting, by using query-specific instead of query-independent clusters, yields the scoring function $\sum_{c: d \in c} p_c(q)p_c(d)$. This model is a result of using Bayes rule upon Eq. 2 and assuming uniform prior distributions for documents and clusters. Our formulation here, on the other hand, is not dependent on these assumptions, and, as it turns out, yields much better re-ranking performance than that of the originally suggested model (Kurland and Lee 2004).

The assumption that a query is independent of a document given a cluster can cause overgeneralization. That is, representing a document *only* via the clusters to which it belongs ignores potentially important information with regard to the document-specific characteristics. Such information can help to estimate the document-query “match”. Hence, we drop this independence assumption, and in addition (i) use the estimate $\lambda p(q|d) + (1 - \lambda) p(q|c)$ for $p(q|d, c)$ where λ is a free parameter, (ii) apply some probability algebra, and (iii) use a language-model-based similarity measure for conditional probabilities in Eq. 1, to derive Kurland and Lee’s best-performing model, **interpolation-t**:⁴

$$Score_{interpolation-t}(d) \stackrel{def}{=} \lambda p_d(q) + (1 - \lambda) \sum_{c \in CI(\mathcal{D}_{init}): d \in c} p_c(q)p_d(c).$$

Note that the interpolation-t algorithm anchors the cluster-based ranking of the aspect-t algorithm to the query by interpolation with the query-similarity score $p_d(q)$. This anchoring makes sense when query-independent clustering is used as in the original proposal of interpolation-t (Kurland and Lee 2004). However, as is the case for the bag-select algorithm from the above, and as we will be shown in Sect. 4, this anchoring has the potential to improve re-ranking effectiveness, even though the clusters are query-specific. This further demonstrates the importance in utilizing document-specific characteristics for ameliorating the overgeneralization caused by the use of clusters as proxies for ranking documents.

We also note that the interpolation-t algorithm can be conceptually viewed as a generalized version of the models of Liu and Croft (2004) and Wei and Croft (2006) that use cluster-based smoothing of document language models. (The former uses k-means clusters and the latter uses LDA clusters (Blei et al. 2003);⁵ note that the interpolation-t algorithm does not require the clusters to be overlapping.)

Document-cluster relationship Both the aspect-t and interpolation-t algorithms use clusters as “representatives” (“proxies”) of their constituent documents. However, if the different clusters are thought of as potentially representing different aspects manifested in

³ The aspect-based models were originally termed “aspect-x” (Kurland and Lee 2004).

⁴ The original name of this algorithm was *interpolation* (Kurland and Lee 2004).

⁵ We hasten to point out that the models in Liu and Croft (2004) and Wei and Croft (2006) operate at the term-level, that is, interpolation is performed upon estimates of term probabilities. In contrast, interpolation-t operates at the score level by fusion of language-model-based similarity scores.

the initial list $\mathcal{D}_{\text{init}}$, then a document can be associated (with varying degrees of strength) with different aspects (clusters) regardless of which clusters it belongs to. Thus, we consider the alternative of smoothing a document language model with the language models of *all* clusters in $Cl(\mathcal{D}_{\text{init}})$ to a degree controlled by the document-cluster language-model-based similarity. Doing so results in a formulation that is more “faithful” to the original probabilistic formulation in Eqs. 1 and 2 than those of the aspect-t and interpolation-t algorithms. (Recall that the latter two use truncation of the summation in Eqs. 1 and 2.) We thereby define the algorithms **aspect-f** and **interpolation-f** using the scoring functions: (“-f” stands for using the full summation in Eqs. 1 and 2)

$$Score_{\text{aspect-f}}(d) = \sum_{c \in Cl(\mathcal{D}_{\text{init}})} p_c(q)p_d(c),$$

and

$$Score_{\text{interpolation-f}}(d) = \lambda p_d(q) + (1 - \lambda) \sum_{c \in Cl(\mathcal{D}_{\text{init}})} p_c(q)p_d(c),$$

respectively.

2.2 Language-model-based similarity induction

In this section we present our estimate for the language-model similarity $p_x(y)$. For language model induction we treat documents and queries as term sequences.

While there are various approaches for representing clusters (Liu and Croft 2006b, 2008), our focus here is on the merits (or lack thereof) of our re-ranking methods. Therefore, we take the standard approach, which was shown to be effective in several applications of cluster-based retrieval (Kurland and Lee 2004, 2006; Liu and Croft 2004), and represent a cluster by the term sequence that results from concatenating its constituent documents; the order of concatenation has no effect since we are only going to define unigram language models that assume term independence.

We use $tf(w \in x)$ to denote the number of times that term w occurs in the text (or text collection) x . The *maximum likelihood estimate* (MLE) of w with respect to x is defined as

$$\tilde{p}_x^{MLE}(w) \stackrel{\text{def}}{=} \frac{tf(w \in x)}{\sum_{w'} tf(w' \in x)}.$$

To cope with the *zero probability problem*, namely, the assignment of zero probability to unseen terms, we adopt the widely used Dirichlet-smoothed estimate (Zhai and Lafferty 2001; Croft and Lafferty 2003):

$$\tilde{p}_x^{[\mu]}(w) \stackrel{\text{def}}{=} \frac{tf(w \in x) + \mu \cdot \tilde{p}_c^{MLE}(w)}{\sum_{w'} tf(w' \in x) + \mu};$$

μ is a free parameter that controls the amount of reliance on corpus statistics. We extend the estimate just described to a term sequence $\mathbf{w} = w_1 w_2 \cdots w_n$ using the term-independence assumption:

$$p_x^{[\mu]}(\mathbf{w}) \stackrel{\text{def}}{=} \prod_{j=1}^n \tilde{p}_x^{[\mu]}(w_j). \tag{3}$$

Using the estimate from Eq. 3 for estimating the similarity $p_x(y)$ will result in longer texts y being assigned lower similarity values than shorter texts are. Also, for very long

texts (as is the case for clusters, for example), the estimate might cause underflow problems. Therefore, we use a previously proposed estimate (Lavrenko et al. 2002; Kurland and Lee 2004, 2005), which is based on the Kullback Leibler divergence $D(\cdot||\cdot)$ (Cover and Thomas 1991)

$$p_x^{KL,\mu}(\mathbf{w}) \stackrel{\text{def}}{=} \exp\left(-D(\tilde{p}_w^{MLE}(\cdot)||\tilde{P}_x^{[\mu]}(\cdot))\right). \quad (4)$$

Using some probability algebra (see, for example, Lafferty and Zhai 2001), it can be shown that the estimate from Eq. 4 is equivalent to

$$p_x^{KL,\mu}(\mathbf{w}) = H(\mathbf{w}) \cdot p_x^{[\mu]}(\mathbf{w})^{\frac{1}{|\mathbf{w}|}}, \quad (5)$$

where H is the entropy function.

Thus, the estimate $p_x^{KL,\mu}(\mathbf{w})$ avoids the length-bias caused by the unigram language model through length normalization. Furthermore, the entropy of a document (language model) was shown to be connected with relevance in the re-ranking setting (Kurland and Lee 2005); hence, our similarity estimate “favors” documents that have a higher “prior” probability of being relevant to the query.

Also, it is important to point out that while the estimates $p_x^{KL,\mu}(\mathbf{w})$ and $p_x^{[\mu]}(\mathbf{w})$ are equivalent for the purpose of ranking documents in response to a fixed query (Lafferty and Zhai 2001), in the re-ranking setting this equivalence does not hold since we estimate similarities between different pairs of text items.

Finally, we note that while the estimate $p_x^{KL,\mu}(\mathbf{w})$ does not form a valid probability distribution, normalizing it for cases wherein one might be needed (e.g., for the distribution of clusters over a document in the aspect models) results in degraded re-ranking performance and therefore we use the estimate as is.

3 Related work

Preece (1973) was perhaps the first to suggest the use of query-specific clusters, although he did not present specific retrieval models for utilizing them.

Willett (1985) proposed to rank query-specific clusters and then to use the constituent documents of the highest-ranked ones to create a document-based ranking. He noted that the limited effectiveness of the approach could be attributed to the correlation-based ranking that was used to rank the clusters in response to the query. Liu and Croft (2004) took a similar approach for re-ranking, but used a language-model-based estimate for the query-cluster similarity; however, the resultant performance did not transcend that of the initial ranking. We compare the re-ranking performance of this cluster-ranking approach with that of the methods from Sect. 2 in Sect. 4.1.

Several researchers showed that if the documents at the top ranks of an initially retrieved list are clustered, then there is a cluster (a.k.a *the optimal cluster*) that if retrieved in its entirety, yields performance that is better than that of the initial ranking (Hearst and Pedersen 1996; Tombros et al. 2002; Kurland 2006). Moreover, such a cluster exists for different clustering approaches: partitioning (Hearst and Pedersen 1996), hierarchical agglomerative clustering (Tombros et al. 2002) and nearest-neighbor (soft) clustering (Kurland 2006, Chapter 7). While automatically detecting the optimal cluster is a difficult challenge (Willett 1985; Liu and Croft 2004; Kurland 2006, 2008; Liu and Croft 2006a; Kurland and Domshlak 2008), this clustering pattern—which gives support to van

Rijsbergen's cluster hypothesis (van Rijsbergen 1979) in the re-ranking setting—helps users to more quickly detect relevant documents if the results are visualized (and navigated) using cluster-based interfaces (Hearst and Pedersen 1996; Leuski 2001).

In work on cluster-based retrieval in the language modeling framework researchers have proposed to smooth a document language model with those of query-independent clusters so as to utilize corpus-context in representing documents (Azzopardi et al. 2004; Kurland and Lee 2004; Liu and Croft 2004; Tao et al. 2006; Wei and Croft 2006). Liu and Croft (2004) examined this cluster-based smoothing approach for re-ranking, having a document language model smoothed with that of the single query-specific (hard) cluster to which it belongs. As stated in Sect. 2, Liu and Croft's model can be viewed as a specific case of the interpolation-t algorithm when implemented with a hard clustering approach. Similarly, the re-ranking model of Lee et al. (2001), who use hierarchical agglomerative clustering, is also a special case of the interpolation-t algorithm: a document is scored by interpolation of its query-similarity score with the query-cluster similarity score of the single cluster to which it belongs; however, the clusters that are used are static query-independent clusters that are related to the query and not query-specific clusters. In Section 4.6 we present the relative merits of using nearest-neighbor overlapping clusters with respect to hard clusters for the interpolation-t and interpolation-f algorithms.

Query-specific clusters reflect inter-document similarities within the initial list. There has been an increasing use of graph-based techniques for modeling these inter-document similarities for document (re-) ranking (Daniłowicz and Baliński 2000; Kurland and Lee 2005; Zhang et al. 2005; Kurland and Lee 2006). The general idea is to identify documents that are *central* with respect to the initial list—i.e., similar to many (central) documents in the list—using graph-based methods, and use this centrality as criterion for ranking. Kurland and Lee (2006) show that it is more effective in general to incorporate both document-based and cluster-based information in the graphs than to use only the former as is the case in Daniłowicz and Baliński (2000) and Kurland and Lee (2005). Specifically, Kurland and Lee (2006) use the HITS (hubs and authorities) algorithm (Kleinberg 1998) over bipartite graphs of documents on the one side and query-specific clusters on the other side (with edge weights representing cluster-document similarities) to find central documents and clusters. They show that document authoritativeness (as induced by HITS) is connected with relevance and that authoritative query-specific clusters contain a high percentage of relevant documents. We compare the principles underlying their methods, and their performance, to those of ours in Sect. 4.5.

In a related vein, Baliński and Daniłowicz (2005) and Diaz (2005) apply score regularization to ensure that similar documents within an initially retrieved list receive similar scores. Recall that the interpolation-f algorithm assigns high scores to documents that are similar both to the query and to clusters that are similar to the query. Now, replacing clusters with documents (i.e., each document serves as a cluster), we get that a score of a document depends on its similarity to the query and on the similarity to the query of documents to which it is similar—the underlying principle of score regularization (Diaz 2005). We study this algorithm in Sect. 4.7.

Finally, it is important to note that a disadvantage of using query-specific clustering is the computational cost involved in creating the clusters. In contrast to offline clustering, wherein the clusters are created once and then used for all queries, with query-specific clustering each query requires a new clustering to be performed upon the list of retrieved documents. Therefore, several researchers proposed fast algorithms for clustering retrieved results (Cutting et al. 1992; Zamir and Etzioni 1998). The focus of the work in this paper, on the other hand, is on the potential effectiveness in exploiting clustering and not on the

efficiency of the clustering method. In fact, as we show in Sect. 4.6, our best performing algorithm (in terms of effectiveness) yields very good precision-at-top-ranks performance with several different clustering methods.

4 Evaluation

4.1 Experimental setup

We conducted our experiments on the following three TREC corpora (Voorhees and Harman 2005):

Corpus	# Of docs	Queries	Disks
AP	242,918	51–64, 66–150	1–3
TREC8	528,155	401–450	4–5 (–CR)
WSJ	173,252	151–200	1–2

These benchmarks were used in previous work on re-ranking (Kurland and Lee 2005, 2006); specifically, the methods in Kurland and Lee (2006) use query-specific-cluster information and will serve as reference comparison for our best performing re-ranking method. The TREC8 document collection is highly heterogeneous and is considered a very challenging benchmark with the 401–450 queries (Hu et al. 2003; Kurland et al. 2005; Voorhees 2005); AP and WSJ, on the other hand, which are composed of only news articles, are considered to be more homogeneous (Voorhees and Harman 2005).

We applied basic tokenization and Porter stemming (Porter 1980) via the Lemur toolkit (www.lemurproject.org), which we also used for language-model induction. Topic titles served as queries.

Since we are interested in the re-ranking effectiveness of our approaches when applied to a relatively short initial list $\mathcal{D}_{\text{init}}$, we focus on evaluation metrics that measure the precision at top ranks of the resultant document list. Specifically, we use the precision at the top 5 and 10 documents (henceforth $\text{prec}@5$ and $\text{prec}@10$, respectively) and the mean reciprocal rank (MRR) of the first relevant document (Shah and Croft 2004). To determine statistically significant differences in performance, we use the Wilcoxon two-tailed test at a confidence level of 95%.

To facilitate the comparison with some recent work on cluster-based re-ranking (Kurland and Lee 2006), we have taken the *exact* same experimental design choices, which are intended for verifying the general validity of the re-ranking principles we explore rather than for engaging in excessive parameter tuning:

- To create the initial list $\mathcal{D}_{\text{init}}$ upon which re-ranking is performed, we use the estimate $p_d^{KL,\mu}(q)$ to rank all documents in the corpus. The value of μ is chosen (using exhaustive search) to optimize the MAP of the top 1,000 retrieved documents with respect to the given set of queries; the idea is to have an initial list of a reasonable quality. In fact, such ranking yields precision at top ranks performance that is statistically indistinguishable from—though lower in absolute terms than—the results obtained by optimizing μ with respect to the evaluation metrics that we focus on. We further discuss this at the below. In the experiments to follow, we set $\mathcal{D}_{\text{init}}$ to be the 50

Table 1 Summary of re-ranking algorithms

bag-select	$p_d(q) \cdot \#(c \in Cl(\mathcal{D}_{init}) : d \in c)$
aspect-t	$\sum_{c \in Cl(\mathcal{D}_{init}):d \in c} p_c(q) \cdot p_d(c)$
aspect-f	$\sum_{c \in Cl(\mathcal{D}_{init})} p_c(q) \cdot p_d(c)$
interpolation-t	$\lambda \cdot p_d(q) + (1 - \lambda) \sum_{c \in Cl(\mathcal{D}_{init}):d \in c} p_c(q) \cdot p_d(c)$
interpolation-f	$\lambda \cdot p_d(q) + (1 - \lambda) \sum_{c \in Cl(\mathcal{D}_{init})} p_c(q) \cdot p_d(c)$

highest-ranked documents according the above criterion and use the term *initial ranking* to refer to the ranking by which \mathcal{D}_{init} was created.

- The smoothing parameter μ , which controls the language-model-based similarity estimate, is set to 2000 in all (re-ranking) algorithms following the recommendation in Zhai and Lafferty (2001), except for when estimating $p_d(q)$ where we use the value chosen for creating the initial list so as to maintain consistency.
- We only optimized settings for k (cluster size) and λ (the interpolation parameter in the interpolation-t and interpolation-f algorithms) with respect to precision at the top 5 documents,⁶ not with respect to all three evaluation metrics employed. While this approach reflects a more realistic experimental setting than one wherein results are presented for different optimized settings, our prec@10 and MRR results are therefore not as high as they could have potentially been. The values of k and λ were selected from {2, 5, 10, 20, 30} and {0, 0.1, ..., 0.9}, respectively. In Sect. 4.3 we study the performance sensitivity of our best performing re-ranking method with respect to the values of k and λ .

In Table 1 we summarize the scoring functions of the re-ranking algorithms that we presented in Sect. 2 for convenience of reference. Whenever a re-ranking method assigns different documents with the same score, we break the ties by document ID.

4.2 From query-independent to query-specific clusters

We first present the re-ranking performance numbers of the bag-select, aspect-t and interpolation-t algorithms in Table 2. Recall that these algorithms are adaptations of methods originally proposed for use with static query-independent clusters.

The first three rows in Table 2 specify reference-comparison data. The *empirical upper bound on re-ranking*, which applies to any algorithm that re-ranks \mathcal{D}_{init} , indicates the performance attained by placing all relevant documents from \mathcal{D}_{init} at the top of the resultant list. The initial ranking, as mentioned above, was produced using $p_d^{KL;\mu}(q)$ with μ chosen to optimize MAP at 1,000. We also present for comparison the performance results of the CQL algorithm (Liu and Croft 2004), which is a cluster-based selection method. The CQL algorithm first ranks the clusters in $Cl(\mathcal{D}_{init})$ by their similarity to the query $p_c(q)$ (ties are broken by cluster ID), and then replaces each cluster with its constituent documents, omitting repeats; documents within a cluster are ordered by their similarity to the query $p_d(q)$. Thus, CQL echoes some previous work on cluster-based re-ranking (Willett 1985) as mentioned in Sect. 3.

Our first observation from Table 2 is that the standard cluster-selection approach that is represented by the CQL method does not yield effective re-ranking performance. (Compare

⁶ If two different parameter settings yield the same prec@5, we choose the one *minimizing* prec@10 so as to provide conservative estimates of expected performance. Similarly, in case of ties for both prec@5 and prec@10, we choose the setting minimizing MRR.

Table 2 Performance numbers of re-ranking algorithms

	AP			TREC8			WSJ		
	prec@5	prec@10	MRR	prec@5	prec@10	MRR	prec@5	prec@10	MRR
upper bound	.876	.788	.930	.944	.850	.980	.896	.800	1.000
init. ranking	.457	.432	.596	.500	.456	.691	.536	.484	.748
CQL	.448	.418	.549*	.500	.432	.723	.504	.454*	.680
bag-select	.507	.494*	.630	.532	.514*	.660	.548	.488	.719
aspect-t	.517*	.496*	.654	.548	.484	.688	.528	.496	.689
interpolation-t	.527*	.499*	.651	.564*	.494	.707	.568	.490	.725

For each evaluation setting, improvements over the initial ranking are given in *italics*; statistically significant differences with the initial ranking are indicated by *; bold highlights the best results per column

CQL's performance with that of the initial ranking.⁷) These results are in line with previous findings (Willett 1985; Liu and Croft 2004). In contrast, the bag-select algorithm, which is a cluster-based selection method that utilizes the structure induced by overlapping clusters, outperforms the initial ranking in most relevant comparisons (corpus \times evaluation measure).

We can also see in Table 2 that both cluster-based smoothing methods, aspect-t and interpolation-t, outperform the initial ranking in most of the relevant comparisons. Furthermore, both algorithms are in general more effective than the CQL and bag-select algorithms, which implies that cluster-based smoothing is more effective for re-ranking than cluster-based selection. This finding is in accordance with reports on utilizing static query-independent clusters for ranking all documents in a corpus (Kurland and Lee 2004).

The interpolation-t algorithm is the best performing re-ranking method among those presented in Table 2. Its performance is on many occasions substantially better than that of the initial ranking (sometimes to a statistically significant degree). Recalling that interpolation-t interpolates the score of the aspect-t algorithm with a query-similarity score, and observing that the former is more effective than the latter in most relevant comparisons, attests that using document-query similarity information in the re-ranking setting can help to improve effectiveness. (Refer back to the discussion in Sect. 2; we hasten to point out, however, that interpolation-t incorporates an additional free parameter (λ) on those used by the aspect-t algorithm.)

4.3 Document-cluster relationship

Recall from Sect. 2 that in both the aspect-t and interpolation-t algorithms clusters can represent (or smooth the language models of) only documents that they contain. We now study the alternative of smoothing a document language model with those of all available clusters (to an extent controlled by the document-cluster similarity) as is the case in the aspect-f and interpolation-f algorithms.

The results in Table 3 clearly indicate that using all the clusters for smoothing is superior to using only the clusters to which a document belongs. (Note that most of the underlined numbers that indicate which of the “-f” and “-t” versions is superior appear in “-f” rows.) Thus, we see that even though the aspect-t and interpolation-t algorithms

⁷ The performance of CQL can be improved if different cluster representations are used (Liu and Croft 2006b, 2008), as is the case for some other cluster-based retrieval algorithms (Kurland and Domshlak 2008). However, experimenting with different cluster representations is out of the scope of this paper.

smooth a document language model with those of (a few overlapping) clusters to which it belongs, it is better to use all clusters for smoothing; this finding further supports the importance of the context drawn from $\mathcal{D}_{\text{init}}$ for representing documents.

We can also see in Table 3 that interpolation-f outperforms aspect-f in most relevant comparisons—as was the case when comparing interpolation-t and aspect-t—showing again the importance of using document-query similarity information for re-ranking. Furthermore, interpolation-f is the best performing re-ranking algorithm among all those considered. Specifically, in terms of prec@5—the metric for which performance was optimized—interpolation-f always improves on the initial ranking by a wide margin that is also statistically significant; in fact, interpolation-f is the only algorithm that achieves statistically significant prec@5 improvement over the initial ranking for all corpora. In addition, the prec@5 performance of interpolation-f is also substantially better than that of document retrieval performed over the *entire* corpus using $p_d^{KL,\mu}(q)$ with μ chosen to optimize prec@5; the prec@5 performance numbers of this optimized baseline are .465, .512, .560, for AP, TREC8, and WSJ, respectively.

Further support to the effectiveness of the interpolation-f algorithm is given in Table 4, wherein we present its MAP performance in comparison to that of the initial ranking. (MAP was calculated at cutoff 50; therefore, the recall of interpolation-f and the initial ranking is equivalent at this cutoff, and consequently, the relative ordering of relevant documents is the only factor affecting MAP performance. In addition, recall that the values of the free parameters of interpolation-f were chosen to optimize prec@5 and *not* MAP, although the former has impact on the latter.) As we can see in Table 4, the MAP performance of interpolation-f is better than that of the initial ranking to a statistically significant degree on AP and WSJ; for TREC8, the MAP performance difference between interpolation-f and that of the initial ranking is not statistically distinguishable.

Table 3 Comparison between the “truncated” (-t) and “full” (-f) versions of the aspect and interpolation algorithms

	AP			TREC8			WSJ		
	prec@5	prec@10	MRR	prec@5	prec@10	MRR	prec@5	prec@10	MRR
init. ranking	.457	.432	.596	.500	.456	.691	.536	.484	.748
aspect-t	.517*	.496*	.654	.548	.484	.688	.528	.496	.689
aspect-f	.537*	.498*	.628	.560*	.504*	.714	.576	.504	.759
interpolation-t	.527*	.499*	.651	.564*	.494	.707	.568	.490	.725
interpolation-f	.537*	.498*	.628	.576*	.496	.687	.592*	.508	.767

Underline: best result in a “block” (corpus × algorithm × evaluation measure). Boldface: best result per column. Statistically significant differences with the initial ranking are marked with “*.”

Table 4 Comparison of the MAP performance (at cutoff 50) of interpolation-f with that of the initial ranking

	AP	TREC8	WSJ
init. ranking	.093	.175	.222
interpolation-f	.100*	.174	.238*

* Marks statistically significant differences with the initial ranking

We further study the interpolation-f algorithm by examining its performance sensitivity to the choice of cluster size (k) and the interpolation parameter (λ). In Fig. 1 we present the algorithm’s prec@5 performance (the metric for which we optimized performance) when either fixing k (and optimizing with respect to λ) or fixing λ (and optimizing with respect to k). We can clearly see that in both cases, and for every considered value of each of the parameters, the performance of interpolation-f is well beyond that of the initial ranking (depicted with horizontal line in all cases). In general, clusters of size 10 and $\lambda \in \{0.6, 0.7\}$ guarantee (near) optimal prec@5 performance over all tested corpora.

4.4 Comparison with pseudo-feedback-based retrieval

The cluster-based methods that we presented exploit information from $\mathcal{D}_{\text{init}}$ for re-ranking it. Pseudo-feedback-based query-expansion methods (Buckley et al. 1994; Ruthven and Lalmas 2003), on the other hand, exploit information from $\mathcal{D}_{\text{init}}$ for defining a query model and use it for ranking the entire corpus. We contrast the two paradigms by comparing the performance of our best-performing method, interpolation-f, to that of using a *relevance model* (Lavrenko and Croft 2001), which is considered to be a state-of-the-art pseudo-feedback-based approach.

We use Lemur’s (www.lemurproject.org) implementation of relevance model number 1 (RM1), which follows the details in Lavrenko and Croft (2003), and which is based on the I.I.D sampling method.

Specifically, following Lavrenko and Croft (2003) we use a Jelinek–Mercer-smoothed document language model to estimate the probability assigned to term w by document $d \in \mathcal{D}_{\text{init}}$:

$$\tilde{p}_d^{JM;[\alpha]}(w) \stackrel{\text{def}}{=} \alpha \tilde{p}_d^{MLE}(w) + (1 - \alpha) \tilde{p}_c^{MLE}(w);$$

α is a free parameter. We then define the probability assigned to w by the relevance model \mathcal{R} as:

$$\tilde{p}_{\mathcal{R}}(w; \alpha) \stackrel{\text{def}}{=} \sum_{d \in \mathcal{D}_{\text{init}}} \tilde{p}_d^{JM;[\alpha]}(w) p(d|q);$$

for query $q = q_1, \dots, q_l$ of length l , $p(d|q)$ is the normalized query likelihood $\prod_i \tilde{p}_d^{JM;[\alpha]}(q_i) / \sum_{d' \in \mathcal{D}_{\text{init}}} \prod_i \tilde{p}_{d'}^{JM;[\alpha]}(q_i)$, which is based on a uniform-distribution assumption for documents in $\mathcal{D}_{\text{init}}$.

An additional step often employed for improving the relevance model’s performance (Connell et al. 2004; Cronen-Townsend et al. 2004; Metzler et al. 2005; Diaz and Metzler 2006). is *term clipping*: we define the probability $\tilde{p}_{\mathcal{R}}(w; \alpha, \beta)$ to be zero, except for the β terms with the highest $\tilde{p}_{\mathcal{R}}(w; \alpha)$, for which we define $\tilde{p}_{\mathcal{R}}(w; \alpha, \beta)$ as the normalized $\tilde{p}_{\mathcal{R}}(w; \alpha)$ so as to have a valid probability distribution over these β terms; β is a free parameter.

Also, as in prior work on relevance models (Abdul-Jaleel et al., 2004; Diaz and Metzler 2006), we examine a variant that interpolates the relevance model with the query-likelihood for avoiding *query drift*; this results in the so-called “relevance model number 3” (RM3), which we will denote \mathcal{IR} (for interpolated relevance-model):

$$\tilde{p}_{\mathcal{IR}}(w; \alpha, \beta, \gamma) \stackrel{\text{def}}{=} \gamma \tilde{p}_q^{MLE}(w) + (1 - \gamma) \tilde{p}_{\mathcal{R}}(w; \alpha, \beta);$$

(γ is a free parameter; $\gamma = 0$ amounts to using only \mathcal{R} .)

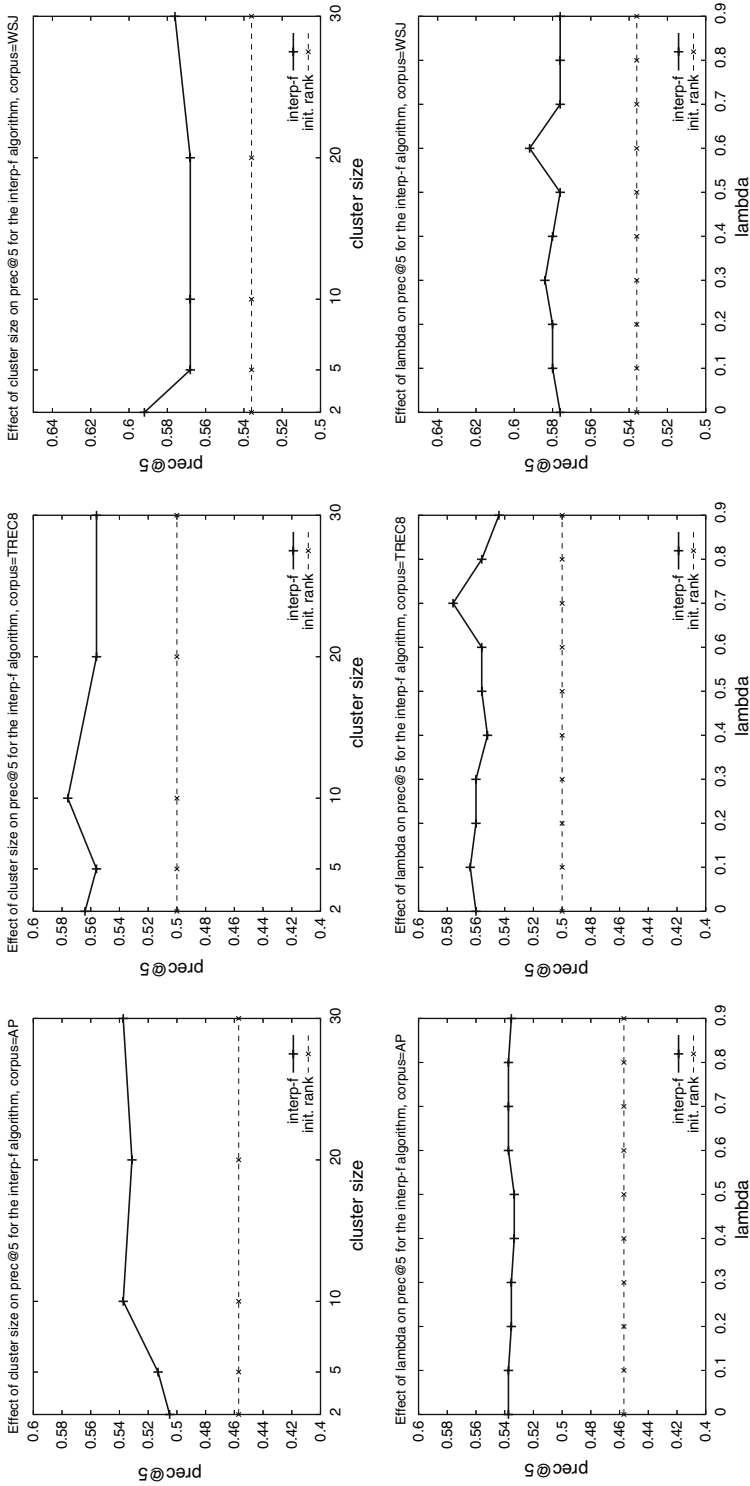


Fig. 1 The prec@5 performance curves of the interpolation-f algorithm. We either fix the cluster size k to a value in {2, 5, 10, 20, 30} (first-row figures), or fix λ (second-row figures) to a value in {0, 0.1, ..., 0.9}; $\lambda = 0$ amounts to the aspect-f algorithm. The prec@5 performance of the initial ranking is drawn for reference as horizontal line in all figures. *Note:* figures are not to the same scale

Then, to score documents using the interpolated relevance model (henceforth referred to as simply “relevance model”), we use the Kullback Leibler divergence $D(\tilde{p}_{TR}(\cdot; \alpha, \beta, \gamma) || \tilde{p}_d^{[\mu]}(\cdot))$, with $\mu = 2,000$ as for all our re-ranking algorithms. (Refer back to Sect. 2.2 for details on the Dirichlet-smoothed estimate $\tilde{p}_d^{[\mu]}(\cdot)$.)

The values of the free parameters that the relevance model is based on, namely α , β and γ , were chosen from $\{0.1, 0.3, 0.5, 0.7, 0.9, 1\}$, $\{25, 50, 75, 100, 500, 1000, 5000, ALL\}$ where “ALL” stands for all terms in the corpus (i.e., no clipping), and $\{0, 0.1, 0.2, \dots, 0.9\}$, respectively. The chosen parameter values are those that result in an optimized $\text{prec}@5$ performance, following the optimization procedure that was described in Sect. 4.1.

In addition to using the relevance model to rank all documents in the corpus—a version which we will refer to as **Rel Model**—as is standard practice, we examine a version that uses the relevance model to only re-rank documents in $\mathcal{D}_{\text{init}}$ —as is the case for our re-ranking algorithms—which we refer to as **Rel Model (re-rank)**. (Performance optimization for each of the two relevance models was carried out independently.) The performance results of the relevance model are presented in Table 5.

As we can see in Table 5, the interpolation-f algorithm outperforms the relevance model in almost all relevant comparisons (corpora \times evaluation metric), whether the relevance model is used to rank all documents in the corpus, or is used to re-rank $\mathcal{D}_{\text{init}}$. (However, the performance differences are not statistically significant.) In addition, recall that the performance of interpolation-f was optimized with respect to two free parameters, while that of the relevance models was optimized with respect to three. Furthermore, while the performance of interpolation-f is better over all corpora to a statistically significant degree than that of the initial ranking with respect to $\text{prec}@5$ —the metric for which we optimized performance—this is not the case for both versions of the relevance model; specifically, for TREC8, both versions of the relevance model do not post a statistically significant $\text{prec}@5$ improvement over the initial ranking, and are substantially inferior to interpolation-f, although not to a statistically significant degree.

4.5 Comparison to graph-based approaches for re-ranking

Recently, Kurland and Lee (2006) proposed a graph-based approach to re-ranking, which utilizes query-specific clusters. They define a bipartite graph wherein vertices are documents from $\mathcal{D}_{\text{init}}$ and clusters from $Cl(\mathcal{D}_{\text{init}})$, and edges are drawn from a cluster to the α

Table 5 Performance comparison of the interpolation-f algorithm with a relevance model (RM3), which is used either to rank all documents in the corpus (Rel Model), or to re-rank documents in $\mathcal{D}_{\text{init}}$ (Rel Model (re-rank))

	AP			TREC8			WSJ		
	prec@5	prec@10	MRR	prec@5	prec@10	MRR	prec@5	prec@10	MRR
init. ranking	.457	.432	.596	.500	.456	.691	.536	.484	.748
Rel Model	.503*	.486*	.585	.536	.462	.653	.588*	.510	.739
Rel Model (re-rank)	.511*	.482*	.598	.536	.470	.649	.588*	.506	.741
interpolation-f	.537*	.498*	.628	.576*	.496	.687	.592*	.508	.767

Boldface marks the best performance in a column; statistically significant difference with the initial ranking is marked with a “*”

documents it most resembles to (α is a free parameter), in a language model sense—i.e., $p_d(c)$; $p_d(c)$ also serves as a weight function for edges. Then, Kleinberg’s HITS algorithm (*hubs* and *authorities*) (Kleinberg 1998) is used to assign documents with authority scores and clusters with hub scores; the idea is to simultaneously identify the documents and clusters that are most *central* with respect to the initial list \mathcal{D}_{init} —i.e., most reflecting the context of \mathcal{D}_{init} and thereby potentially most relevant to the query—and to score documents by their centrality (i.e., authority scores).

To re-anchor centrality information to the query, Kurland and Lee (2006) also suggest to score a document by scaling its authority score with its query-similarity score $p_d^{KL,\mu}(q)$. Thus, this approach ranks high documents that are both similar to the query and similar to (many) clusters that are similar to (many) *documents*. Note that this is reminiscent of the ranking method of the interpolation-f algorithm, which assigns high scores to documents that are both similar to the query and similar to many clusters that are similar to the *query*. The main difference between the approaches is the way by which clusters are relatively weighted. In the case of the interpolation-f algorithm, cluster’s “importance” is determined based on its similarity to the query. In the graph-based approach, on the other hand, cluster’s importance is determined based on its similarity to central documents in \mathcal{D}_{init} .

We use the exact same implementation and optimization details as those described in Kurland and Lee (2006) to score documents by either their authority scores ($auth(d)$) or by scaling the authority scores with the query-similarity score ($auth(d) \cdot p_d^{KL,\mu}(q)$). Specifically, the cluster size k was set to $\{2, 5, 10, 20, 30\}$ as for our re-ranking algorithms, and the graph-out-degree α (i.e., the number of edges with a non-zero weight that connect a cluster to documents) was set to $\{2, 4, 9, 19, 29, 39, 49\}$; the free parameters (k and α) values were determined using the optimization procedure described in Sect. 4.1, which is exactly the same as that in Kurland and Lee (2006).

We present in Table 6 the performance numbers of scoring a document by either of the two mentioned authority-based approaches, along with the performance numbers of the interpolation-f algorithm.

In comparing the interpolation-f algorithm with using only authority scores, we see in Table 6 that interpolation-f is noticeably better in terms of $prec@5$ —the metric for which performance is optimized—for both TREC8 and WSJ. (For AP, using only authority scores results in somewhat better performance than that of interpolation-f.) In comparison to scaling the authority scores with query-similarity scores, interpolation-f posts $prec@5$ performance that is at least as good as that of the former for all corpora, with favorable noticeable difference on WSJ. It is also important to note that in contrast to the graph-based

Table 6 Performance comparison of the interpolation-f algorithm with re-ranking by authority scores (as induced over cluster-document graphs) and by scaling the authority scores by the document-query similarity score

	AP			TREC8			WSJ		
	prec@5	prec@10	MRR	prec@5	prec@10	MRR	prec@5	prec@10	MRR
init. ranking	.457	.432	.596	.500	.456	.691	.536	.484	.748
$auth(d)$.541*	.501*	.669	.544	.452	.674	.564	.514	.746
$auth(d) \cdot p_d^{KL,\mu}(q)$.537*	.493*	.630	.572*	.490	.702	.572	.510	.771
interpolation-f	.537*	.498*	.628	.576*	.496	.687	.592*	.508	.767

Boldface: best result per column; “*” indicates a statistically significant difference with the initial ranking

methods, interpolation-f attains statistically significant improvements over the initial ranking with respect to $\text{prec}@5$ for *all* corpora.

4.6 Alternative clustering schemes

Heretofore, we have focused on nearest-neighbor overlapping clusters. However, as mentioned in Sect. 2, the aspect and interpolation algorithms are not committed to a specific clustering approach. Indeed, using the interpolation-t algorithm with hard clusters, for example, amounts (modulo smoothing details) to the CBDM model⁸ of Liu and Croft (2004), which scores document d by $\lambda p_d(q) + (1 - \lambda)p_c(q)$, where c is the single hard cluster to which d belongs; Liu and Croft use CBDM for the re-ranking task with hierarchical agglomerative clustering algorithms. Using the interpolation-f algorithm with hard clusters, on the other hand, results in smoothing a document language model with those of hard clusters to which it does not belong, to an extent controlled by the document-cluster similarity. If we think of clusters as exhibiting different (query-related) aspects manifested in the initial list, then such a smoothing approach lets a document be represented by several of these “aspects” rather than by (potentially) a single one. Hence, this approach could be considered as ameliorating the overgeneralization caused by associating a document with a single cluster (van Rijsbergen 1979).

We now study the effect of the clustering approach on the performance of the interpolation algorithm for both its variants, namely, interpolation-t and interpolation-f.

The first hard clustering approach that we consider is hierarchical agglomerative clustering, which was the main focus of previous work on re-ranking (Willett 1985; Leuski 2001; Tombros et al. 2002; Liu and Croft 2004). We use a bottom-up approach and the following criteria for merging clusters: single link, complete link, average distance, distance between centroids and the Ward criterion (El-Hamdouchi and Willett 1986); we refer to these criteria using the abbreviations **agg-single**, **agg-comp**, **agg-avg**, **agg-centroid** and **agg-ward**, respectively. To produce a list of non-overlapping clusters, we stop the merging process when the number of clusters is a value in $\{2, 5, 10, 25\}$ so as to roughly result in an average cluster size equivalent to that used in the nearest-neighbor clustering method that we have utilized so far.

We also use the **k-means** clustering algorithm in deference to some previous work on visualization of retrieved results, which utilizes partitioning algorithms (Cutting et al. 1992; Hearst and Pedersen 1996). We set k to a value in $\{2, 5, 10, 25\}$ to comply with the choice made for the agglomerative clustering methods from above.

The clustering algorithms from above require a symmetric similarity measure, and are usually implemented using a vector space representation. We therefore use a standard tf.idf representation and the cosine similarity function that was used in previous work on re-ranking with hard clusters (Willett 1985; Leuski 2001; Tombros et al. 2002; Liu and Croft 2004). For completeness of comparison, we also implement a nearest-neighbor clustering method using the same vector space representation and the cosine measure, and set the number of clusters k to $\{2, 5, 10, 20, 30\}$, as was the case with the original language model implementation; we use **nn-VS** to denote the vector-space-based nearest-neighbor clustering and **nn-LM** to denote the original language-model-based nearest-neighbor clustering.

⁸ In its original form, the CBDM method works at the term level, in contrast to the interpolation-t algorithm that fuses language-model-based similarity scores.

Table 7 Performance numbers of the interpolation-t and interpolation-f algorithms with different clustering methods

	AP			TREC8			WSJ		
	prec@5	prec@10	MRR	prec@5	prec@10	MRR	prec@5	prec@10	MRR
<i>(a) Performance numebrs of the interpolation-t algorithm</i>									
init. ranking	.457	.432	.596	.500	.456	.691	.536	.484	.748
nn-LM	.527*	.499*	.651	.564*	.494	.707	.568	.490	.725
nn-VS	.519*	.474	.644	.524	.446	.748	.584	.492	.761
agg-single	.521*	.488*	.620	.528	.490*	.662	.580	.530	.783
agg-comp	.493	.467	.600	.524	.468	.674	.544	.488	.740
agg-avg	.525*	.490*	.619	.540	.504*	.675	.572	.532*	.765
agg-centroid	.523*	.487*	.592	.528	.492*	.662	.592	.530	.773
agg-ward	.491	.465	.587	.504	.446	.694	.584	.508	.715
k-means	.475	.459	.579	.532	.486*	.720	.568	.502	.743
<i>(b) Performance numbers of the interpolation-f algorithm</i>									
init. ranking	.457	.432	.596	.500	.456	.691	.536	.484	.748
nn-LM	.537*	.498*	.628	.576*	.496	.687	.592*	.508	.767
nn-VS	.523*	.490*	.642	.572*	.476	.760	.608*	.544*	.758
agg-single	.493	.482*	.582	.516	.468	.704	.552	.476	.747
agg-comp	.513*	.480*	.613	.564*	.508*	.750	.584	.534	.771
agg-avg	.503	.482*	.599	.524	.468	.738	.568	.506	.767
agg-centroid	.493	.482*	.584	.508	.466	.673	.556	.494	.748
agg-ward	.525*	.491*	.612	.536	.484	.685	.592	.530	.762
k-means	.509	.465	.584	.556*	.482	.701	.580	.520	.753

Boldface: best result per column; **: significant difference with the initial ranking

In Table 7 we present the performance numbers of the interpolation-t and interpolation-f algorithms when utilizing the different clustering methods.⁹ Table 8 then summarizes the relative performance patterns of interpolation-t and interpolation-f: each entry depicts in non ascending order of performance the algorithms that improve on the initial ranking with a relative difference of 2.5% or more; a hat (^) indicates that the difference is statistically significant.

Our first observation with respect to Table 7 is that for all clustering methods in most of the relevant comparisons (evaluation measure × corpora) both versions of the interpolation algorithm (interpolation-t and interpolation-f) improve on the initial ranking, thereby demonstrating its effectiveness as a general cluster-based re-ranking approach.

We can also see in Table 7b that for the interpolation-f algorithm, overlapping nearest-neighbor clusters are more effective than hard clusters. Specifically, both nearest-

⁹ The presented results for all clustering approaches are based on the optimization criterion that was described in Sect. 4.1.

Table 8 Comparison of the interpolation-t (T) and interpolation-f (F) algorithms when utilizing different clustering methods

	nn-LM	nn-VS	agg-single	agg-comp	agg-avg	agg-centroid	agg-ward	k-means	
AP	prec@5	$\hat{F}\hat{T}$	$\hat{F}\hat{T}$	$\hat{T}F$	$\hat{F}T$	$\hat{T}F$	$\hat{T}F$	$\hat{F}\hat{T}$	FT
	prec@10	$\hat{T}\hat{F}$	$\hat{F}T$	$\hat{T}\hat{F}$	$\hat{F}T$	$\hat{T}\hat{F}$	$\hat{T}F$	$\hat{F}\hat{T}$	FT
	MRR	TF	TF	T	F	T		FT	
TREC8	prec@5	$\hat{F}\hat{T}$	$\hat{F}T$	TF	$\hat{F}T$	TF	T	FT	$\hat{F}T$
	prec@10	FT	F	$\hat{T}\hat{F}$	$\hat{F}T$	$\hat{T}\hat{F}$	\hat{T}	FT	$\hat{T}\hat{F}$
	MRR		FT		F	F			T
WSJ	prec@5	$\hat{F}T$	$\hat{F}T$	TF	F	TF	TF	FT	FT
	prec@10	F	\hat{F}	T	F	$\hat{T}\hat{F}$	T	FT	FT
	MRR	F		T	F	F	T		

Each entry of the table depicts in non-ascending order of performance the algorithm(s) among the two (if any) that post a 2.5% (or more) relative performance improvement over the initial ranking (a hat (“^”) indicates that the improvement is significant). Bold highlights the best performing algorithm per entry

neighbor-based implementations of interpolation-f (nn-VS and nn-LM) are the only methods among those in Table 7 that yield prec@5—the metric for which performance was optimized—performance that is better to a statistically significant degree than that of the initial ranking for all corpora. For the interpolation-t algorithm (refer to Table 7a), the best performance numbers are often obtained by a nearest-neighbor clustering approach (refer to the boldfaced numbers), although the nn-VS version is not the top-performing clustering method among those that are based on a vector-space representation.

The comparison of the interpolation-t and interpolation-f algorithms in Table 8 shows that except for the cases of agg-single, agg-avg, and agg-centroid, interpolation-f is superior to interpolation-t (“F” is positioned to the left of “T”). This finding gives further support to the argument from above regarding the potential in using interpolation-f to ameliorate the overgeneralization caused by associating a document with, and representing it by, the single cluster to which it belongs. Additional exploration reveals that for the agg-single, agg-avg and agg-centroid cases, the clustering usually results in one large cluster along with very few small ones, each of which often contains a single document. Thus, it seems that in these cases the interpolation-t algorithm “benefits” from separating what could be considered as “outliers” from the rest of the documents, while the interpolation-f algorithm does not, since for each document all clusters are considered.

All in all, we believe that the most important message arising from the above analysis is that the interpolation algorithm is a highly effective paradigm for re-ranking that can utilize different clustering methods, whether they result in (soft) overlapping clusters or hard clusters; nearest-neighbor (overlapping) clusters, however, do seem to be a better choice in general than hard clusters for cluster-based smoothing in the interpolation algorithm.

4.7 Clusters-mediated similarity versus distinct document similarity

The aspect-f algorithm assigns document d the score $\sum_{c \in Cl(D_{init})} p_c(q)p_d(c)$. The interpolation-f algorithm assigns d the score $\lambda p_d(q) + (1 - \lambda) \sum_{c \in Cl(D_{init})} p_c(q)p_d(c)$. In these two scoring functions, clusters only play the role of smoothing: their language models are used to smooth d 's language model so as to provide “context” from the list D_{init} .

Table 9 Performance comparison of the aspect-f and interpolation-f algorithms with nearest-neighbors (in LM space) clusters (“nn-LM”) versus singleton (“single.”) clusters, wherein each document serves as a cluster

	AP			TREC8			WSJ		
	prec@5	prec@10	MRR	prec@5	prec@10	MRR	prec@5	prec@10	MRR
init. ranking	.457	.432	.596	.500	.456	.691	.536	.484	.748
aspect-f (single.)	.499	.477	.622	.508	.448	.653	.520	.484	.687
aspect-f (nn-LM)	.537*	.498*	.628	.560*	.504*	.714	.576	.504	.759
interpolation-f (single.)	.515	.497*	.615	.536	.482	<u>.698</u>	.564	.510	.696
interpolation-f (nn-LM)	.537*	.498*	.628	.576*	.496	.687	.592*	.508	.767

Underline: best performance within a block (algorithm × cor-pus × evaluation metric). Boldface: best performance per column; *: statistically significant difference with the initial ranking

We therefore ask now the following question: if clusters are indeed required only for providing context within the list \mathcal{D}_{init} , can single documents play the same role? We study this question by simply defining *singleton* clusters, i.e., each document in \mathcal{D}_{init} serves as a cluster. Then, the scoring functions of the aspect-f and interpolation-f algorithms become $\sum_{d_i \in \mathcal{D}_{init}} p_{d_i}(q)p_{d_i}(d_i)$ and $\lambda p_a(q) + (1 - \lambda) \sum_{d_i \in \mathcal{D}_{init}} p_{d_i}(q)p_{d_i}(d_i)$, respectively.

Note that in this version the aspect-f algorithm assigns relatively high scores to documents that are similar (to a large extent) to many other documents in \mathcal{D}_{init} that are similar to the query, and interpolation-f integrates this information with direct similarity to the query.

Table 9 presents the performance results of using documents as singleton clusters in both the aspect-f and interpolation-f algorithms. Our first observation is that while the aspect-f algorithm (with singleton clusters) only sometimes outperforms the initial ranking, the interpolation-f (with singleton clusters) does so in almost all of the relevant comparisons.

When comparing the performance of the aspect-f and interpolation-f algorithms when language-model-based nearest-neighbor clusters (nn-LM) are used, as in our original implementation, to that obtained by using singleton clusters (i.e., documents), we clearly see in Table 9 that the former is a much better approach for both algorithms (aspect-f and interpolation-f) than the latter. (Refer to the underlined numbers.) This finding gives further support to the importance of clusters in providing context from the initial list \mathcal{D}_{init} .¹⁰

Finally, the fact that the interpolation-f algorithm almost always outperforms the aspect-f algorithm when singleton clusters are used (as is the case for using nearest-neighbor clusters) shows again the importance of using document-query similarity information in the re-ranking setting.

5 Conclusions

We showed that algorithms that were originally designed for using static query-independent clusters for ranking an entire corpus in response to a query can be adapted to utilize query-specific clusters for effectively re-ranking documents in an initially retrieved list so

¹⁰ A similar conclusion with respect to the superiority of clusters to documents for providing (query-independent) corpus context was made in Kurland et al. (2005).

as to improve precision at the top ranks. The best-performing algorithm that we developed consistently outperforms both the initial document ranking and a state-of-the-art pseudo feedback method. In further exploration we studied the effect of various—both hard and soft—clustering algorithms on the effectiveness of cluster-based smoothing for the re-ranking task, and showed the importance of clusters in providing context from the initial list by comparison to using single documents to this end.

Acknowledgments The author thanks Lillian Lee for many valuable discussions and comments. Part of the work that is described in this paper was done while the author was at Cornell University. The paper is based upon work supported in part by the National Science Foundation under grant no. IIS-0329064 and by a research award from Google. Any opinions, findings, and conclusions or recommendations expressed are those of the author and do not necessarily reflect the views of any sponsoring institutions or the U.S. government.

References

- Abdul-Jaleel, N., Allan, J., Croft, W. B., Diaz, F., Larkey, L., Li, X., Smucker, M. D., & Wade, C. (2004). UMass at TREC 2004—novelty and hard. In *Proceedings of the Thirteenth Text Retrieval Conference (TREC-13)* (pp. 715–725).
- Azzopardi, L., Girolami, M., & van Rijsbergen, K. (2004). Topic based language models for ad hoc information retrieval. In *Proceedings of International Conference on Neural Networks and IEEE International Conference on Fuzzy Systems* (pp. 3281–3286).
- Baliński, J., & Daniłowicz, C. (2005). Re-ranking method based on inter-document distances. *Information Processing and Management*, 41(4), 759–775.
- Blei, D. M., Ng, A. Y., & Jordan, M. I. (2003). Latent dirichlet allocation. *Journal of Machine Learning Research*, 3, 993–1022.
- Buckley, C., Salton, G., Allan, J., & Singhal, A. (1994). Automatic query expansion using SMART: TREC3. In *Proceedings of the Third Text Retrieval Conference (TREC-3)* (pp. 69–80).
- Connell, M., Feng, A., Kumaran, G., Raghavan, H., Shah, C., & Allan, J. (2004). UMass at TDT 2004. TDT2004 System Description.
- Cover, T. M., & Thomas, J. A. (1991). Elements of Information Theory. Wiley series in telecommunications. Wiley-Interscience, New York.
- Croft, W. B. (1980). A model of cluster searching based on classification. *Information Systems*, 5, 189–195.
- Croft, W. B. (1995). What do people want from information retrieval? D-Lib magazine.
- Croft, W. B., & Lafferty, J. (Eds.). (2003). Language modeling for information retrieval. In *No. 13 in Information Retrieval Book Series*. Kluwer.
- Cronen-Townsend, S., Zhou, Y., & Croft, W. B. (2004). A language modeling framework for selective query expansion. Tech. Rep. IR-338, Center for Intelligent Information Retrieval, University of Massachusetts.
- Cutting, D. R., Karger, D. R., Pedersen, J. O., & Tukey, J. W. (June 1992). Scatter/Gather: A cluster-based approach to browsing large document collections. In *15th Annual International SIGIR* (pp. 318–329). Denmark.
- Daniłowicz, C., & Baliński, J. (2000). Document ranking based upon Markov chains. *Information Processing and Management*, 41(4), 759–775.
- Diaz, F. (2005). Regularizing ad hoc retrieval scores. In *Proceedings of the Fourteenth International Conference on Information and Knowledge Management (CIKM)* (pp. 672–679).
- Diaz, F., & Metzler, D. (2006). Improving the estimation of relevance models using large external corpora. In *Proceedings of SIGIR* (pp. 154–161).
- El-Hamdouchi, A., & Willett, P. (1986). Hierarchic document clustering using ward's method. In *Proceedings of SIGIR* (pp. 149–156).
- Griffiths, A., Luchhurst, H. C., & Willett, P. (1986). Using interdocument similarity information in document retrieval systems. *Journal of the American Society for Information Science (JASIS)*, 37(1), 3–11, reprinted in K. S. Jones & P. Willett (Eds.), *Readings in information retrieval* (pp. 365–373), Morgan Kaufmann, 1997.
- Hearst, M. A., & Pedersen, J. O. (1996). Reexamining the cluster hypothesis: Scatter/Gather on retrieval results. In *Proceedings of SIGIR* (pp. 76–84).

- Hofmann, T., & Puzicha, J. (1998). Unsupervised learning from dyadic data. Tech. Rep. TR-98-042, International Computer Science Institute (ICSI).
- Hu, X., Bandhakavi, S., & Zhai, C. (2003). Error analysis of difficult TREC topics. In *Proceedings of SIGIR* (pp. 407–408), poster.
- Kleinberg, J. (1998). Authoritative sources in a hyperlinked environment. In *Proceedings of the 9th ACM-SIAM Symposium on Discrete Algorithms (SODA)* (pp. 668–677), extended version in *Journal of the ACM*, 46, 604–632, 1999.
- Kurland, O. (2006). Inter-document similarities, language models, and ad hoc retrieval. Ph.D. thesis, Cornell University.
- Kurland, O. (2008). The opposite of smoothing: A language model approach to ranking query-specific document clusters. In *Proceedings of SIGIR*.
- Kurland, O., & Domshlak, C. (2008). A rank-aggregation approach to searching for optimal query-specific clusters. In *Proceedings of SIGIR*.
- Kurland, O., & Lee, L. (2004). Corpus structure, language models, and ad hoc information retrieval. In *Proceedings of SIGIR* (pp. 194–201).
- Kurland, O., & Lee, L. (2005). PageRank without hyperlinks: Structural re-ranking using links induced by language models. In *Proceedings of SIGIR* (pp. 306–313).
- Kurland, O., & Lee, L. (2006). Respect my authority! HITS without hyperlinks utilizing cluster-based language models. In *Proceedings of SIGIR* (pp. 83–90).
- Kurland, O., Lee, L., & Domshlak, C. (2005). Better than the real thing? Iterative pseudo-query processing using cluster-based language models. In *Proceedings of SIGIR* (pp. 19–26).
- Lafferty, J. D., & Zhai, C. (2001). Document language models, query models, and risk minimization for information retrieval. In *Proceedings of SIGIR* (pp. 111–119).
- Lavrenko, V., Allan, J., DeGuzman, E., LaFlamme, D., Pollard, V., & Thomas, S. (2002). Relevance models for topic detection and tracking. In *Proceedings of the Human Language Technology Conference (HLT)* (pp. 104–110).
- Lavrenko, V., & Croft, W. B. (2001). Relevance-based language models. In *Proceedings of SIGIR* (pp. 120–127).
- Lavrenko, V., & Croft, W. B. (2003). Relevance models in information retrieval. In W. B. Croft & J. Lafferty (Eds.), *No. 13 in Information Retrieval Book Series* (pp. 11–56).
- Lee, K.-S., Park, Y.-C., & Choi, K.-S. (2001). Re-ranking model based on document clusters. *Information Processing and Management*, 37(1), 1–14.
- Leuski, A. (2001). Evaluating document clustering for interactive information retrieval. In *Proceedings of the Tenth International Conference on Information and Knowledge Management (CIKM)* (pp. 33–40).
- Liu, X., & Croft, W. B. (2004). Cluster-based retrieval using language models. In *Proceedings of SIGIR* (pp. 186–193).
- Liu, X., & Croft, W. B. (2006a). Experiments on retrieval of optimal clusters. Tech. Rep. IR-478, Center for Intelligent Information Retrieval (CIIR), University of Massachusetts.
- Liu, X., & Croft, W. B. (2006b). Representing clusters for retrieval. In *Proceedings of SIGIR* (pp. 671–672), poster.
- Liu, X., & Croft, W. B. (2008). Evaluating text representations for retrieval of the best group of documents. In *Proceedings of ECIR* (pp. 454–462).
- Metzler, D., Diaz, F., Strohman, T., & Croft, W. B. (2005). Using mixtures of relevance models for query expansion. In *Proceedings of the Fourteenth Text Retrieval Conference (TREC)*.
- Ponte, J. M., & Croft, W. B. (1998). A language modeling approach to information retrieval. In *Proceedings of SIGIR* (pp. 275–281).
- Porter, M. F. (1980). An algorithm for suffix stripping. *Program*, 14(3), 130–137, reprinted in K. S. Jones & P. Willett (Eds.), *Readings in information retrieval*, Morgan Kaufmann, 1997.
- Preece, S. E. (1973). Clustering as an output option. In *Proceedings of the American Society for Information Science* (pp. 189–190).
- Ruthven, I., & Lalmas, M. (2003). A survey on the use of relevance feedback for information access systems. *Knowledge Engineering Review*, 18(2), 95–145.
- Shah, C., & Croft, W. B. (2004). Evaluating high accuracy retrieval techniques. In *Proceedings of SIGIR* (pp. 2–9).
- Song, F., & Croft, W. B. (1999). A general language model for information retrieval (poster abstract). In *Proceedings of SIGIR* (pp. 279–280).
- Tao, T., Wang, X., Mei, Q., & Zhai, C. (2006). Language model information retrieval with document expansion. In *Proceedings of HLT/NAACL* (pp. 407–414).
- Tombros, A., Villa, R., & van Rijsbergen, C. (2002). The effectiveness of query-specific hierarchic clustering in information retrieval. *Information Processing and Management*, 38(4), 559–582.
- van Rijsbergen, C. J. (1979). *Information retrieval* (2nd ed.). Butterworths.

- Voorhees, E. M. (2005). Overview of the TREC 2005 robust retrieval task. In *Proceedings of the Fourteenth Text Retrieval Conference (TREC)*.
- Voorhees, E. M., & Harman, D. K. (2005). *TREC: Experiments and evaluation in information retrieval*. The MIT Press.
- Wei, X., & Croft, W. B. (2006). LDA-based document models for ad-hoc retrieval. In *Proceedings of SIGIR* (pp. 178–185).
- Willett, P. (1985). Query specific automatic document classification. *International Forum on Information and Documentation*, 10(2), 28–32.
- Zamir, O., & Etzioni, O. (1998). Web document clustering: a feasibility demonstration. In *Proceedings of SIGIR* (pp. 46–54).
- Zhai, C., & Lafferty, J. D. (2001). A study of smoothing methods for language models applied to ad hoc information retrieval. In *Proceedings of SIGIR* (pp. 334–342).
- Zhang, B., Li, H., Liu, Y., Ji, L., Xi, W., Fan, W., Chen, Z., & Ma, W.-Y. (2005). Improving web search results using affinity graph. In *Proceedings of SIGIR* (pp. 504–511).