

Fast Parallel Top-K Retrieval of Verbose Queries

Gali Sheffi^{1,2}, Dmitry Basin¹, Edward Bortnikov¹, David Carmel³, and Idit Keidar^{1,2}

¹ Yahoo Research

² Technion

³ Amazon

Abstract. Top-k document retrieval is emerging as a performance bottleneck in query processing as verbose queries are becoming mainstream, while state-of-the-art algorithms fail to process such long queries in real time. To date, attempts to accelerate top-k evaluation via approximate result computation and intra-query parallelism have produced limited results.

We present Sparta – a practical parallel algorithm that exploits multi-core hardware to run approximate top-k retrieval within interactive latency bounds. Its design is inspired by Fagin et al.’s seminal Threshold Algorithm for aggregation in databases. Sparta scales through lightweight coordination and context sharing among concurrent threads. The resulting algorithm is highly scalable. When applied to the 50M-document public ClueWeb09B dataset, Sparta processes 12-term queries 3.6 times faster than the state-of-the-art. Sparta further significantly improves throughput.

Keywords: Parallel computing, web-search, top-k

1 Introduction

Internet users are changing the ways in which they interact with search engines. While early days web queries were short (2.4 terms on average [?]), modern search experiences (query suggestion, reformulation, conversational interfaces, etc.) stimulate their users to submit much longer queries. For example, more than 5% of queries submitted via voice search on mobile devices exceed 10 terms [?]. At the same time, usability studies (e.g., [?]) show that users are extremely sensitive to end-to-end delays beyond 500 ms, and any excess delay beyond 250 ms leads to material degradation in their experience. Maintaining this *service-level agreement (SLA)* is becoming a major challenge as more queries become longer.

Most search systems today evaluate queries via a two-stage process [?]. The first stage is *top-k retrieval*: it roughly matches the top-k documents most relevant to the query (typically, hundreds to thousands) based on some simple relevance scoring function like tf-idf or BM25 [?]. The second stage then re-ranks the documents via some sophisticated (e.g., machine-learned) function, to produce the final result list (typically, a few tens of documents). The first stage sifts through huge volumes of data and dominates the execution time. We therefore focus on the latency induced by this stage.

The retrieval stage is typically executed in a backend tier on dedicated search nodes, each of which retrieves the top-k results from a local index shard. With the advance of

storage and memory technologies, the size of index shards increases, thus increasing the load on each search node. A promising approach to handle this load is to parallelize query evaluation on multiprocessor hardware. The current trend in server architecture, which favors parallelism over sequential speed, makes this approach particularly attractive. We explore this approach in this paper.

Obtaining the exact top-k matches out of a large corpus is typically too slow to meet stringent SLA requirements. Given that the top-k selection is only the first query processing step, perfect results are often not essential, as later query processing stages can be satisfied by approximate results [?, ?, ?]. Based on this observation we focus on *approximate* (sometimes called *non-safe*) query evaluation, tuned to achieve a certain high recall (e.g., 97% or more).

State-of-the-art sequential algorithms for approximate top-k retrieval, e.g., *Block-Max WAND (BMW)* [?], serve traditional (short) web queries with impressive speed. However, these algorithms do not scale well to verbose queries (as pointed out in [?] and confirmed by our experiments). Parallelizing these algorithms can improve performance by 30%–40% [?], but our experiments show that such a parallel BMW implementation still has limited scalability, falling short of expected SLAs on long queries.

We present *Sparta* (Scalable PARallel Threshold Algorithm) – a novel concurrent algorithm that substantially improves search time for verbose queries over large search indices. Its query latencies fit well within real-time SLAs on standard server hardware.

Sparta’s design is inspired by the seminal *Threshold Algorithm (TA)* by Fagin et al. [?] for retrieving the top-k objects from a database based on an aggregation of attributes that may reside in multiple nodes. In the web setting, term posting lists represent attributes and the ranking function is a linear aggregation of term scores. TA has two variants [?]: (1) *Random Access (RA)*, which relies on random access to posting lists; and (2) *No Random Access (NRA)*, which only traverses them sequentially. The former is ineffective in the context of real-time search because it requires (1) a secondary random-access index (which doubles the required index space), and (2) random I/O to the on-disk secondary index, which incurs high overhead when the index is too big to be kept entirely in RAM. We therefore base our algorithm on NRA.

Transforming NRA into an efficient concurrent algorithm is challenging because coordination around shared state can become a major bottleneck unless carefully designed. On the one hand, sharing state among threads is essential in order to benefit from TA’s early stopping feature. Indeed, our experiments show that a shared-nothing (partitioned) parallelization of NRA performs two times worse than a single-threaded implementation. On the other hand, a naïve attempt to parallelize NRA using shared memory also results in worse performance than the sequential algorithm. Sparta judiciously shares information among threads while significantly reducing the synchronization overhead and memory overlap among them, leading to major performance gains.

Our results show that Sparta scales well with query length. E.g., on the 50M-document TREC ClueWeb09B dataset [?], Sparta can serve verbose 12-term queries within less than 200 ms, whereas a parallel execution of RA processes them in 400 ms, and a parallel execution of BMW takes more than 600 ms. Sparta’s throughput on a production-grade query mix (with the query length distribution measured in [?]) is twice that of parallel BMW.

Summing up, we present a practical scalable parallel algorithm that significantly improves the state-of-the-art latency of long query processing. The paper proceeds as follows: Section 2 defines the top-k retrieval problem and associated metrics. Section 3 gives background on existing algorithms. Section 4 presents Sparta, and Section 5 evaluates it. Finally, Section 6 discusses related work and Section 7 concludes the paper.

2 Top-k Retrieval

We focus on the fundamental primitive of top-k retrieval, commonly used by search engines for initial selection of documents over which more refined search is performed [?]. Given a query q , the primitive retrieves the top k scored documents in the corpus according to a scoring function $score(D, q)$. More specifically, a *query* is given as a list of terms (to search for). Given an m -term query $q = t_1, \dots, t_m$ and a document D , the score of D for query q is $score(D, q) \triangleq \sum_{i=1}^m ts(D, t_i)$ where $ts(D, t_i)$ is the term score of term t_i to document D .

An *exact* top-k retrieval primitive returns a list of the k documents with the highest scores for the query. An *approximate* solution returns a list of k results that approximates the top k . Let L be the exact list of top-k documents, and let A be the list returned by an approximate solution. The quality (or accuracy) of A is measured by its *recall*, which is the fraction of L included in A .

3 Background

We provide background on state-of-the-art top-k algorithms and Fagin et al.’s TA [?].

3.1 State-of-the-art Top-k Algorithms

Search algorithms use a preprocessed inverted index of the corpus. The index is organized according to terms. For each term, the index holds a *posting list* listing all documents associated with that term. Top-k retrieval algorithms typically traverse posting lists sequentially; multiple lists may be scanned simultaneously. For big datasets, only a small portion of the index can reside in RAM at any given time. However, the I/O overhead is low because contiguous chunks of the lists are infrequently fetched from disk into memory.

In order to avoid scoring a huge number of documents per query, state-of-the-art algorithms reduce the number of evaluated documents while identifying the top scored results. *MaxScore* [?,?], *WAND* [?], and *Block-Max WAND (BMW)* [?] are popular examples of such algorithms, widely used in production systems. They simultaneously scan all relevant posting lists in order of document id, evaluating the full score of each document before moving to the next one. We therefore refer to these as *document-order* algorithms. Such algorithms track the top-k documents among those scored so far (typically, in a heap). A variable Θ – called the *threshold* – holds the score of the k -th (lowest-ranked) document in the heap; any document whose score is below this threshold is not a candidate for the final top-k list. As long as the heap contains fewer than k documents, Θ remains zero.

The disadvantage of document-order algorithms is that high-scoring documents may be discovered late in the traversal because document ids are not correlated with query relevance, and so the algorithm does not always produce useful partial results before it completes. This is mitigated by *score-order* algorithms (sometimes called *impact-order*), e.g., JASS [?], which traverse posting lists in decreasing order of term score. These algorithms accumulate the score of each document throughout the traversal, and thus document scores may be inaccurate at first and improve over time. Score-order algorithms have been shown to be slower but have more predictable performance than document-based ones [?].

3.2 The Threshold Algorithm

TA [?] was originally presented in a database setting, where the partial scores of an item (term scores in our context) reside at different nodes. We cast it here in the IR setting, where partial scores are obtained from posting lists rather than nodes, and query evaluation occurs on a single machine, possibly using multiple threads accessing shared memory.

TA is a score-order algorithm. To evaluate a query $q = t_1, \dots, t_m$, it dynamically traverses the posting lists of the m query terms in an interleaved manner. An *upper bound* vector, $UB[m]$, bounds the term scores of documents that were not yet visited in each term's posting list. Figure 1 shows an example posting list traversal and the corresponding values in UB . Here, $m = 3$ and the scores of the last traversed items in each list are $UB[1] = 38$, $UB[2] = 32$, and $UB[3] = 41$.

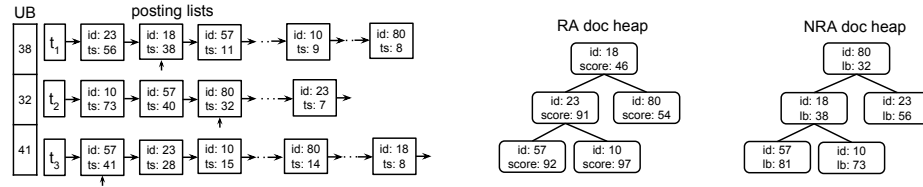


Fig. 1. Traversal example in RA and NRA variants of the Threshold Algorithm (TA). Posting lists are sorted by decreasing term score. Vertical arrows depict iterator positions. UB holds upper bounds on the terms' contributions to scores of untraversed documents. The RA document heap is ordered by full document score, and the NRA heap by partial score (lower bound).

TA also maintains a threshold Θ , i.e., the score of the k -th document in the heap. It stops when no candidate can exceed the threshold score. Fagin et al. present two flavors of the algorithm, which we now describe.

Random Access (RA) The RA variant assumes that given a document id, we can use random access in order to obtain all its term scores in order to compute its total score. It thus computes the full score for every document it encounters. If the score is higher than the threshold, it is inserted into the heap. Then, Θ and the corresponding term's UB

are updated. The algorithm stops when the following *upper bound stopping condition* holds:

$$UBStop \triangleq \sum_{i=1}^m UB[i] \leq \Theta. \quad (1)$$

At this point, no non-traversed document can achieve a high enough score to be included in the heap. RA's output is the set of documents in the heap, ordered by their scores.

RA is an exact algorithm as it returns the top-k results. In our evaluation, we implement an approximate variant of RA by stopping whenever the heap does not change for some parameter Δ ms. Note that the Threshold Algorithm is amenable to such early stopping because it traverses posting lists in order of score. High scoring documents are likely to become evident early, and finding new high scoring documents becomes less and less likely as the algorithm progresses. This is in contrast with document-order algorithms like BMW, where finding high scoring documents remains equally likely throughout the execution.

Unfortunately, random access is costly, in particular for large data sets that do not reside fully in RAM. Whereas a sequential posting list traversal requires infrequent I/O (at the end of each data block) and exhibits cache-friendly locality of access, each random document access entails an I/O request and a cache miss. RA has an additional drawback in the IR setting, as it needs to maintain a secondary index by document id in addition to the posting lists sorted by term score, which doubles its footprint.

No Random Access (NRA) The alternative NRA method refrains from computing the full score for each traversed document, and instead maintains lower and upper bounds for candidate documents based on *partially* computed scores. For a document D and a term t_i , we define the upper bound $UB(D, t_i)$ to be the term score $ts(D, t_i)$ if it has already been encountered, and otherwise $UB[i]$, which provides an upper bound on t_i 's term score. We similarly define its lower bound $LB(D, t_i)$ to be the term score if it is known, and zero otherwise. We then aggregate these scores to compute the document's upper and lower bounds:

$$UB(D) \triangleq \sum_{i=1}^m UB(D, t_i); \quad LB(D) \triangleq \sum_{i=1}^m LB(D, t_i).$$

For example, in Figure 1, $UB(D_{57}) = 38 + 40 + 41 = 119$ and $LB(D_{57}) = 40 + 41 = 81$, whereas its actual score is $11 + 40 + 41 = 92$.

NRA maintains the top-k heap according to the document *lower bounds*, and Θ holds the smallest value among them. Its output is the set of documents in the heap, sorted by LB.

NRA's safe variant stops when (1) the *UBStop* stopping condition of RA holds, and (2) all the visited documents that are not in the heap have *upper bounds* lower than or equal to Θ . These two conditions are complementary: (1) ensures that no non-traversed documents are among the final top-k, whereas (2) ensures the same for traversed documents that are not among the current top-k. While NRA does return the exact top-k results, unlike RA, it does not necessarily preserve the order among them, since some

returned documents may be partially scored. As in RA, the approximate variant stops after the heap has not changed for Δ ms.

4 Sparta

Sparta is a parallel algorithm that adapts NRA to shared-memory multiprocessor hardware platforms. Like our implementation of RA and NRA described above, it can be safe, or can be configured to provide approximate results by stopping after the heap does not change for some Δ time.

Section 4.1 describes the algorithm’s data structures. Section 4.2 explains how we divide the work involved in query processing among threads. Section 4.3 describes synchronization around the shared data structures. Finally, Section 4.4 discusses the properties of our algorithm.

4.1 Sparta Data Structures

name	value
<i>DocType</i>	$\langle \text{int id, int score}[m], \text{int LB} \rangle$
<i>docHeap</i>	init empty
Θ	init 0
$UB[m]$	init ∞
$UBStop$	$\sum_{i=1}^m UB[i] \leq \Theta$
$UB(D)$	$\sum_{i=1}^m (D.score[i] > 0 ? D.score[i] : UB[i])$
<i>docMap</i>	init empty
<i>heapUpdTime</i>	init now
<i>done</i>	init false
<i>termMap</i> [m]	init pointer to <i>docMap</i>
<i>tmpDocMap</i>	init empty

Table 1. Sparta’s data structures and initial values.

Table 1 defines the data structures used by Sparta and Figure 2 illustrates them. As in NRA, the algorithm maintains the current top-k results in a heap, *docHeap*, and its lowest value in Θ . It keeps m pointers to the next elements to traverse in all posting lists (not listed in Table 1) and an array UB of upper bounds on non-traversed term scores. UB is used for computing the documents’ upper bounds as well as for checking NRA’s stopping conditions.

The hash map *docMap* maps document ids encountered thus far to document *DocType* objects. A *DocType* holds a vector of term scores observed thus far for this document as well as a lower bound on the document score computed as their sum. NRA’s first stopping condition is checked by the macro *UBStop*, and the second stopping condition is checked as follows:

$$\forall D \in docMap \setminus docHeap : UB(D) \leq \Theta. \quad (2)$$

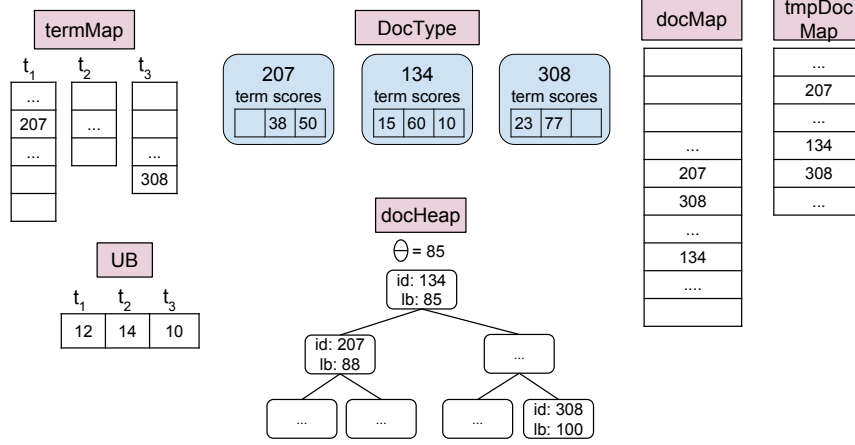


Fig. 2. Sparta data structures. The `docMap` keeps track of partially scored documents; `tmpDocMap` and `termMap` are local partial copies of `docMap`.

The stopping conditions are evaluated by a *cleaner* task, which also checks the heap’s latest update time `heapUpdTime` and sets the *done* flag once the algorithm can stop.

In addition, to reduce the synchronization overhead and improve cache locality, Sparta uses two local data structures that hold partial copies of the global `docMap`, namely the `termMap` array with a (local) hash map per term, and the `tmpDocMap` used by a dedicated maintenance thread. The role of these will become evident when we discuss synchronization and locality below.

4.2 Splitting the Work

A naïve attempt to parallelize NRA would be to share the data structures among all threads. (Note that sharing state is important because the partial document scores, and consequently the lower bounds, are affected by multiple threads that generate term scores). This approach leads to high contention, primarily around `docMap`.

To reduce contention, we consider the point in time when the first stopping condition (Equation 1) holds. From this time on, no new document’s score can surpass the lower bound of any document in the shared `docHeap`. Therefore, adding new documents to `docMap` is no longer helpful (a similar observation was made in [?]). On the other hand, it is possible to shrink `docMap` by removing documents whose upper bounds have been exceeded by Θ . This is exploited by our concurrent implementation, which stops sharing `docMap` among the threads once it is sufficiently small, thereby eliminating the synchronization overhead altogether.

Sparta’s pseudocode appears in Algorithm 1. It exploits up to m *worker* threads per query, but can run with fewer threads if less are available. We divide posting list traversals to segments of size `segSize`, and use a job queue to allocate posting list segments to threads (line 2). The `PROCESSTERM(i)` function processes the next segment of term t_i .

A thread that finishes its assigned segment inserts into the queue a new task for scanning the next segment in the same term’s posting list (line 26). Thus, we progress on all posting lists at the same rate modulo the segment size. In case m threads are available, a large segment size can be used.

In addition to posting list-traversing tasks, a task executing the CLEANER function (lines 40–49) performs maintenance on the *docMap*. This task is invoked once Equation 1 holds and so *docMap* no longer grows. The cleaner serves two purposes. First, as its name suggests, it removes entries that ceased to be top-k candidates from *docMap*. Since Sparta is memory-intensive, a smaller *docMap* allows it to run much faster; (a similar observation, in a sequential setting, was made in [?]). Second, it detects the stopping conditions (line 47). In the approximate version, it checks whether the heap has not changed for Δ time (the exact version is obtained by setting $\Delta = \infty$). It also checks the condition of Equation 2: once *docMap* is the same size as *docHeap* we know that the two are identical because *docMap* always includes all *docHeap* entries. At this point, *docHeap* holds the top-k scored results, and stopping is safe. Once the algorithm stops, the main thread returns the heap’s contents.

4.3 Synchronization

Note that, *docHeap*, *UB*, *docMap*, and *DocType* objects referenced by them are accessed concurrently by multiple threads. We need to protect such access to avoid inconsistencies. On the other hand, reducing contention is crucial for performance. Moreover, Sparta is a memory-intensive algorithm, and in order to keep the memory access latencies low, it is paramount to exploit the CPU hardware cache, in particular the core-private L1 caches. We now explain how we synchronize access to each of the shared variables in a way that reduces contention and improves cache utilization.

Since at most one thread processes each term, no races arise around updating *UB* entries, and no lock is needed. However, all threads read all *UB* entries, and therefore frequent updates can lead to frequent cache misses, and in turn, poor performance. In order to reduce the number of cache misses, instead of updating *UB* after each document evaluation, the workers update it lazily, at the end of a segment traversal (line 25). Since upper bounds can only decrease whereas Θ can only increase, such lazy updates do not affect correctness.

Updates of *docHeap* and Θ are protected by a shared lock (lines 28 and 39), which serializes all updates. To avoid races around evaluating a *DocType*’s lower bound and inserting it into *docHeap*, we update the lower bound in a lazy manner while holding the global lock on *docHeap*: Every thread that adds a document to the heap updates the lower bounds of all heap documents (lines 31-33).

Before the first stopping condition (Equation 1) holds, multiple workers update *docMap* concurrently. We therefore protect each hash bucket by a granular lock, which performs better than the generic Java concurrent hashmap [?].

The cleaner task starts removing elements from *docMap* after it is guaranteed that no new entries are added to it; such removals substantially improve the term processing performance. Nevertheless, allowing the cleaner to constantly update *docMap* would lead to frequent cache invalidations at the tasks that read the map. To avoid frequent cache misses, the global map is kept read-only most of the time, while the cleaner works

on a local copy: it repeatedly builds a new map *tmpDocMap*, holding *docMap* entries whose upper bounds are higher than Θ as well as ones that are included in *docHeap* (whose upper bounds may be exactly Θ); recall that other *docMap* entries no longer need to be kept. Once *tmpDocMap* is ready, the cleaner replaces *docMap* with it via a single pointer swing (flipping the global reference).

Access to *docMap* is a principal performance bottleneck, since it is frequently read by all worker threads. Initially, it is too large to fit into local caches, and so the parallel execution inherently requires global memory accesses. But thanks to the cleaner’s work, *docMap* shrinks in the course of the execution. Moreover, not all *docMap* entries are relevant for all terms – if D ’s term score for t_i has already been computed, then a thread handling term t_i does not need to access D . Thus, the relevant subset of *docMap* for each term eventually becomes small enough to fit in its local cache. As long as the thread continues to access the global *docMap*, it experiences massive cache misses every time the cleaner replaces the global *docMap*. But once it becomes small enough to fully fit in the local cache, there is no need to keep using the global copy.

To this end, Sparta associates a local map replica, *termMap*, with each posting list. *termMap* is created by the worker that currently owns that posting list once *docMap*’s size drops below a threshold Φ , in our implementation, $\Phi = 10\text{K}$ entries. In lines 11-12, we scan *docMap*, and copy to *termMap* the references to those *DocType* objects that do not contain the score for the worker’s term yet. Once a *termMap* has been created, every worker that handles its posting list uses it. Note that since every posting list is accessed by a single worker at any given time, no synchronization is required.

4.4 Analysis

Sparta accesses posting lists in the same manner as NRA does, and stops only when NRA’s stopping conditions hold. Thus, like NRA, its exact version ($\Delta = \infty$) is safe, and returns the top-k results.

In terms of performance, NRA was shown to be *instance-optimal* when random access is impossible [?], namely, its number of accesses to posting list entries is asymptotically close the optimum for every problem instance. This property holds for pNRA as long as the rates in which different threads access different posting lists are within constant multiples of each other [?], because in this case a thread that “runs ahead” without knowing it should stop only accesses a constant factor more entries than the algorithm needs to, which the asymptotic analysis ignores. Sparta further differs from pNRA in deferring updates to UB until the end of the segment, which may further delay stopping for *segSize* additional posting list accesses. Since *segSize* is constant, Sparta is asymptotically instance-optimal under the same assumptions as pNRA.

5 Evaluation

We compare the performance of Sparta to the following multiprocessor algorithms:

pBMW [?] – parallel BMW, the best-in-class parallelization of a document-order top-k algorithm;

pJASS [?] – a recent parallelization of the JASS [?] score-order retrieval algorithm;
pRA – a parallel implementation of RA;
sNRA – a shared-nothing parallelization of NRA; and
pNRA – a naïve shared-state parallel implementation of NRA.

Although our primary focus is on approximate algorithms, for completeness, we also experiment with their exact counterparts.

In order to crystallize the comparison among the core algorithms, we abstract away other contributors to the wall-clock latency, e.g., index compression. As recent studies show, given state-of-the-art compression techniques, the impact of decompression on end-to-end performance is marginal (e.g., up to 6% with QMX-D4 compression [?]).

We focus on disk-resident search indexes. We also experimented with RAM-resident indexes, and in all cases, all algorithms except pRA got similar results, which is not surprising given that the algorithms traverse posting lists sequentially. These results are omitted.

In what follows, Section 5.1 describes the experiment setup, Section 5.2 explains how the algorithms are implemented, and Section 5.3 presents our results.

5.1 Experiment Setup

We study the algorithms in terms of query latency and throughput attainable at a single multi-core server. We use mid-tier industry-standard hardware – a 12-core Intel Xeon E5620 with 24GB RAM and 1TB SSD drive.

The benchmarking environment and the algorithms are implemented in Java. A *benchmark driver* draws queries from an input queue and submits them to the algorithm being tested, which uses a thread pool for intra-query parallelism. The driver controls the pool size. When testing latency, the entire thread pool is used by a single query. In the throughput evaluation mode, queries are scheduled first-come-first-served, and a new query is scheduled for execution (i.e., assigned threads) once there are idle threads with no outstanding work from currently executing queries. All queries scheduled for execution equally share the thread pool.

In all experiments, the appropriate index (either in document order or in score order) is pre-built offline and stored on disk uncompressed as a collection of binary files, each storing a shard of data partitioned by term. The benchmark environment memory-maps the contents of these files via the MappedByteBuffer API [?]. Prior to each experiment, we flush the file system’s page cache so all pages are physically read from disk during the experiment.

We experiment with the TREC dataset, Category B (ClueWeb09B), which is widely used for information retrieval research [?]. This dataset includes approximately 50M web documents and takes up roughly 30GB of original content, uncompressed.

We use the popular Lucene open-source search engine [?] for preprocessing the index to generate posting lists; this includes text tokenization, posting list maintenance, and term statistics retrieval. We score documents using a standard tf-idf score function with document length normalization [?].

We draw queries from the public AOL search log [?]. For each number of terms from 1 to 12, we independently sample 100 queries of this length uniformly at random

from the AOL log. We also experimented with a query log of another commercial web search engine; this experiment produced statistically similar results, and so we omit them here.

We use $k = 1000$, based on the assumption that simple tf-idf retrieval is the first phase of multi-stage ranking, which may require large values of k for effectiveness [?]. (Indeed, Crane et al. [?] report that the classical AP and RBP quality metrics are close to optimal with this choice of k , for multiple datasets.) Experiments with $k = 100$ produced qualitatively similar results.

5.2 Implementation

Posting lists are stored as contiguous uncompressed arrays; pRA also stores its secondary index (document id to position in the posting list mapping) in the same form. Term scores (namely, tf-idf) are stored in the posting lists as integers, scaled by 10^6 and rounded as in [?]. Using integer arithmetic instead of floating-point significantly speeds up document evaluation.

The specific algorithm implementations are as follows.

State-of-the-art parallel search algorithms

pBMW Our implementation of pBMW closely follows the description in [?]. The algorithm partitions the execution of the sequential BMW [?] among multiple threads. Each thread handles a distinct subset of documents, and computes a local top-k result. The algorithm then merges the partial results to obtain the final top-k.

Similarly to Sparta, pBMW’s threads obtain jobs from a common job queue. Here, a job defines a range of document ids to scan. We set the number of jobs to be twice the number of worker threads, and assign equal-size ranges to all threads. This partition results in well-balanced executions in which the whole worker pool is utilized most of the time.

Each thread maintains a thread-local heap with the current top-k documents. (We also experimented with a shared heap and got inferior results; a similar finding was reported in [?].) Similarly, each thread T maintains a local threshold Θ_T for filtering heap insertions; Θ_T is *at least* the lowest score in the local heap, but may be higher due to the progress of other threads. Thread T periodically compares Θ to its local Θ_T and promotes the smaller of the two to $\max(\Theta_T, \Theta)$. This way, slower workers catch up with faster ones.

pBMW splits posting list segments into blocks, and uses block-level statistics to prune the search [?]. We experimented with multiple block sizes and selected 64, which yielded the best performance. The approximate version’s pruning aggressiveness is controlled by a parameter $f \geq 1$, which multiplies Θ to relax the threshold for score upper bounds [?]. For $f = 1$, the algorithm is exact.

pJASS Our implementation of pJASS follows the description in [?]. It traverses all posting lists in parallel, in score order, and accumulates the encountered scores per-document in *docMap*. Each document is protected by a lock, and a thread that encounters a document locks it, adds the partial score from the term it traversed, and then

Sparta	pNRA	sNRA	pRA	pBMW	pJASS
860	13,291	5,553	480	750	54,343

Table 2. Average query latency (in ms) of 12-term queries with exact algorithms using 12 threads. None of the algorithms meets real-time SLAs.

unlocks it. The algorithm stops after a predefined fraction of postings (denoted p) are scanned.

Parallel TA variants

sNRA *sNRA* is a shared-nothing parallelization of NRA, where the index is partitioned to 12 shards by document id. Each thread finds the top- k documents in its shard by running NRA independently with thread-local data structures. When all threads complete, their lists are merged and the global top- k documents are kept.

pNRA *pNRA* is a naïve shared-state parallelization of NRA that does not employ Sparta’s optimizations. Namely, it uses a shared document map, which it does not clean, and it updates the term upper bounds upon every document evaluation. As in Sparta, a dedicated task checks the stopping condition. (Distributed stopping detection yielded worse results).

pRA Our implementation of *pRA* maintains its results in a shared heap (experiments did not show any advantage to using local heaps). Note that the algorithm’s multiple worker threads may encounter postings of the same document independently, and consequently score that document and try to insert it into the heap multiple times. The implementation guarantees the uniqueness of insertion (only the first one takes effect).

Since RA’s stopping detection is lightweight, we do not dedicate a task to it. Instead, all workers check the *UBStop* condition and monitor the time elapsed since the last heap update, and notify each other if they decide to stop.

5.3 Results

For an algorithm A , we consider the following variants: exact (denoted A -exact), high-recall (denoted A -high), and low-recall (denoted A -low). Note that the approximate algorithms’ parameters (Δ , f , and p) affect the recall but do not directly control it; our high recall instances are ones that empirically achieve a recall of 96% or higher.

Exact Algorithms Our first experiment shows that none of the exact algorithms meets a real-time SLA for verbose queries. Table 2 depicts the mean processing latencies of 12-term queries with 12 worker threads (i.e., a single query fully exploits the multi-core CPU).

We revisit the algorithms’ execution dynamics – namely, how fast they accrue their results – as we study approximate algorithms in the sequel.

Approximate Algorithms With exact instances of all algorithms failing to match real-time requirements, we turn to focus on the approximate instances. We parameterize Sparta, pRA, pNRA, and sNRA with $\Delta = 10$ ms. This yields high recall in all four algorithms, and satisfactory performance (meeting SLA requirements) in Sparta. We instantiate pBMW with $f = 5$ for high recall and $f = 10$ for low recall. Finally, pJASS is instantiated with $p = 0.005$ for low recall and $p = 0.02$ for high recall (using $p = 0.1$, as suggested, e.g., in [?], produced unacceptably high latencies).

Accuracy. Table 3 depicts the empirical accuracy results for 12-term queries.

Sparta-high	pRA-high	pNRA-high	sNRA-high	pBMW-high	pBMW-low	pJASS-high	pJASS-low
97.5%	98.5%	98.5%	99%	97.5%	80%	96%	93%

Table 3. Recall of approximate algorithms for 12-term queries.

Latency. Figure 3 depicts the scaling of single-query latency of the approximate algorithms as the number of terms scales from 1 to 12. The number of workers in each test is equal to the number of terms (for Sparta, pRA, and pNRA, this allows maximal parallelism).

Figures ?? and ?? present, respectively, the mean and 95% latencies of queries in the high-recall algorithms. The latter captures the so-called “tail latency” of the slowest 5% of the queries. Figures ?? and ?? present the same for Sparta-high and the low-recall variants of the algorithms.

When compared to other high recall algorithms, Sparta outperforms all of them in terms of both average latency and the 95th percentile, on all query lengths. The margin is noticeable for verbose queries (5+ terms). Sparta’s average query latency scales perfectly, remaining below 180 ms for all query lengths. Its 95% latency is below 500 ms.

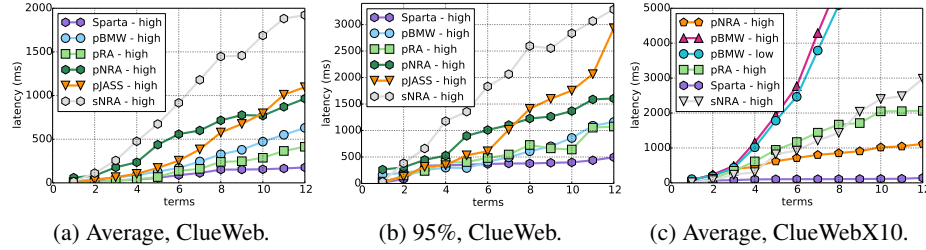


Fig. 3. Top- k ($k = 1000$) query latency scaling (average and 95%) with the number of query terms. The intra-query parallelism is equal to the number of terms in all algorithms.

In contrast, pRA exhibits much weaker scalability. For example, for 12-term queries it is more than 2x slower than Sparta. We explain this by the cost of evaluating the complete document scores in pRA, which requires intensive access to the secondary index. This translates to volumes of random I/O that cannot be sustained even with modern

SSD hardware. Note that this trend is the reverse of the one observed for pRA-exact, which outperforms Sparta-exact (Table 2). That is, Sparta spends much more work than pRA in order to collect the remaining 2.5% of the exact result set. We revisit this phenomenon below.

By sacrificing recall, pBMW and pJASS manage to improve performance and bring their tail latencies for short queries below Sparta’s tail latency with high recall. Nevertheless, Sparta’s average latency is still smaller than theirs for all query lengths, as is its tail latency for long queries. For example, for 12-term queries, pBMW-high completes within 630 ms on average, pBMW-low, which sacrifices 20% of the recall for performance, only succeeds to improve this latency by 15%, pJASS-high completes after more than a second on average, and pJASS-low improves this to 292 ms on average.

The shared-nothing and unoptimized parallelizations of NRA are weaker than all other alternatives: pNRA’s average latency for 12-term queries is 1 second, and sNRA’s is 1.7 seconds. These results emphasize the necessity of sharing information among threads (unlike sNRA) on the one hand, and the importance of Sparta’s locality optimizations, (which are missing in pNRA), on the other. Specifically, the background cleaning and local copies of *docMap* and the lazy updates of *UB* allow Sparta to benefit from local access to data that resides in hardware caches. We omit pNRA and sNRA from further discussion.

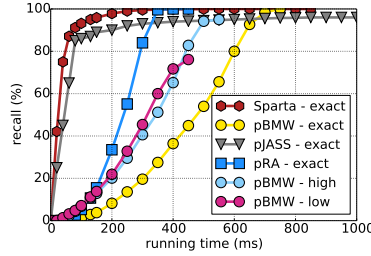


Fig. 4. Recall dynamics with elapsed time, for 12-term queries, with 12 worker threads.

Recall dynamics. In order to understand how the top-k results get accrued by the different algorithms, we zoom in on the dynamics of query recall over the running time. We focus on 12-term queries in a 12-worker configuration. The results are presented in Figure 4. Because the approximate versions of Sparta, pRA, and pJASS are identical to the respective exact versions until they stop, we show the dynamics of the exact versions only. The same is not true for pBMW, where f impacts the algorithm’s results from the outset. Hence, we show the dynamics of all three instances of pBMW.

We see that Sparta’s recall growth is the fastest. For instance, it surpasses 80% recall in less than 50 ms, and 90% recall in less than 100. But over time, its returns diminish, and most of the work becomes unproductive. Whereas pRA takes much longer to converge because it needs to fully score each encountered document, its concluding phase is faster because most relevant documents have complete scores. pBMW scans the postings in the order of document ids, which is unrelated to document scores, and

hence accumulates the true hits at a near-linear rate. Obviously, the convergence rate of pBMW-high and pBMW-low is faster than that of pBMW-exact. The first two accrue results at similar rates until pBMW-low stops at approximately 80% recall. pJASS’s behavior is similar to Sparta’s, but it is a bit slower and fails to reach 100% recall within a minute of execution.

Parallelism. We next consider 12-term queries with a number of threads varying from 1 to 12. The average latencies appear in Figure 5. The left-most data point in each curve corresponds to the performance of the respective sequential algorithm.

Sparta requires some level of parallelism in order to achieve real-time speed – its sequential latency is 840 ms, which is above typical SLA requirements. Most of the gain is achieved at low-parallelism levels (2 threads suffice). The same is true for pJASS, which hardly improves when afforded more threads, due to unequal thread workloads. On the other hand, for pBMW, much higher parallelism is essential – its latency is inversely proportional to the number of threads. Thus, Sparta is not only faster than pBMW, but also requires less resources, which benefits throughput as we next show.

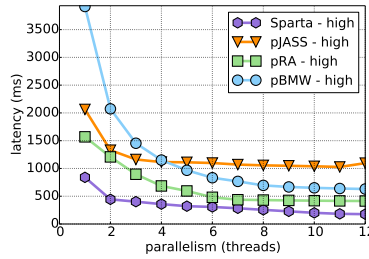


Fig. 5. Average query latency scaling with intra-query parallelism, for 12-term queries, high recall algorithms.

Throughput. Finally, we compare the throughput (in queries per second) provided by the different algorithms. To this end, we generate a workload with the query size distribution reported in [?], where the average query length is 4.2 (std: 2.96), and more than 5% of the queries are 10 terms or longer. The queries are generated as follows: we first sample a query length ℓ from the distribution in [?], and then select a query uniformly at random among all the length- ℓ queries in the complete set of 1200 AOL queries.

Table 4 depicts the results of running this query mix on a shared worker pool of 12 threads. Here too, Sparta improves over its competitors. Its pronounced advantage over pBMW is thanks to a combination of Sparta’s speed and lower resource utilization.

6 Related Work

Verbose queries challenge standard top-k processing techniques in terms of runtime latency. Huston and Croft [?] evaluated several sequential query processing techniques for verbose queries, concluding that the most effective one is to simply reduce the length

Sparta-high	pRA-high	pBMW-high	pJASS-high
12.5	10.9	5.95	10.8

Table 4. Average throughput (in queries per second) of the approximate algorithms on a query distribution measured for voice queries in production.

of the query by omitting stop-words or “stop-structure” expressions. In this work we ignore the query pre-processing phase and consider the query as a bag of words given after textual analysis.

Crane et al. [?] showed that document-order algorithms are susceptible to tail queries that may take orders of magnitude longer than the median query; approximate query evaluation in WAND and BMW does not significantly reduce the variance. Moreover, they showed that score-order algorithms are less sensitive to tail queries due to their effective early termination capability.

Mackenzie et al. [?] suggest pJASS— a parallel version of the score-order algorithm JASS [?]. JASS’s virtue is its simplicity – it performs very little processing per-posting. In order to achieve early termination, JASS applies a heuristic to limit the fraction of processed terms. Similarly to Sparta, pJASS focuses on approximate top-k retrieval (its exact variant is inefficient). We show that Sparta achieves better performance and better recall thanks to multiple optimizations – e.g., continuous result-set pruning that leads to better locality, fine-tuned termination, and careful synchronization. Sparta’s additional advantage is its small RAM footprint, whereas pJASS intentionally avoids pruning and maintains a huge in-memory document map throughout the query evaluation.

Some previous works, e.g., [?,?], studied parallel computation of conjunctive queries via posting list intersection. Note that this problem is different from (and easier than) the problem considered in this paper, where a top-scored document does not necessarily include all query terms.

Other works [?,?] have parallelized state-of-the-art sequential algorithms like WAND and BMW by sharding the document space, computing the top-k in each shard independently, and finally merging the results. Implementations differ in whether threads share a common heap and threshold Θ or not. A global threshold is tighter than the threads’ local thresholds, hence less work is done by each of the threads as more documents can be safely skipped. On the other hand, additional overhead is induced by the synchronization (e.g., using locks) needed to guarantee exclusive updates of the shared heap. Experimental results [?] have shown the superiority of the local-heap approach. The pBMW implementation used in our experiments follows this approach, but periodically shares the Θ values among the threads for improved performance.

Jeon et al. [?] presented an adaptive resource management algorithm that chooses the degree of parallelism at runtime for each query, based on predicting high-latency queries. Such efforts are orthogonal to the performance improvement we achieve via parallelization of (long) queries.

Other works [?,?] have explored using GPU hardware for information retrieval; [?] focused on adaptively choosing whether to use a CPU or a GPU based on the query’s difficulty, and [?] focused on optimizing throughput rather than latency. In contrast, our work leverages standard server-grade multi-threaded CPUs.

The Threshold Algorithm and its variants [?,?] have been extensively studied by the database community, and have been applied in many relational database systems (for a comprehensive survey see [?]). Mamoulis et al. [?] observed two main phases during NRA processing – the “growing phase”, where the candidate list grows, and the “shrinking phase” where no new documents can end up in the top-k results, after the first stopping condition is met. They used different data structures for the two phases in order to minimize the number of accesses and the memory requirements. Gursky et al. [?] also noticed the bottleneck in NRA computation derived from NRA’s needs to maintain an extremely large number of partially scored candidates. They proposed several optimization methods for candidate list maintenance to speed up the search. One of their suggested approaches is to periodically remove irrelevant candidates from the candidate list, which we also do in Sparta.

Yuan et al. [?] observed that the number of accesses to the sorted lists by NRA could be further reduced by selectively performing the sorted accesses to the different lists (instead of in parallel). They proposed a selection policy that prioritizes the accesses to the sorted lists and cuts down unnecessary accesses. They showed significant cutoff in the number of accesses with respect to the original NRA. However, as the authors pointed out, the effectiveness of this approach in terms of run-time latency still has to be explored.

In the IR setting, a few works [?,?,?] have experimented with TA on web data using standard IR metrics. Bast et al. [?] optimized the TA scheduling method based on a cost model for sequential and random accesses. Theobald et al. [?] extended TA for XML query languages. Another work by Theobald et al. [?] introduced an approximate TA algorithm based on probabilistic arguments: When scanning the posting lists in descending order of local scores, various forms of derived bounds are employed to predict when it is safe, with high probability, to skip candidate items hence trading off accuracy for sorted access. Applying similar probabilistic pruning rules for Sparta may prove beneficial and is left for future work.

7 Conclusions

We presented Sparta – a practical algorithm to provide approximate top-k retrieval of verbose queries within interactive latency bounds. Sparta leverages the efficiency and early-stopping properties of the seminal Threshold Algorithm (TA). It forgoes the need for random access and duplicate indices by using the “lazy” scoring approach of TA’s NRA variant. We parallelized the algorithm on shared-memory multi-core hardware while optimizing memory footprints, memory access patterns, inter-thread data sharing, and synchronization.

Sparta yields sub-180 ms average latencies on standard hardware for queries of up to 12 terms when applied to a 50M-document dataset, and can therefore support modern search experiences – which induce long queries – within real-time SLA requirements. It does so while producing a highly accurate approximation of the exact results (a recall of above 97.5%). For comparison, its state-of-the-art parallel competitors, pBMW and pJASS, required 640 ms and above 1 second, respectively, to provide similar accuracy.

Algorithm 1 Sparta algorithm.

```

1: for  $i = 1$  to  $m$  do  $\triangleright$  processing  $m$ -term query
2:   add  $\text{PROCESSTERM}(i)$  to job queue
3:   spawn up to  $m$  threads to run jobs from queue
4:   wait until  $UBStop$   $\triangleright$  all candidates are in  $docMap$ 
5:   add  $\text{CLEANER}()$  to job queue
6:   wait until done
7:   return  $docHeap$ 

8: procedure  $\text{PROCESSTERM}(i)$ 
9:   if  $UBStop \wedge |docMap| < \Phi \wedge termMap[i] = docMap$  then
 $\triangleright docMap$  is shrinking and small – create a local copy for term  $i$ 
10:     $termMap[i] \leftarrow$  new hash map
11:    for all  $D \in docMap$  s.t.  $D.score[i] = 0$  do
12:      add  $D$  to  $termMap[i]$ 
13:    for  $j = 1$  to  $segSize$  do
14:      if done then return
15:       $\langle id, score \rangle \leftarrow$  next entry in  $i$ th posting list
16:       $D \leftarrow termMap[i](id)$ 
17:      if  $D = \perp$  then  $\triangleright$  document missing
18:      if  $\neg UBStop$  then  $\triangleright$  hash incomplete
19:      create new document object  $D$ 
20:      add  $D$  to  $termMap[i](id)$ 
21:      else continue
22:       $D.score[i] \leftarrow$  score  $\triangleright$  update term score
23:      if  $\sum_{j=1}^m D.score[j] > \Theta$  then
24:         $\text{UPDATE\_HEAP}(D)$ 
25:       $UB[i] \leftarrow$  score  $\triangleright$  update term's upper bound

26:   add  $\text{PROCESSTERM}(i)$  to job queue

27: procedure  $\text{UPDATE\_HEAP}(D)$ 
28:   lock  $docHeap$ 
29:   if  $D \notin docHeap$  then
30:     insert  $D$  to  $docHeap$ 
31:     for all  $d \in docHeap$  do
32:        $d.LB \leftarrow \sum_{j=1}^m d.score[j]$ 
33:       move  $d$  to correct place in heap

34:   if  $|docHeap| > k$  then
35:     remove lowest scored doc
36:   if  $|docHeap| = k$  then  $\triangleright$  set  $\Theta$  to  $k^{th}$  lowest score
37:      $\Theta \leftarrow$  lowest score in  $docHeap$ 
38:      $HeapUpdateTime \leftarrow$  current time
39:   unlock  $docHeap$ 

40: procedure  $\text{CLEANER}$ 
41:   if  $|docMap| > \Phi$  then  $\triangleright$  shrink  $docMap$ 
42:      $tmpDocMap \leftarrow$  new hash map
43:     for every doc  $D \in docMap$  do
44:       if  $UB(D) > \Theta \vee D \in docHeap$  then
45:         add  $D$  to  $tmpDocMap$   $\triangleright D$  is still relevant
46:       replace  $docMap$  by  $tmpDocMap$ 

 $\triangleright$  check algorithm's stopping conditions; in exact version  $\Delta = \infty$ 
47:   if  $(|docMap| = |docHeap|) \vee HeapUpdateTime + \Delta < \text{now}$  then
48:     done  $\leftarrow$  true
49:   else add  $\text{CLEANER}()$  to job queue

```
