Boosting Search Performance Using Query Variations

RODGER BENHAM, RMIT University
JOEL MACKENZIE, RMIT University
ALISTAIR MOFFAT, The University of Melbourne
J. SHANE CULPEPPER, RMIT University

Rank fusion is a powerful technique that allows multiple sources of information to be combined into a single result set. However, to date fusion has not been regarded as being cost-effective in cases where strict perquery efficiency guarantees are required, such as in web search. In this work we propose a novel solution to rank fusion by splitting the computation into two parts – one phase that is carried out offline to generate pre-computed centroid answers for queries with broadly similar information needs, and then a second online phase that uses the corresponding topic centroid to compute a result page for each query. We explore efficiency improvements to classic fusion algorithms whose costs can be amortized as a pre-processing step, and can then be combined with re-ranking approaches to dramatically improve effectiveness in multi-stage retrieval systems with little efficiency overhead at query time. Experimental results using the ClueWeb12B collection and the UQV100 query variations demonstrate that centroid-based approaches allow improved retrieval effectiveness at little or no loss in query throughput or latency, and with reasonable pre-processing requirements. We additionally show that queries that do not match any of the pre-computed clusters can be accurately identified and efficiently processed in our proposed ranking pipeline.

1 INTRODUCTION

Rank fusion is used to combine knowledge from different result sets into a single highly-effective answer page. The fusion can be score-based, in which the retrieval scores of documents are aggregated; or rank-based, in which documents are assigned a weighting based solely on their positions in the separate lists. In both cases, the new top-k result set is derived by re-sorting the documents after aggregate weightings are computed. Vogt and Cottrell [71] describe the effects that allow fusion to produce a more effective response: taking advantage of diversity in document representation (skimming); building consensus among ranked lists (chorus); and catering for differences in quality of rankers (dark horse). Belkin et al. [8] show that a simple unsupervised fusion of all related queries on the same system yields greater effectiveness than fusing one query issued to many better-performing systems. Fusing the output of a one-shot query issued to many IR systems has also received attention - for example, Vogt [69] empirically shows that there is an implicit upper-bound of systems that should be fused before diminishing returns on effectiveness are experienced. Vogt and Cottrell [70] make a case for fusion in web search based on parallelizing "fast but inaccurate IR systems", and then combining the lists to obtain results commensurate with a single high-performance system. Although fusion of results returned from different commercial web search-engines can be effective in practice [37], the approach is generally perceived as being inefficient.

Here we introduce a new approach to efficient online re-ranking which directly leverages data fusion. The key idea is to amortize the cost of fusion, and use a pre-processing phase that computes query cluster centroids. In addition, to efficiently compute offline clusters, we introduce a cost-sensitive fusion technique that employs a single heap to simultaneously evaluate all query variations. Similar queries can be mined or generated [11], then aggregated together using a variety

^{© 2018} Copyright held by the owner/author(s).

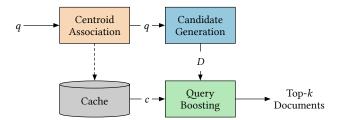


Fig. 1. The user enters a query q, which is both evaluated against the document collection, and used to search a cache of query clusters. If a match is found in the cache, cluster centroid c and the top-k document list D are combined by the query boosting method, yielding a final top-k result page.

of well-known techniques [57, 64, 73]. For reproducibility, the experiments we report here employ the publicly available UQV100 test collection [5].

Figure 1 summarizes the proposed architecture. Pre-computed *centroid rankings* for identified query clusters are held in a searchable cache. When a new query arrives, an association process identifies a matching cluster, and the centroid ranking is retrieved; at the same time, the query is processed to produce a top-k ranking. The centroid ranking and the query ranking are then fused, and a combined top-k ranking computed. The goal is for that fused ranking to represent the combined consensus of all of the queries in the pre-processed cluster, together with the documents specifically identified by the query that was issued.

Given this proposed combination of offline and online processing, we present several methods for combining the centroid ranking and the original query ranking. These include balanced interleaving between the query ranking and the centroid ranking, as occurs in online retrieval experiments [63]; carrying out a weighted CombSUM [34] between the two rankings; and employing a re-ranking approach in which the common-to-both documents are placed at the head of the result page, ordered by their position in the centroid.

Contributions. In particular, in this paper we:

- Investigate the wall-clock time of single-pass rank fusion and the parallel method described by Vogt and Cottrell [70];
- Propose a novel cost-effective single-pass rank fusion technique;
- Describe a novel query fusion architecture that efficiently re-ranks queries in an online setting; and
- Validate our results using the ClueWeb12B corpus and the query variations provided by Bailey et al. [5].

In what follows, Section 2 introduces related work; Section 3 describes the experimental setup; Section 4 explores a number of techniques for fusing multiple query variations *online*; Section 5 shows how the fusion can be computed *offline* and proposes various rankers that can utilize such pre-computed data; Section 6 outlines some of the shortcomings of our investigation and outlines future work; and Section 7 concludes the article.

2 BACKGROUND

2.1 Rank Fusion

It has been known for many years that combining the ranked retrieval outputs of query variations representing the same information need can improve retrieval effectiveness. Belkin et al. [8] showed a representative sample of topic descriptions to ten experienced searchers, generating a pool of five Boolean queries for each topic. The outputs of the query variations were combined using an

unweighted sum of retrieval scores. This aggregation of ranked-retrieval scores was later named CombSUM by Fox and Shaw [34] in their investigation into different rank fusion techniques. Buckley and Walz [15] showed that the same technique could be applied to ranked retrieval. Bailey et al. [6] studied these phenomena further in their exploration of fusion and query variation consistency. Several recent follow-on studies on query variations and rank fusion have shown that combining multiple individual rankings consistently boosts effectiveness [10–13].

Kozorovitsky and Kurland [43] showed that inter-document similarities can be used in a fusion framework to reward documents that are similar to the head of the result list. Other work has explored the role of fusion in diversification [46]; automatic generation of query perturbations [76]; the relationship between fusion and clustering [43]; and boosting tail query performance using supervised rank fusion [39]. No previous studies have explored the resource implications of these techniques in a large-scale search environment.

The use of supervised rank fusion with query variations has also been examined. Sheldon et al. [65] describe a supervised data fusion method named LambdaMerge, which optimizes a retrieval effectiveness metric based on fusion over user query reformulations and a wide range of document features. Lee et al. [45] extend the LambdaMerge framework from data fusion to collection fusion, where query-list features in a collection are averaged and used as query-vertical features. These approaches, as well as those of Huo et al. [39], are closely aligned with our own. Other fusion techniques are also possible [2, 44].

The relationship between supervised fusion and learning-to-rank is an important issue in its own right, but one which is orthogonal to this work. Our approach here does not require supervised learning to achieve competitive results against strong baselines that were selected based on their well-established ability to efficiently and effectively rank documents in a web search environment. We will explore supervised approaches in future work.

Xue and Croft [76] show that query perturbations can be generated using automatic methods with some success. Sheldon et al. [65] describe a supervised data fusion method named LambdaMerge, which optimizes a retrieval effectiveness metric based on fusion over user query reformulations and a wide range of document features. Lee et al. [45] extend the LambdaMerge framework from data fusion to collection fusion, where query-list features in a collection are averaged to be utilized as query-vertical features. Liang et al. [46] also consider fusion, showing that it helps with diversification when queries have multiple interpretations.

2.2 Efficient Index Traversal

The most commonly used structure for top-k document retrieval is the *inverted index*. Each unique term t is represented by a *postings list*, a sequence of document identifier (docid)/term-frequency $(d_{t,i}, f_{t,i})$ pairs, one for every document in which term t appears, where $d_{t,i}$ is the docid of the ith document containing t, and $f_{t,i}$ is the corresponding within-document term frequency. Inverted indexes provide efficient and scalable access to the necessary statistics for document ranking [78]. When a query is received, the postings lists associated with the query terms are fetched, and combined to rank and return the top-k documents.

The way in which the postings lists are iterated, known as the *index traversal strategy*, has a large impact on efficiency, and different postings layouts are amenable to different traversal strategies. Here we focus on the popular *document-ordered* index layout, which is most commonly used with *Document-at-a-Time* (DAAT) query processing strategies such as WAND [14] and more recent block-based variants (BMW) [19, 28, 31, 55]. These approaches tend to be more efficient than are DAAT MAXSCORE and *term-at-a-time* (TAAT) approaches [67, 68], particularly for short queries, the most common scenario in web search. However, for long queries or large candidate sets, the case is less clear-cut [23, 33, 54]; moreover, fusion over query variations often leads to very long queries.

That is, both WAND and MAXSCORE have advantages as well as disadvantages depending on the length of the query, the number of postings to process, the term selectivity, and a range of other factors.

2.3 Batch Query Processing

Batch processing is a commonly used technique in the database community [21], but is exploited less often in web search. The majority of queries received by a search system are from interactive users and must be processed efficiently, with query latency an important contributor to user satisfaction [4]. Ding et al. [30] exploit the existence of a subclass of queries that can be processed in batches to minimize the cost of search. These queries include operations such as cache updates, internal testing, and index mining. Ding et al. also show how batches of queries can be processed in the context of large-scale search systems, optimizing I/O and CPU costs. Their key observation is that if queries can be batched, then reordering allows common intersections to be cached, reducing net evaluation cost. The task of pre-computing an inverted index is also, of course, a batch-processing operation. It is the existence of a pre-computed index that allows queries to be resolved within millisecond response times over massive document collections, and has allowed web search to become the critically important tool that it now is.

2.4 Caching for Large-Scale Search

In order to consistently meet *service level agreements* (SLAs), search engines must of necessity avoid redundant computation. Caching is a simple and common approach for increasing query throughput at the cost of additional space consumption. Caches can be deployed at many levels of storage, including in-memory or on-disk [3, 72]. In IR there are two major approaches to caching. *List caching* involves storing commonly accessed postings lists in fast-access memory [18, 72]. For example, if the postings lists making up the index are stored on a SSD, a list caching strategy may opt to keep to *n* most accessed postings in main-memory. Alternatively, *result caching* involves storing a query along with the relevant results (or some proxy thereof) that were returned for the query [32, 36]. In practice, both list and result caching are useful, and are deployed in tandem [3, 72]. Most caches utilize historical data such as static query logs or sliding windows of recent queries to build models of *what* to cache, and *when* to cache it, and can also be personalized on a per-user basis [50].

The architecture proposed in Figure 1 provides another form of caching to allow subsequent operations to be made more efficient. We explain that architecture more fully in Section 5.

3 METHODOLOGY

Before providing details of the new techniques in Section 4 and 5, we first describe the experimental framework that is employed.

3.1 Hardware and Software

Our experiments are conducted on an idle Red Hat Enterprise Linux Server with 256 GiB of RAM and two Intel Xeon E5-2690 v3 CPUs, each with 12 physical cores. All algorithms were implemented with C++11 and compiled with GCC 6.3.1 using the highest optimization settings. Where multi-threading was used, up to 48 threads were spawned using the C++ STL threading libraries. All algorithms were implemented as components within the state-of-the-art VBMW code-base described by Mallia et al. [55]¹, and in support of reproducibility, are also made available for others (see Section 7).

¹https://github.com/rossanoventurini/Variable-BMW

Table 1. The ClueWeb12B and UQV100 resources used. Note that the queries were stopped and Krovetz stemmed, reducing the number of distinct queries compared to that reported by Bailey et al. [5].

Documents	52,343,021
Topics	100
Total queries	10,835
Unique queries	4,175
Mean unique queries per topic	41.75
Hold-out queries	500

3.2 Collections and Indexes

We conduct our experiments across the 52 million document ClueWeb12B corpus and employ the UQV100 query collection [5] and its 100 single-faceted topics derived from the multi-faceted TREC 2013 and 2014 Web Tracks. The collection contains 10,835 query variations, sourced from crowd-workers who were presented with a narrative "backstory" for each topic, and asked to formulate a query in response. For more details in regard to the size and homogeneity of the clusters we refer the reader to Bailey et al. [5]. Moffat [58] has also explored properties of this collection; and Moffat et al. [59] discuss the wider implications of query variations.

To support the required experiments, we split the 10,835 UQV100 queries into two sets: a *training* set, used to build the query variation clusters, and a *testing*, or *hold-out* set, used to measure the final performance of the proposed approaches. The hold-out set was created by selecting five unique query variants per topic (that is, queries appearing only a single time in the UQV100 set), yielding a set of 500 queries across the 100 topics. Each hold-out query was drawn randomly from the corresponding topic's single-instance variations, without replacement. This hold-out approach differs from the simple hold-out method described by Fuhr [35], as all baselines and new techniques are evaluated against the total universe of topics, with at least five query impressions per-topic. The arrangement also avoids the limitations observed in other train-test splits in a single-query-per-topic evaluation scenario, where results are biased by the topic-effect of the generated split. The training set here represents the most commonly seen query variations for each topic, and hence can be regarded as being representative of what could be mined from logs in a production system. Table 1 summarizes the situation. Note that the number of distinct queries is less than in the underlying UQV collection [5] because of our use of a Krovetz stemmer and a stop list.

We used Indri 5.11^2 to index the collection, and then converted the inverted index into the format expected by the VBMW code-base. Before building the VBMW index, we reordered the docid space using the recursive graph bisection³ approach of Dhulipala et al. [27], as it has been shown to substantially improve index compression. The average block size of our VBMW index is approximately 40 integers per block, a result of binary searching for the parameter λ as discussed by Mallia et al. [55]. The final index was compressed using the Partitioned Elias-Fano mechanism [62].

3.3 Evaluation Metrics

We use two metrics to evaluate effectiveness: the recall-based NDCG approach [40] at a fixed cutoff depth of 10; and the utility-based RBP method [60] applied to each full ranking, with persistence $\phi = 0.8$ and hence an expected viewing depth of five. These metrics and cutoffs were selected based on the judgment depth of the UQV100 collection [48], and result in relatively low score uncertainty.

²https://www.lemurproject.org/indri.php

³https://github.com/mpetri/recursive_graph_bisection

Significance is computed using the Bonferroni-corrected paired t-test, and is denoted by \dagger for p < 0.05, and by \ddagger for p < 0.001. Significance was always tested with respect to the strongest available baseline.

4 REAL-TIME FUSION OF QUERY VARIANTS

This section considers the cost of computing fused answer rankings at query time, assuming (through this section) that it is to be done without the use of the pre-computed centroids that were foreshadowed in Figure 1.

4.1 Efficiently Processing Variations

Given a cluster of queries, including the query just entered by the user, the goal is to process them all, fuse their results, and return a single SERP (search engine results page), all the while noting that web-search systems typically impose strict per-query resource budgets [26, 41, 77].

Parallel Fusion (PF). The simplest approach is to spawn a process for each unique query variation in the cluster, and execute the queries in parallel. Once all threads have returned their top-k results, a rank-fusion algorithm is applied, to assemble the final SERP. This approach is viable provided there are sufficient CPU cores available, and requires that each thread operate to the same response-time requirement as the original query. If latency as a critical component of the SLA, queries that run longer than a fixed time threshold could be abandoned and not considered for the fusion process [41, 77], thereby sacrificing effectiveness to stay within the resource constraint. We measured three variants of the parallel fusion approach, denoted "PF-a", where a is a query processing strategy, one of VBMW, WAND, and MAXSCORE.

Single Pass Daat. A drawback of the parallel approach is that many similar queries are processed concurrently, and hence that some of the corresponding postings lists are processed many times, without any commonality being exploited. An alternative is to perform all scoring operations in a single Daat pass across the inverted index, concurrently building a top-k heap for each unique query variant. That is, an empty top-k heap is constructed for each unique query variant, and the postings lists for all terms are iterated in parallel, selecting as the *pivot* the minimum document ID across all of the cursors. At each processing step, all postings lists are advanced to align with the pivot, with variables tracking the current set of scores of the pivot document with respect to the terms appearing in each query. Once all aligned postings have been processed, the document scores are checked against the corresponding heaps, each of which is updated if necessary. Finally, when all postings cursors are exhausted, the set of heaps contain the top-k results for the set of query variations, and can be fused to create the required single SERP. We refer to this approach as "SP-Exhaustive".

Single Pass CombSUM. We now propose an efficient single-pass approach for computing the CombSUM [34] rank fusion score for a set of query variations, allowing improved efficiency through *dynamic pruning*. Given a set of ranked lists of documents, and a positive numeric score for each document in each list, the CombSUM score for a document d is the sum of the scores of d, computed over its appearances in the ranked lists. If there are ℓ lists, L_1 to L_ℓ , and the score of some document d in the i th of the lists is given by $s_{i,d}$ (with $s_{i,d} \equiv 0$ if $d \notin L_i$), then

CombSUM(d) =
$$\sum_{i=1}^{\ell} s_{i,d}$$
.

Now consider each of the component scores $s_{i,d}$, and the query q_i that led to it. If the scoring computation is an *additive* one, then

$$s_{i,d} = \sum_{t \in q_i} F(t,d),$$

where F(t, d) is the term-document score contribution associated with the term t in the document d according to the chosen retrieval model. Taking these together gives

CombSUM(d) =
$$\sum_{i=1}^{\ell} \left(\sum_{t \in q_i} F(t, d) \right).$$
 (1)

Now define $n_t \equiv |\{q_i \mid 1 \le i \le \ell \land t \in q_i\}|$, the number of input queries containing term t; and $Q \equiv \bigcup_{1 \le i \le \ell} q_i$, the union of the queries. Equation 1 can then be rewritten as

CombSUM(d) =
$$\sum_{t \in Q} n_t \cdot F(t, d)$$
, (2)

making it clear that for additive similarity scoring mechanisms, the CombSUM score for a set of query variations can be computed by forming the union Q of the queries, counting term frequencies n_t across the variations, and then evaluating a single "super query" against the index of the collection and taking a linear sum of the individual contributions F(t,d). This requires that all of the scores are positive. Similarity models that are not additive, or that yield negative scores, may not be used in this way. Note also that the usual CombSUM process of scaling the document scores into the range $[0\dots 1]$ would prevent the simplification shown in Equations 1 and 2. In the following experiments, we use the BM25 similarity model, and do not normalize the individual ranking scores.

In order to employ Equation 2, and allow dynamic pruning techniques to be safely applied to the super query, the index traversal process must be slightly modified, with the upper-bound scores, U_t , also multiplied by n_t . For the block-based VBMW approach, we must also supply the n_t multiplier to each block-max score, $U_{b,t}$, on-the-fly. Query processing does not differ in any other way. This approach is generalizable to all safe-to-k dynamic pruning traversal strategies; and hence we again test three variations, denoted "SP-CS-a", with a one of VBMW, WAND, and MAXSCORE.

4.2 Experiment: Real-Time Fusion

To test these approaches, we take all query variations for each topic, and measure the cost of computing a fused top-100 result list. This involves combining 42 query variations per topic on average, each computed to depth k=1,000, and then (except in the case of the SP-CS approaches) fusing the results. Three indicators are reported: the number of CPU cycles consumed; the number of postings scored; and response latency. In the case of the parallel approaches, the first two are summed over all threads. Note that the CPU measurements ignore the slight processing overhead generated by the forking and locking activities inherent in parallel execution.

Figure 2 shows the results. The first (top left) pane shows the total CPU time required. The three SP-CS methods are the most efficient in terms of processing cost, with the SP-CS-MAXSCORE approach slightly better than the other two. In the second pane (top right), the SP-CS-MAXSCORE implementation processes more postings on average than either the SP-CS-WAND and SP-CS-VBMW approaches – the latter two reduce the number of postings, at the cost of more non-posting processing. Finally, the third pane (bottom left) shows elapsed wall-clock time. The three PF approaches are the fastest, due to their extensive use of parallelism. Even though each query's latency is dictated by the slowest-running variation, execution on average is fast; and the use of suitable aggregation policies [77] can improve the *tail-latency* [53] of such approaches. Note, however, that this speed

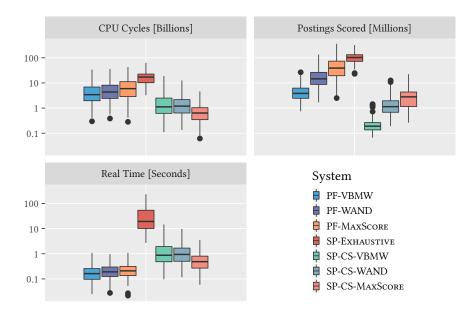


Fig. 2. Efficiency of rank-fusion algorithms, assuming that the starting point is a set of query variations. The panes show total cost in CPU cycles; total cost in terms of postings processed; and query latency.

comes at a resource cost, as shown in the first pane. If latency and overall workload are both critical concerns, parallelism should instead be added to the SP-CS versions, by splitting the collection across the pool of available processors.

4.3 How Many Queries Should Be Fused?

Another way of reducing processing costs is to fuse fewer queries. To determine the impact that the number of variants has on processing costs, random samples (with replacement) were drawn from each UQV100 query cluster. Those samples were then executed and fused, recording the cost in CPU cycles and the effectiveness of the final SERP. The selection process was carried out incrementally, with one variation drawn and measured, then a second added to it and measured, and so on; with that entire sequence repeated ten times, and the values recorded being the averages over those ten runs. Table 2 and Figure 3 show the results.

Table 2 makes it clear that adding variants increases effectiveness, but that the gains diminish as more variants are added. This outcome is at least partially a consequence of the methodology used, since "with replacement" means that there is an increasing probability of a previously-selected query being drawn again. Bailey et al. [6] also observe that, in general, adding more distinct variants improve the effectiveness of rank fusion, irrespective of metric or fusion mechanism.

Figure 3 shows that SP-CS-MAXSCORE retains its computational advantage across the range of variation set sizes, and that all of the SP-CS approaches are cheaper than the SP-EXHAUSTIVE and PF methods. The reason for this advantage is that the SP-CS approaches are largely determined by the number of unique *terms* that are required for processing, whereas the other methods are more closely tied to the number of unique *queries*. Since query variations often contain many similar terms, the SP-CS approaches scale better.

Table 2. Fusion effectiveness as the number of query variants is increased. Variants are selected from the query clusters at random, with replacement; measured as averages over a set of ten such sequences. The values in parentheses are RBP residuals, recording the maximum extent of the RBP score uncertainty.

Num. variants	NDCG@10	RBP $\phi = 0.8$
1	0.182	0.426 (+0.067)
2	0.210	0.469 (+0.081)
5	0.236	0.514 (+0.043)
10	0.254	0.537 (+0.028)
20	0.256	0.540 (+0.021)
50	0.261	0.542 (+0.018)
100	0.262	0.550 (+0.017)

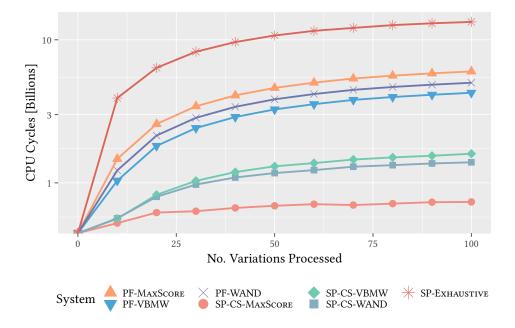


Fig. 3. Average per-topic cost in CPU cycles of approaches for generating fused rankings, as a function of the number of query variations being fused. Sampling is random from each cluster, with replacement, averaged over ten independent sequences. All of these methods obtain the same effectiveness, as listed in Table 2. Note the logarithmic vertical scale.

4.4 Discussion

Three approaches for efficiently fusing a set of query variations have been described and measured. The SP-EXHAUSTIVE and PF approaches can be used with any rank-fusion algorithm, but are more costly than the SP-CS approaches, which are based on CombSUM and additive functions such as BM25. The PF approaches have the lowest latency, but use nearly an order of magnitude more resources than the implementation of the SP-CS methods that was measured here. As a result, the SP-CS approaches are more scalable with respect to the number of query variations, and, as already noted, can be parallelized to reduce their latency. That is, rank fusion across query variations

can be computed such that latency is small enough for online use, but only if substantial total computational resources are available. Better techniques are required if the overall resource cost must also be managed carefully.

Although this section has focused on fusion using the cost-effective CombSUM method described above, there are other fusion techniques that can also be applied, some of which are easier than others to transform into single-pass implementations. Benham and Culpepper [10] compare the effectiveness of CombSUM with six other unsupervised fusion techniques on the ClueWeb12B corpus using the UQV100 query set, and demonstrate that CombSUM is competitive with the CombMNZ [34], RBC [6], and RRF [22] approaches. The next section shows how even more complex combinations of fusion over query variations and systems can produce highly effective search results.

5 USING PRE-COMPUTED CENTROIDS

The previous section introduced a reduced-cost query fusion technique and showed that it is more efficient than computing a separate ranked list for each query variation. Even so, if aggressive service level agreements governing response time and throughput must be complied with, on-the-fly query fusion may not be a viable approach. The key issue then becomes: is it possible to compute fused lists in an offline manner, and use those cached intermediate results to boost online query performance? With this question in mind, we now turn to the arrangement that was sketched in Figure 1, in which pre-generated query centroids are stored, and used during online query processing to adjust the ranking generated for each individual query as it is executed.

5.1 Computing Query Centroids

Fused centroid lists can be pre-computed using the methods described in Section 4 at the same time as pseudo-documents based on the union of the terms in each cluster are constructed and themselves indexed. Since this is a pre-computation, total resource requirement is the appropriate cost measure, and not latency; and with the cost equation further moderated by the expectation that the pre-processing time can be amortized over multiple subsequent queries that refer to that cluster. Hence, any desired fusion technique can be used, with no restriction to BM25 and/or CombSUM.

If query centroid data in the form of indexed pseudo-documents and consensus fused rankings are stored in main memory or on SSD, it can be searched and retrieved quickly as queries are processed. In the experiments described shortly, the fused query consensus rankings (the cluster ranking *centroids*) are stored in main memory, using a doubly-linked circular list that preserves the ranked order of the fused set; with document identifiers also maintained in an O(1) average-time hash table. With these structures, the UQV100 data required 49.1 kiB per topic on average to store a ranked list containing k = 1000 documents (centroid lists were always computed and stored to a length of 1,000 documents). That space requirement can be reduced to 11.5 kiB per topic if the subsequent fusion is restricted to methods that are rank-based, where knowledge of document identifiers and ranks alone is sufficient without document scores.

5.2 Boosting Effectiveness Using Centroids

Three different approaches to joining the centroid ranking and a query ranking – with the goal of boosting overall effectiveness – were explored and measured. While not an exhaustive list of possibilities, these three methods demonstrate that improvements in risk-sensitivity and retrieval effectiveness are achievable using relatively simple and inexpensive techniques.

Plain Interleaving. The first approach is inspired by online retrieval experiments in which results from two different systems are presented as a single list. In a balanced interleave, elements are

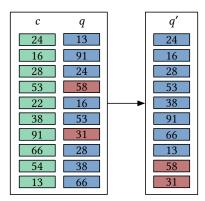


Fig. 4. Reference list re-ranking, where c is the k-document pre-computed centroid ranking, q is the k-document run generated for the user query, and q' is their k-document fused result, in this example with k=10 for both of the lists. Documents in green are from the centroid list, documents in blue are common to it and the query list, and the documents in red are the ones from q that were not in c.

chosen in alternating fashion from two results, with the first chosen randomly. Instead, we bias towards the query centroid and always take it first, assuming that it aggregates more information about the information need, and is more likely to have high-ranking relevant documents. We denote this approach as "Interleave", and at each step the next highest ranked document not already included in the output list is taken from one of the two lists being joined in a strictly alternating pattern, and added to the fused run. A possible benefit of the interleaving model is its connection to A/B testing, where click logs may be helpful in deciding whether to preference should be given to the user query or the fused result centroid.

Linear Combination. Another option is to adopt the approach proposed by Vogt and Cottrell [71], and compute a per-document weighted sum of the min-max re-scaled centroid set and the re-scaled user query answer set, then use those scores as a descending-order sort key to form the SERP. To implement this approach we applied a weight of δ to the query centroid and $(1-\delta)$ to the user query, and set δ to 0.5 as a starting point. An exploration of the impact of varying δ on retrieval effectiveness and risk-sensitivity is carried out below; when $\delta > 0.5$, the method approaches the bias exhibited in the Interleave method. Unlike the Interleave method, a linear combination does not guarantee an equal contribution from either the cluster centroid or the documents retrieved in response to the actual query. This is because inclusion in the final results list is sensitive to both the union of the two sets of documents, and also to the scores of those documents. We denote this approach as "LC".

Reference List Re-Ranking. This approach intersects the query's ranking and the centroid ranking, adopting the ordering supplied by the centroid run. Any remaining documents from the query's run are added after the tail of the intersection set, in the same order as they appear in the fused list. Both the centroid ranking order and the query's ranking order are respected, but with the documents that were in both listed first. We denote this approach as "Ref-Reorder". Figure 4 provides an example in which two runs of length k=10 are joined, favoring the documents that appear in both, and accepting the ordering of the consensus ranking. Note that in the experiments described shortly the centroid list was always 1,000 documents long, and only the query's ranking length was altered during the experimentation.

5.3 Centroid Construction Approach

For the following experiments, we construct query centroids by fusing the training, or "held-in", set of queries for each cluster using . This allows us to establish how much improvement can be gained through building clusters using the proposed cost-effective single-pass query centroid generation scheme introduced in Section 4.1. Note, however, that this is not a necessary part of the arrangement, and any construction approach can be used to generate centroids, not just fusing BM25 results using CombSUM. Other options are discussed in Section 6.

5.4 Evaluation

The baselines used in this section include two proximity-based algorithms and one learning-to-rank model.

Term Dependency Models. We use two term-dependency approaches that represent strong baselines. The first is the L2P approach proposed by Lu et al. [49], a bigram ranking model that linearly combines the BM25 score of a document with sequential bigrams. Unlike other term dependency models, L2P does not require global statistics for term dependencies, greatly improving overall efficiency. We also use a field-weighted variation of the sequential dependency model (SDM) [56] that operates across the document body, document title, and inlink text. For both of these approaches, the ClueWeb09B collection was used in conjunction with the TREC 2009–2012 Web Track topics to tune the parameters.

Learning-to-Rank. As a learning-to-rank baseline, denoted "LTR", the LightGBM⁴ framework was used to build a LambdaRank model [16]. Instead of training on just the hold-out set, we train our model on both the topics in the hold-out set as well as the most common five variants from each topic. This additional exposure leads to a more robust model. We used ten-fold cross-validation to train and test the LtR model across a large set of features, including approximately 150 query-dependent unigram, 150 query-dependent bigram, 15 static document priors, 100 document-dependent unigram, and 5 document-dependent bigram features. The features were generated using the publicly available system described by Chen et al. [20]. We performed an ablation study of the features produced by this system, and used LightGBM for re-ranking. An exhaustive description of features and tools to reproduce this baseline will be released after this paper is accepted. We refer the interested reader to Liu [47] and to Macdonald et al. [51, 52] for further information on building a competitive LtR system.

Risk Sensitivity. In addition to effectiveness, we consider the risk-sensitivity of all proposed methods [29]. Although improved average effectiveness is desirable, this may not be indicative that the centroid-based boosting algorithm is generalizing well. Risk-sensitive retrieval metrics show that a method is improving the quality of user result pages without making them worse. Wins, ties, and losses (W/T/L) are shown for each of the approaches and baselines with respect to the BM25 run on the held-out user query. Ties are defined to occur if the experimental run score is more than 10% different to the BM25 baseline score. Dinçer et al. [29] introduced a risk-sensitive retrieval metric named TRisk – a studentized URisk which accepts an evaluation metric and a baseline run. It is computed by providing a run as input, taking the sum of all improvements over all topics, less the sum of all losses with a linearly-scaled α value to vary the impact of losses. TRisk scores less than –2 indicate that the experimental run exhibits statistically significant harm compared to the baseline, and greater than +2 indicates an improvement in risk-sensitivity with statistical confidence. We use these two measures to evaluate risk-sensitivity, and set the parameter

⁴https://github.com/Microsoft/LightGBM

Table 3. Mean query boosting time and total query execution time (milliseconds per query) and percentage
overhead due to boosting when retrieving the top k documents for user queries computed with BM25. The
query centroid rankings were always 1,000 documents long.

Length	Method	Boosting	Total	Overhead
k = 10	Interleave	1.019	23.479	+4.5%
	LC	1.187	23.647	+5.3%
	Ref-Reorder	0.092	22.552	+0.4%
k = 100	Interleave	1.118	37.126	+3.1%
	LC	1.264	37.272	+3.5%
	Ref-Reorder	0.147	36.155	+0.4%
k = 1,000	Interleave	2.050	68.857	+3.1%
	LC	1.999	68.807	+3.0%
	Ref-Reorder	0.548	67.355	+0.8%

 α = 3, assessing score decreases to be four times more damaging than the benefit derived from any numerically identical score increases.

5.5 Results

The methods introduced in the previous section are compared to baselines in the context of a performance-sensitive query environment. Three key usability aspects are monitored – efficiency, effectiveness, and risk-sensitivity.

Efficiency. Table 3 provides boosting times for BM25 queries using the three proposed methods, as well as the end-to-end time and the percentage overhead relative the original BM25 query. The Interleave and LC approaches show similar performance, and when k=1,000, both techniques take an average of 2 milliseconds to complete, a 3% overhead on executing the BOW query using BM25. The Ref-Reorder approach is approximately four times faster, taking 0.5 milliseconds when k=1,000. In all cases the end-to-end latency from query submission to final top-k is dominated by the initial query BM25 stage, which is efficient in practice, and none of the boosting approaches add more than 6% overhead to the cost of generating the input query ranking.

Effectiveness. Table 4 lists effectiveness scores for the three centroid-based fusion methods, and compares them to the four baseline approaches, using NDCG@10 and RBP $\phi=0.8$ as evaluation metrics, and averaging across the 500 single-instance held-out queries. The number of wins, ties, and losses compared to a BM25 evaluation of the same 500 queries is also reported, together with TRisk $_{\alpha=3}$, which penalizes score losses by a factor of four compared to score gains. Recall that TRisk less than -2 and greater than +2 imply statistical significance.

Of the methods tested, the Ref-Reorder mechanism is the most effective, with an NDCG score of 0.243 and RBP score of 0.527 (and in the latter case with a small residual, indicating that the judgments are a good fit for the runs that were scored). This approach also has the lowest risk-sensitivity among all centroid-based boosting methods, with only 93 of the queries performing worse than the user query alone for NDCG (101 for RBP); note, however, that these results are contingent on the correct centroid having been identified, an assumption that is explored below.

The $TRisk_{\alpha=3}$ values in Table 4 provide a somewhat different picture, and suggest that Interleave boosting is the most desirable if losses relative to the baseline BM25 approach are penalized more

Table 4. Effectiveness of online fusion measured using two effectiveness metrics, compared to three baselines. Significance is measured with respect to LTR, the strongest of those baselines. The scores listed in the "Mean" column are averages across the 500 held-out query variations, five for each of the 100 topics; the "W/T/L" numbers are the respective counts of those 500 queries for which the method exceeds the BM25 baseline by more than 10%, is within 10% of the baseline, and is less than 10% of the baseline.

Metric	Method	Mean	W/T/L	$TRisk_{\alpha=3}$
NDCG@10	User Query (BM25)	0.170^{\ddagger}	_	_
	L2p	0.169^{\ddagger}	160/236/104	-6.992
	SDM+FIELDS	0.180^{\dagger}	247/76/177	-6.801
	LтR	0.200	281/53/166	-4.640
	Interleave	0.223^\dagger	305/108/87	3.829
	LC	0.231^{\ddagger}	322/94/ 84	3.614
	Ref-Reorder	0.243^{\ddagger}	345 /62/93	0.812
$RBP \ \phi = 0.8$	User Query (BM25)	$0.401 \ (+0.088)^{\dagger}$	_	_
	L2p	$0.400 \; (+0.089)^{\dagger}$	136/264/100	-8.457
	SDM+FIELDS	0.417 (+0.153)	215/117/168	-7.557
	LтR	0.435 (+0.205)	258/71/171	-6.790
	Interleave	$0.486 \ (+0.037)^{\ddagger}$	287/150/ 63	4.178
	LC	0.504 (+0.033)‡	290/130/80	3.762
	Ref-Reorder	0.527 (+0.053) [‡]	304 /95/101	1.454

than gains are rewarded, even though when measured by average per-query scores it is the least useful boosting approach.

The LC approach listed in Table 4 employed a parameter δ of 0.5, placing equal weight on the two rankings being combined. Figure 5 shows that as δ is varied the linear combination of the scores can be adjusted to counterbalance risk-sensitivity and effectiveness. For example, shifting to a weighting based 70% on the query-centroid result set and 30% on the user query set retrieved by BM25 yields an improved risk-reward trade-off compared to evenly weighting both lists. Note that the weightings mentioned in Figure 5 are not intended to be prescriptive. In order to test how generalizable these weights are, query variation collections formed using a similar methodology to UQV100 would need to become available to test on other collections.

Figure 6 provides another perspective on these issues, and shows the extent to which NDCG@10 shifts up and down compared to the baseline BM25 ranking on a query-by-query basis. The top three panels represent effectiveness profiles for the three other baseline systems, while the lower three panels demonstrate the profiles of the three boosting techniques described here. The same decreasing-score ordering of queries is used in each of the six panes, with the BM25 scores shown as black dots and the corresponding scores for that method shown as colored dots, and with (in all of the panes) the five other scores for each query shown as light grey background points. Of the three baseline systems, L2p follows the BM25 line closely, as it is derived directly from BM25 with a small additional weight attributed to the sequential dependency scores. The SDM+FIELDs approach uses a language model and does not always agree with BM25, but remains a fair comparison. The LTR effectiveness with respect to BM25 exhibits higher variance the other two baselines, presumably since optimization over many features can induce bigger wins, but can also bigger losses too. For the three boosting methods, Interleave and LC are operating with similar risk-sensitivity, as evidenced by their similar TRisk scores in Table 4. Ref-Reorder wins most consistently, but is also

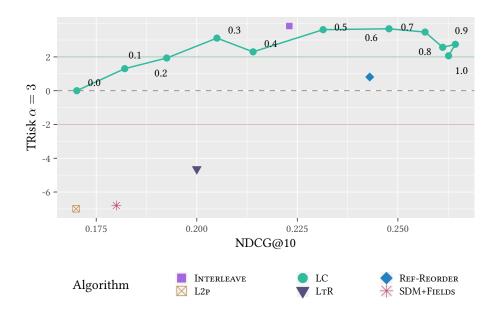


Fig. 5. Varying the weighting of the linear combination (method LC) between the query centroid and the user's original query can reduce risk and improve retrieval effectiveness. The numbers adjacent to the LC line indicate the corresponding values of δ ; effectiveness tends to be greater when $\delta > 0.5$ and more emphasis is given to the cluster centroid run. (Recall that TRisk scores less than -2 and greater than +2 indicate statistical significance.)

vulnerable to quite notable score degradations compared to Interleave and LC, and explains why its TRisk score is lower. A similar pattern of performance was observed when the same plots were generated for RBP.

5.6 Online Re-Ranking using Real Queries

We now turn our attention to a more fundamental question. Using pre-computed clusters appears to be a powerful approach if every query can be neatly categorized against a previously-identified information need. But in a real system, a reasonable percentage of incoming queries can be expected to differ from all known clusters, or worse, to be matched against the wrong cluster. How do we handle queries that do not match any cluster, and minimize the impact on effectiveness of failed query matchings?

Centroid Identification. We first ask whether incoming queries can be reliably mapped to an existing query cluster. Wen et al. [74] showed that combining query keywords and cross-reference similarity using document hierarchies from click-through data can give a precision and recall of approximately 95% in an experiment involving 20,000 query clusters. A 5% error rate is encouraging, given that precision can be traded against recall, to guard against the potential impact of incorrect cluster matches.

To replicate this evaluation, we built an index of *pseudo-documents*, each containing the union of all query terms belonging to one of the UQV100 topics. To this set of 100 documents we added 10,000 synthetically-formed query clusters using the anchor text of inlinks to Wikipedia articles on

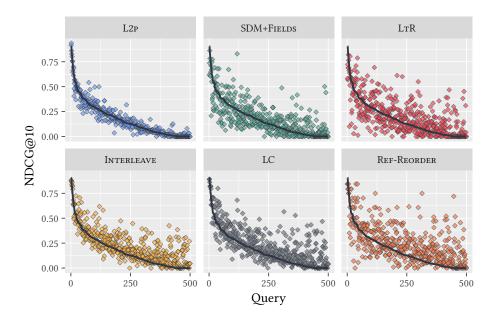


Fig. 6. Difference between NDCG@10 effectiveness score for the baseline BM25 query and for six other evaluation approaches. The 500 held-out queries are ordered by their BM25 scores (the black line of dots in each graph) and the corresponding scores from each of the different evaluation approaches are shown, one set in each pane.

the ClueWeb09 corpus, adopting the approach of Dang and Croft [25], who showed that anchor text can be used as a reliable substitute for user queries. Each pseudo-document contains the union of the query terms (rather than allowing duplicates) so as to not bias cluster selection. Upon receiving a new query, the pseudo-document index is searched using BM25 to find the top-scoring pseudo-document. A range of policies for forming synthetic centroids from inlink data were employed to aid with pre-processing, sanitization, and diversity; code for these steps and a manually curated stop-list are available in our codebase (see Section 7). Across the 10,000 clusters, 221,989 queries were used (41,624 unique); the largest and smallest clusters contained 442 queries and 11 queries respectively. Finally, to test the effect of increasing the number of "distractor" clusters, 0, 100, 1,000, 5,000 and 10,000 were added to the UQV100 clusters, and the success rates for the five held-out UQV queries per topic were computed, where success consisted of selecting the correct UQV100 cluster in the presence of the distractors. These success rates were 97%, 97%, 94%, 91%, and 89% respectively, broadly in line with the results of Wen et al. [74].

Incorrect Centroid Association. The results in Table 4 assume that the centroid identification technique is perfect, which is an unlikely scenario. To quantify the effect of incorrect matchings, we ran ten trials in which each query was assigned to the wrong cluster with varying probabilities, so that the failure profiles of our online combination approaches can be compared in the presence of more realistic centroid matching. We denote the error rate as ϵ and explore values of $\epsilon=0.05$ (with 5% of queries assigned the wrong cluster, the rate attained by Wen et al. [74]) and a somewhat pessimistic $\epsilon=0.2$ (with 20% of queries being misassigned). The latter is worse than the rate observed in the experiments with 10,000 distractor centroids reported above.

Table 5. Effectiveness of online fusion approaches, where ϵ represents the rate in which incorrect clusters are selected. Significance is measured with respect to LTR, the strongest of the four baselines that were employed. Return Cluster Centroid (RCC) is tabulated with a perfect cluster matching rate of $\epsilon=0.0$ as an initial reference point. The scores listed in the "Mean" column are averages across the 500 held-out query variations, five for each of the 100 topics; the "W/T/L" numbers are the respective counts of those 500 queries for which the method exceeds the BM25 baseline by more than 10%, is within 10% of the baseline, and is less than 10% of the baseline.

Error Rate	Method	Mean	W/T/L	TRisk $_{\alpha=3}$
		NDCG@10		
$\epsilon = 0.0$	RCC	0.263 [‡]	375/40/85	3.075
$\epsilon = 0.05$	RCC	0.249^{\ddagger}	334/63/103	0.468
	Interleave	0.217^{\dagger}	298/108/94	2.221
	LC	0.225^{\dagger}	309/103/88	2.365
	Ref-Reorder	0.240^{\ddagger}	343 /66/91	0.596
$\epsilon = 0.2$	RCC	0.209	278/63/159	-4.339
	Interleave	0.199	277/94/129	-2.382
	LC	0.203	275/95/130	-1.984
	Ref-Reorder	0.212	339/75/86	0.679
-		RBP $\phi = 0.8$		
$\epsilon = 0.0$	RCC	0.553 (+0.011) [‡]	440 /20/ 40	11.231
$\epsilon = 0.05$	RCC	0.523 (+0.063) [‡]	289/101/110	-0.356
	Interleave	$0.472 (+0.061)^{\ddagger}$	269/163/ 68	1.991
	LC	$0.489 \ (+0.059)^{\ddagger}$	271/141/88	1.828
	Ref-Reorder	$0.519 \ (+0.056)^{\ddagger}$	301 /100/99	1.384
$\epsilon = 0.2$	RCC	0.441 (+0.215)	235/78/187	-7.302
	Interleave	0.435 (+0.128)	222/170/108	-4.448
	LC	0.444 (+0.137)	218/144/138	-4.313
	Ref-Reorder	0.498 (+0.061) [‡]	290 /119/ 91	1.072

Table 5 shows the NDCG and RBP effectiveness attained in this non-perfect evaluation scenario, as well as their respective risk-reward trade-offs compared to a BM25 baseline, as quantified by TRisk. As a further reference point, Table 5 also includes the effectiveness score of returning the matched centroid run without processing the input query from the user, denoted as RCC (Return Cluster Centroid). We also include RCC for the perfect matching scenario $\epsilon=0.0$, to exhibit how intolerant to failure this strategy is when faced with the possibility of erroneous cluster matching. In an ideal world where all possible clusters have been pre-computed and can be matched perfectly, RCC represents the best score that could be achieved using the approach shown in Figure 1. The performance of all of the approaches degrades as ϵ increases, and they are vulnerable to inaccuracies if the clustering error might be high. The key point is that if a search engine receives an exact match well-known query, it can simply return a cached result, which is what is commonly done in practice [3, 18, 32, 36, 50, 72]; what we are exploring here is the consequence of also allowing approximated matches to be exploited.

Overall, Ref-Reorder is the most robust technique, as was also shown in Table 4, and remains significantly better than LTR in terms of effectiveness, even when 5% of queries are assigned to the wrong cluster. In the two lower halves of the table, with the error rate increased to $\epsilon=0.2$, Ref-Reorder remains statistically significantly better than the LTR baseline. In contrast to the results of Table 4, this method now has the most wins and least number of losses of any of the online fusion approaches considered on both evaluation metrics. Observe also that the higher error rate leads to increased RBP residuals, but with the Ref-Reorder average residual increasing by least. Even with unrealistically large error rates such as $\epsilon=0.5$ (not shown in Table 5), the NDCG score for Ref-Reorder is 0.205 (and for RBP, is 0.462 (+0.075)). That is, even with a very high error rate in terms of cluster identification, retrieval effectiveness remains close to the original BM25 run prior to the application of boosting (Table 4).

6 DISCUSSION

We have presented three novel query effectiveness boosting techniques, built around the idea (Figure 1) of pre-computed query cluster centroids. We now consider some of the issues associated with this proposal, and a range of avenues for possible extension.

Improving Real-Time Rank Fusion. Although the real-time rank fusion approaches outlined in Section 4 can be implemented in a low-latency manner, they nevertheless require substantial amounts of computation. For this reason, we opted to build query centroids offline using a cost-effective rank-fusion approach based on CombSUM (Section 5), and deploy online boosting approaches that make use of those centroids. It would be interesting to further develop these approaches to make them more scalable and less resource intensive. One possibility would be to employ large-scale distributed architectures, with parallel fusion conducted across multiple index server nodes and multiple (perhaps even all) clusters at a time. Another interesting avenue for future work is to develop ways in which non-additive fusion techniques such as RBC [6] and RRF [22] can be incorporated into the same framework. As a motivation, recent work has shown that cost-effective supervised rank fusion techniques can outperform state-of-the-art learning-to-rank models [61].

Simulating Query Intent and Data Privacy. To demonstrate our results in a laboratory setting, we used the UQV100 test collection, in which the clustering is a direct consequence of the data collection process; then, in order to measure the impact of incorrect cluster identification, we also carried out a failure analysis (Section 5.6). Even so, future exploration is required in order to verify the reliability and robustness of the approaches we propose, including whether it is helpful to form centroids with respect to the current step in a user information foraging activity; and how effective the new techniques are outside of test-collection settings.

To validate our observation against a full-scale search engine, query logs and click-graphs would be required to form these clusters using automatic methods. That data is not currently publicly available, and perhaps never will be after the concerns raised in connection with Cambridge Analytica [17] and the earlier AOL log release [7]. Without such data, academics are constrained to exploring performance improvements using data that is publicly available, as we have done here. The upside of using public resources such as the UQV100 queries – no matter how limited they may be in scope – is that the experiments are reproducible, and do not rely on access to private information that may have been gathered from users without adequate consent being given. Fortunately, previous research from industry research labs has shown that query clustering by intent is not only possible, but that it is being used in several different contexts, which we discuss now.

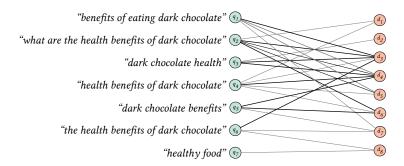


Fig. 7. A bipartite click-graph, showing the associations of document clicks from queries. The thickness of each line represents the frequency of clicks for that query and document pair.

Forming Query Centroids Automatically. A range of authors have carried out experiments in connection with grouping queries by intent (that is, by information need). Relevant applications include query rewriting [9, 11, 38], and forming query clusters using query logs and click-graphs [24, 42, 66]. A team of Microsoft researchers showed that web document click-graphs can be used to generate "virtual" queries by associating documents that are semantically close on the click-graph [75]. This technique was used to create gating features for query variations in LambdaMerge [65]. Craswell and Szummer [24] further demonstrated that random walks on a click graph can form effective query clusters in the domain of image search, noting that this is perhaps not unsurprising, given the prior work of Xue et al. [75].

Similarly, a team of Google researchers [42] write that:

Weighted bigraph clustering capitalizes on organic search results to construct a bipartite graph with a set of queries and a set of URLs as nodes. Edge weights of the graph are computed with the impression and click data of (query, URL) pairs from a Bayesian perspective and are used to induce query (URL) pairwise similarities. Due to information embedded in Google search results, this method is superb in grouping semantically close queries together.

Figure 7 shows an illustrative click-graph which can be combined with random walks to induce query variations directly from large query logs; using such a structure, an at-scale exploration of the ideas we have introduced here would be a useful further step in terms of validating the new approach.

Query Centroid Trade-Offs. In the current experiments, query clusters were generated through a simple CombSUM fusion of BM25 ranked lists constructed from query variations. This simple approach allows for cost-effective cluster computation, as explored in Section 4; with improved query effectiveness a realized goal. Curating more refined (and also more expensive) query clusters might further improve the effectiveness of the proposed boosting approaches.

In a preliminary test of that hypothesis, we built another set of query clusters based on the RBC fusion approach of Bailey et al. [5], who considered rank fusion across both query variations and across different systems – a methodology they refer to as *double fusion*. As input, we ran a number of ranking models from Indri, including BM25, SDM+FIELDS, BM25 with RM3 query expansion, and SDM+FIELDS with RM3 query expansion. We also added several models from Terrier⁵, including PL2; PL2 with Bo1 query expansion; four instances of BM25 with different parameterizations (two with Bo1 query expansion, and two without); and four instances of InL2 with different parameterizations

⁵http://terrier.org/

Table 6. Further effectiveness gains attained using a *double fusion* [5] query centroid construction strategy. The baseline here is taken to be the corresponding boosted mechanism from Table 4 using (only) the BM25-based runs across the training queries; statistical significance and TRISK in this table are measured relative to that starting point. The scores listed in the "Mean" column are averages across the 500 held-out query variations, five for each of the 100 topics; the "W/T/L" numbers are the respective counts of those 500 queries for which the method exceeds the baseline more than 10%, is within 10% of the baseline, and is less than 10% of the baseline.

Metric	Method	Mean	W/T/L	TRisk $_{\alpha=3}$
NDCG@10	Interleave LC	0.233 [†] 0.243 [‡]	166/182/152 212 /148/ 140	-7.092
	LC Ref-Reorder	0.243° 0.252^{\dagger}	182/157/161	-6.163 -7.638
RBP $\phi = 0.8$	Interleave LC Ref-Reorder	0.497 (+0.037) [†] 0.519 (+0.032) [†] 0.539 (+0.049)	152/232/ 116 177 /204/119 176/181/143	-6.602 -6.519 -8.075

(two with query expansion, and two without). There were 14 runs used in total. The resultant ranked lists constructed from each of the training query variations associated with each UQV100 topic was fused together using CombSUM to create a final query cluster result, and then measured using the 500 held-out queries, using the methodology described in Section 3 and used in Section 4.

Table 6 shows the effectiveness of three boosting approaches using these more expensive clusters. The baseline used to ascertain risk-reward profile and statistical significance are the corresponding boosted mechanisms presented in Table 4. Higher quality query centroids translate to improved effectiveness, validating our hypothesis. We plan to explore the interesting problem of additivity and fusion effectiveness across both systems and queries in future work. Note, however, the large negative values associated with this change, indicating harm to the baseline for the TRisk $\alpha=3$ metric when compared to the fused BM25 starting point. The gains recorded are, at this stage, quite unevenly distributed across topics, and more work is required in this area.

Beyond Single-Faceted Information Needs. The UQV100 test collection is, by design, composed of topics for 100 single-facet information needs. Multi-faceted search would invariably provide additional challenges. If a suitable threshold in cluster association scores cannot be met to confidently associate a query with a single centroid, a fusion of many faceted query centroids may be required to resolve the user's information need. An alternative approach could involve building diversified query centroids for information needs that are typically diverse, allowing diversification to be implicitly included via the mechanisms we have introduced here. Search result diversification [1] is an important problem in its own right, and fusion techniques have already been shown to be highly effective for this problem [46]. We plan to also explore this problem further in future work.

7 CONCLUSIONS

We have explored three different strategies to improve search effectiveness for real-time query streams utilizing pre-computed query centroids. In general, centroid-based re-ranking techniques offer a highly efficient mechanism for boosting query effectiveness – requiring on average only 2 milliseconds to reorder 1,000 documents, while improving effectiveness significantly. We further

show that on-the-fly rank fusion is viable, and can be reasonably efficient using current state-of-the-art dynamic pruning techniques, but if aggressive SLAs on query performance are enforced, carrying out fusion at run-time remains costly.

We then demonstrated that query level fusion can instead be used to combine similar queries *offline*, making it a practical alternative in high-performance search engines. To validate this idea, we explore how query clusters can be used to improve the effectiveness of incoming queries using three different approaches, namely re-ranking, interleaving and a linear combination of the cluster and the user query. Experiments using the ClueWeb12B UQV100 collection show that the new approaches we describe provide competitive efficiency, and, at the same time, effectiveness improvements over strong baselines in a performance-focused query processing framework.

Software. In the interests of reproducibility, our codebase is available at https://github.com/rmit-ir/centroid-boost.

Acknowledgements. This work was supported by the Australian Research Council's *Discovery Projects* Scheme (DP170102231) and a grant from the Mozilla Foundation.

REFERENCES

- [1] R. Agrawal, S. Gollapudi, A. Halverson, and S. Ieong. 2009. Diversifying search results. In *Proc. Conf. on Web Search and Data Mining (WSDM)*. 5–14.
- [2] Y. Anava, A. Shtok, O. Kurland, and E. Rabinovich. 2016. A probabilistic fusion framework. In *Proc. International Conf. on Theory of Information Retrieval (ICTIR)*. 1463–1472.
- [3] R. Baeza-Yates, A. Gionis, F. Junqueira, V. Murdock, V. Plachouras, and F. Silvestri. 2008. Design trade-offs for search engine caching. *ACM Trans. on the Web* 2, 4 (2008), 1–28.
- [4] X. Bai, I. Arapakis, B. B. Cambazoglu, and A. Freire. 2017. Understanding and leveraging the impact of response latency on user behaviour in web search. ACM Trans. on Information Systems 36, 2 (2017), 1–42.
- [5] P. Bailey, A. Moffat, F. Scholer, and P. Thomas. 2016. UQV100: A test collection with query variability. In *Proc. ACM Conf. on Research and Development in Information Retrieval (SIGIR)*. 725–728.
- [6] P. Bailey, A. Moffat, F. Scholer, and P. Thomas. 2017. Retrieval consistency in the presence of query variations. In Proc. ACM Conf. on Research and Development in Information Retrieval (SIGIR). 395–404.
- [7] M. Barbaro and T. Zeller. 2006. A face is exposed for AOL searcher No. 4417749. https://nytimes.com/2006/08/09/technology/09aol.html. (Aug. 2006). Accessed: 2018-11-08.
- [8] N. J. Belkin, C. Cool, W. B. Croft, and J. P. Callan. 1993. The effect of multiple query variations on information retrieval system performance. In *Proc. ACM Conf. on Research and Development in Information Retrieval (SIGIR)*. 339–346.
- [9] M. Bendersky, D. Metzler, and W. B. Croft. 2012. Effective query formulation with multiple information sources. In *Proc. Conf. on Web Search and Data Mining (WSDM)*. 443–452.
- [10] R. Benham and J. S. Culpepper. 2017. Risk-reward trade-offs in rank fusion. In *Proc. Australasian Document Computing Symp. (ADCS)*. 1:1–1:8.
- [11] R. Benham, J. S. Culpepper, L. Gallagher, X. Lu, and J. Mackenzie. 2018. Towards efficient and effective query variant generation. In *Proc. Conf. on Design of Experimental Search & Information Retrieval Systems (DESIRES)*. 62–67.
- [12] R. Benham, L. Gallagher, J. Mackenzie, T. T. Damessie, R-C. Chen, F. Scholer, A. Moffat, and J. S. Culpepper. 2017. RMIT at the 2017 TREC CORE track. In *Proc. Text Retrieval Conf. (TREC)*.
- [13] R. Benham, L. Gallagher, J. Mackenzie, B. Liu, X. Lu, F. Scholer, A. Moffat, and J. S. Culpepper. 2018. RMIT at the 2018 TREC CORE track. In *Proc. Text Retrieval Conf. (TREC)*.
- [14] A. Z. Broder, D. Carmel, M. Herscovici, A. Soffer, and J. Y. Zien. 2003. Efficient query evaluation using a two-level retrieval process. In *Proc. ACM International Conf. on Information and Knowledge Management (CIKM).* 426–434.
- [15] C. Buckley and J. Walz. 1999. The TREC-8 query track. In Proc. Text Retrieval Conf. (TREC).
- [16] C. Burges, R. Ragno, and Q. V. Le. 2006. Learning to Rank with nonsmooth cost functions. In Proc. Conf. on Neural Information Processing Systems (NIPS). 193–200.
- [17] C. Cadwalladr and E. Graham-Harrison. 2018. Revealed: 50 million Facebook profiles harvested for Cambridge Analytica in major data breach. https://theguardian.com/news/2018/mar/17/cambridge-analytica-facebook-influence-us-election. (March 2018). Accessed: 2018-11-08.

- [18] B. B. Cambazoglu, F. P. Junqueira, V. Plachouras, S. Banachowski, B. Cui, S. Lim, and B. Bridge. 2010. A refreshing perspective of search engine caching. In Proc. Conf. on the World Wide Web (WWW). 181–190.
- [19] K. Chakrabarti, S. Chaudhuri, and V. Ganti. 2011. Interval-based pruning for top-k processing over compressed lists. In *Proc. International Conf. on Data Engineering (ICDE)*. 709–720.
- [20] R-C. Chen, L. Gallagher, R. Blanco, and J. S. Culpepper. 2017. Efficient cost-aware cascade ranking in multi-stage retrieval. In Proc. ACM Conf. on Research and Development in Information Retrieval (SIGIR). 445–454.
- [21] F. M. Choudhury, J. S. Culpepper, Z. Bao, and T. Sellis. 2018. Batch processing of top-k spatial-textual queries. ACM Trans. on Spatial Algorithms and Systems 3, 4 (2018), 1–40.
- [22] G. V. Cormack, C. L. A. Clarke, and S. Buettcher. 2009. Reciprocal rank fusion outperforms Condorcet and individual rank learning methods. In Proc. ACM Conf. on Research and Development in Information Retrieval (SIGIR). 758–759.
- [23] M. Crane, J. S. Culpepper, J. Lin, J. Mackenzie, and A. Trotman. 2017. A comparison of Document-at-a-Time and Score-at-a-Time query evaluation. In *Proc. Conf. on Web Search and Data Mining (WSDM)*. 201–210.
- [24] N. Craswell and M. Szummer. 2007. Random walks on the click graph. In *Proc. ACM Conf. on Research and Development in Information Retrieval (SIGIR)*. 239–246.
- [25] V. Dang and W. B. Croft. 2010. Query reformulation using anchor text. In *Proc. Conf. on Web Search and Data Mining (WSDM)*. 41–50.
- [26] J. Dean and L. A. Barroso. 2013. The tail at scale. Comm. ACM 56, 2 (2013), 74-80.
- [27] L. Dhulipala, I. Kabiljo, B. Karrer, G. Ottaviano, S. Pupyrev, and A. Shalita. 2016. Compressing graphs and indexes with recursive graph bisection. In Proc. Conf. on Knowledge Discovery and Data Mining (WSDM). 1535–1544.
- [28] C. Dimopoulos, S. Nepomnyachiy, and T. Suel. 2013. Optimizing top-k document retrieval strategies for Block-Max indexes. In *Proc. Conf. on Web Search and Data Mining (WSDM)*. 113–122.
- [29] B. T. Dinçer, C. Macdonald, and I. Ounis. 2014. Hypothesis testing for the risk-sensitive evaluation of retrieval systems. In *Proc. ACM Conf. on Research and Development in Information Retrieval (SIGIR)*. 23–32.
- [30] S. Ding, J. Attenberg, R. Baeza-Yates, and T. Suel. 2011. Batch query processing for web search engines. In *Proc. Conf. on Web Search and Data Mining (WSDM)*. 137–146.
- [31] S. Ding and T. Suel. 2011. Faster top-k document retrieval using Block-Max indexes. In *Proc. ACM Conf. on Research and Development in Information Retrieval (SIGIR)*. 993–1002.
- [32] T. Fagni, R. Perego, F. Silvestri, and S. Orlando. 2006. Boosting the performance of web search engines: Caching and prefetching query results by exploiting historical usage data. ACM Trans. on Information Systems 24, 1 (2006), 51–78.
- [33] M. Fontoura, V. Josifovski, J. Liu, S. Venkatesan, X. Zhu, and J. Zien. 2011. Evaluation strategies for top-k queries over memory-resident inverted indexes. *Proc. Conf. on Very Large Databases (VLDB)* 4, 12 (2011), 1213–1224.
- [34] E. A. Fox and J. A. Shaw. 1994. Combination of multiple searches. Proc. Text Retrieval Conf. (TREC) (1994), 243-252.
- [35] N. Fuhr. 2018. Some common mistakes in IR evaluation, and how they can be avoided. In *Proc. ACM Conf. on Research and Development in Information Retrieval (SIGIR)*. 32–41.
- [36] Q. Gan and T. Suel. 2009. Improved techniques for result caching in web search engines. In Proc. Conf. on the World Wide Web (WWW). 431–440.
- [37] E. J. Glover, S. Lawrence, W. P. Birmingham, and C. L. Giles. 1999. Architecture of a metasearch engine that supports user information needs. In *Proc. ACM International Conf. on Information and Knowledge Management (CIKM)*. 210–216.
- [38] Y. He, J. Tang, H. Ouyang, C. Kang, D. Yin, and Y. Chang. 2016. Learning to rewrite queries. In *Proc. ACM International Conf. on Information and Knowledge Management (CIKM)*. 1443–1452.
- [39] S. Huo, M. Zhang, Y. Liu, and S. Ma. 2014. Improving tail query performance by fusion model. In *Proc. ACM International Conf. on Information and Knowledge Management (CIKM)*. 559–658.
- [40] K. Järvelin and J. Kekäläinen. 2002. Cumulated gain-based evaluation of IR techniques. ACM Trans. on Information Systems 20, 4 (2002), 422–446.
- [41] S. Kim, Y. He, S-W. Hwang, S. Elnikety, and S. Choi. 2015. Delayed-dynamic-selective (DDS) prediction for reducing extreme tail latency in web search. In *Proc. Conf. on Web Search and Data Mining (WSDM)*. 7–16.
- [42] J. Kong, A. Scott, and G. M. Goerg. 2016. Improving semantic topic clustering for search queries with word co-occurrence and bigraph co-clustering. *Google Inc* (2016).
- [43] A. K. Kozorovitsky and O. Kurland. 2011. Cluster-based fusion of retrieved lists. In *Proc. ACM Conf. on Research and Development in Information Retrieval (SIGIR)*. 893–902.
- [44] O. Kurland and J. S. Culpepper. 2018. Tutorial / Fusion in information retrieval. In *Proc. ACM Conf. on Research and Development in Information Retrieval (SIGIR)*. 1383–1386.
- [45] C.-J. Lee, Q. Ai, W. B. Croft, and D. Sheldon. 2015. An optimization framework for merging multiple result lists. In Proc. ACM International Conf. on Information and Knowledge Management (CIKM). 303–312.

- [46] S. Liang, Z. Ren, and M. de Rijke. 2014. Fusion helps diversification. In Proc. ACM Conf. on Research and Development in Information Retrieval (SIGIR). 303–312.
- [47] T.-Y. Liu. 2009. Learning to rank for information retrieval. Foundations & Trends in Information Retrieval 3, 3 (2009), 225–331.
- [48] X. Lu, A. Moffat, and J. S. Culpepper. 2016. The effect of pooling and evaluation depth on IR metrics. *Information Retrieval* 19, 4 (2016), 416–445.
- [49] X. Lu, A. Moffat, and J. S. Culpepper. 2016. Efficient and effective higher order proximity modeling. In *Proc. International Conf. on Theory of Information Retrieval (ICTIR)*. 21–30.
- [50] H. Ma and B. Wang. 2012. User-aware caching and prefetching query results in web search engines. In Proc. ACM Conf. on Research and Development in Information Retrieval (SIGIR). 1163–1164.
- [51] C. Macdonald, R. L. T. Santos, and I. Ounis. 2013. The whens and hows of learning to rank for web search. *Information Retrieval* 16, 5 (2013), 584–628.
- [52] C. Macdonald, R. L. T. Santos, I. Ounis, and B. He. 2013. About learning models with multiple query-dependent features. *ACM Trans. on Information Systems* 31, 3 (2013), 11:1–11:39.
- [53] J. Mackenzie. 2017. Managing tail latencies in large scale IR systems. In Proc. ACM Conf. on Research and Development in Information Retrieval (SIGIR). 1369.
- [54] J. Mackenzie, J. S. Culpepper, R. Blanco, M. Crane, and J. Lin. 2018. Query driven algorithm selection in early stage retrieval. In Proc. Conf. on Web Search and Data Mining (WSDM). 396–404.
- [55] A. Mallia, G. Ottaviano, E. Porciani, N. Tonellotto, and R. Venturini. 2017. Faster BlockMax WAND with variable-sized blocks. In Proc. ACM Conf. on Research and Development in Information Retrieval (SIGIR). 625–634.
- [56] D. Metzler and W. B. Croft. 2005. A Markov random field model for term dependencies. In Proc. ACM Conf. on Research and Development in Information Retrieval (SIGIR). 472–479.
- [57] D. Metzler, S. Dumais, and C. Meek. 2007. Similarity measures for short segments of text. In Proc. European Conf. in Information Retrieval (ECIR). 16–27.
- [58] A. Moffat. 2016. Judgment pool effects caused by query variations. In Proc. Australasian Document Computing Symp. (ADCS). 65–68.
- [59] A. Moffat, P. Bailey, F. Scholer, and P. Thomas. 2017. Incorporating user expectations and behavior into the measurement of search effectiveness. *ACM Trans. on Information Systems* 35, 3 (2017), 24:1–24:38.
- [60] A. Moffat and J. Zobel. 2008. Rank-biased precision for measurement of retrieval effectiveness. ACM Trans. on Information Systems 27, 1 (2008), 2.1–2.27.
- [61] A. Mourão and J. Magalhães. 2018. Low-complexity supervised rank fusion models. In *Proc. ACM International Conf. on Information and Knowledge Management (CIKM)*. 1691–1694.
- [62] G. Ottaviano and R. Venturini. 2014. Partitioned Elias-Fano indexes. In *Proc. ACM Conf. on Research and Development in Information Retrieval (SIGIR)*. 273–282.
- [63] F. Radlinski, M. Kurup, and T. Joachims. 2008. How does clickthrough data reflect retrieval quality?. In Proc. ACM International Conf. on Information and Knowledge Management (CIKM). 43–52.
- [64] F. Scholer and H. E. Williams. 2002. Query association for effective retrieval. In Proc. ACM International Conf. on Information and Knowledge Management (CIKM). 324–331.
- [65] D. Sheldon, M. Shokouhi, M. Szummer, and N. Craswell. 2011. LambdaMerge: Merging the results of query reformulations. In Proc. Conf. on Web Search and Data Mining (WSDM). 795–804.
- [66] J. Shen, M. Karimzadehgan, M. Bendersky, Z. Qin, and D. Metzler. 2018. Multi-task learning for email search ranking with auxiliary query clustering. In Proc. ACM International Conf. on Information and Knowledge Management (CIKM). 2127–2135.
- [67] T. Strohman, H. Turtle, and W. B. Croft. 2005. Optimization strategies for complex queries. In Proc. ACM Conf. on Research and Development in Information Retrieval (SIGIR). 219–225.
- [68] H. R. Turtle and J. Flood. 1995. Query evaluation: Strategies and optimizations. *Information Processing & Management* 31, 6 (1995), 831–850.
- [69] C. C. Vogt. 2000. How much more is better? Characterizing the effects of adding more IR systems to a combination. In *Proc. Recherche d'Information etses Applications (RIAO)*. 457–475.
- [70] C. C. Vogt and G. W. Cottrell. 1999. Adaptive combination of evidence for information retrieval. University of California, San Diego.
- [71] C. C. Vogt and G. W. Cottrell. 1999. Fusion via a linear combination of scores. Information Retrieval 1, 3 (1999), 151-173.
- [72] J. Wang, E. Lo, M. L. Yiu, J. Tong, G. Wang, and X. Liu. 2014. Cache design of SSD-based search engine architectures: An experimental study. ACM Trans. on Information Systems 32, 4 (2014), 1–26.

- [73] J. R. Wen, J. Y. Nie, and H. J. Zhang. 2001. Clustering user queries of a search engine. In *Proc. Conf. on the World Wide Web (WWW)*. 162–168.
- [74] J. R. Wen, J. Y. Nie, and H. J. Zhang. 2002. Query clustering using user logs. ACM Trans. on Information Systems 20, 1 (2002), 59–81.
- [75] G-R. Xue, H-J. Zeng, Z. Chen, Y. Yu, W-Y. Ma, W. Xi, and W. Fan. 2004. Optimizing web search using web click-through data. In *Proc. ACM International Conf. on Information and Knowledge Management (CIKM)*. 118–126.
- [76] X. Xue and W. B. Croft. 2013. Modeling reformulation using query distributions. ACM Trans. on Information Systems 31, 2 (2013), 6:1–6:34.
- [77] J-M. Yun, Y. He, S. Elnikety, and S. Ren. 2015. Optimal aggregation policy for reducing tail latency of web search. In *Proc. ACM Conf. on Research and Development in Information Retrieval (SIGIR)*. 63–72.
- [78] J. Zobel and A. Moffat. 2006. Inverted files for text search engines. Comput. Surveys 38, 2 (2006), 6.1-6.56.

Received November, 2018