





Deterministic Cache-Oblivious Funnelselect

Gerth Stølting Brodal  

Aarhus University, Denmark

Sebastian Wild  

University of Liverpool, UK

Abstract

In the multiple-selection problem one is given an unsorted array S of N elements and an array of q query ranks $r_1 < \dots < r_q$, and the task is to return, in sorted order, the q elements in S of rank r_1, \dots, r_q , respectively. The asymptotic deterministic comparison complexity of the problem was settled by Dobkin and Munro [JACM 1981]. In the I/O model an optimal I/O complexity was achieved by Hu *et al.* [SPAA 2014]. Recently [ESA 2023], we presented a *cache-oblivious* algorithm with matching I/O complexity, named *funnelselect*, since it heavily borrows ideas from the cache-oblivious sorting algorithm *funnelsort* from the seminal paper by Frigo, Leiserson, Prokop and Ramachandran [FOCS 1999]. Funnelselect is inherently randomized as it relies on sampling for cheaply finding many good pivots. In this paper we present *deterministic funnelselect*, achieving the same optimal I/O complexity cache-obliviously without randomization. Our new algorithm essentially replaces a single (in expectation) reversed-funnel computation using random pivots by a recursive algorithm using multiple reversed-funnel computations. To meet the I/O bound, this requires a carefully chosen subproblem size based on the entropy of the sequence of query ranks; deterministic funnelselect thus raises distinct technical challenges not met by randomized funnelselect. The resulting worst-case I/O bound is $O(\sum_{i=1}^{q+1} \frac{\Delta_i}{B} \cdot \log_{M/B} \frac{N}{\Delta_i} + \frac{N}{B})$, where B is the external memory block size, $M \geq B^{1+\varepsilon}$ is the internal memory size, for some constant $\varepsilon > 0$, and $\Delta_i = r_i - r_{i-1}$ (assuming $r_0 = 0$ and $r_{q+1} = N + 1$).

2012 ACM Subject Classification Theory of computation \rightarrow Design and analysis of algorithms

Keywords and phrases Multiple selection, cache-oblivious algorithm, entropy bounds

Digital Object Identifier 10.4230/LIPIcs.SWAT.2024.17

Funding Gerth Stølting Brodal: Independent Research Fund Denmark, grant 9131-00113B.

Sebastian Wild: Engineering and Physical Sciences Research Council grant EP/X039447/1.

1 Introduction

We present the first optimal deterministic cache-oblivious algorithm for the multiple-selection problem. In the multiple-selection problem one is given an unsorted array S of N elements and an array R of q query ranks in increasing order $r_1 < \dots < r_q$, and the task is to return, in sorted order, the q elements of S of rank r_1, \dots, r_q , respectively; (see Figure 1 for an example).

On top of immediate applications, like finding a set of quantiles in a numerical data set, the multiple-selection problem is of interest as it gives a natural common generalization of (single) selection by rank (using a single query rank $r_1 = r$) and fully sorting an array (corresponding

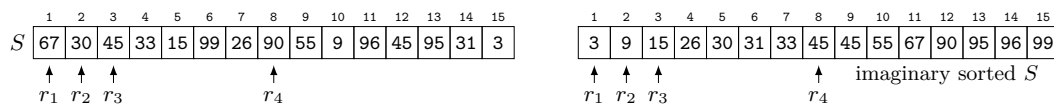


Figure 1 Example input with $N = 15$, $q = 4$ and $R[1..q] = [1, 2, 3, 8]$. The expected output 3, 9, 15, 45 is obvious from the sorted array (right). (The sorted array is for illustration only; the goal of efficient multiple-selection algorithms is to avoid ever fully sorting the input.)



© Gerth Stølting Brodal and Sebastian Wild;
licensed under Creative Commons License CC-BY 4.0

19th Scandinavian Symposium and Workshops on Algorithm Theory (SWAT 2024).

Editor: Hans L. Bodlaender; Article No. 17; pp. 17:1–17:12

Leibniz International Proceedings in Informatics



LIPICs Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

to selecting every index as a query rank, i.e., $q = N$ and $r_i = i$ for $i = 1, \dots, N$). It thus allows us to quantitatively study the transition between these two foundational problems, which are of different complexity and each have their distinct set of algorithms. For example, the behavior of selection and sorting with respect to external memory is quite different: For single selection, the textbook median-of-medians algorithm [4] simultaneously works with optimal cost in internal memory, external memory, and the cache-oblivious model (models are defined below). For sorting, by contrast, the introduction of each model required a substantially modified algorithm to achieve optimal costs: Standard binary mergesort is optimal in internal memory, but requires $\approx M/B$ -way merging to be optimal in external memory, where M is the internal memory size and B the external memory block size, measured in elements [1]; achieving the same cache obliviously, i.e., without knowledge of B and M , requires the judiciously chosen buffer sizes from the recursive constructions of funnelsort [12].

Since multiple selection simultaneously generalizes both problems, it is not surprising that also here subsequent refinements were necessary going from internal to external to cache-oblivious; the most recent result being our algorithm funnelselect [7]. However, all algorithms mentioned above for single selection and sorting are *deterministic*. By contrast, funnelselect is inherently relying on randomization and known deterministic external-memory algorithms [2, 15] are crucially relying on the knowledge of M and B . Prior to this work it thus remained open whether a single deterministic cache-oblivious algorithm exists that smoothly interpolates between selection and sorting without having to resort to randomization.

In this paper, we answer this question in the affirmative. Our algorithm *deterministic funnelselect* draws on techniques from cache-oblivious sorting (funnelsort) and existing multiple-selection algorithms, but it follows a rather different approach to our earlier randomized algorithm [7] and previous (cache-conscious) external-memory algorithms. A detailed comparison is given below.

1.1 Model of computation and previous work

Our results are in the cache-oblivious model of Frigo, Leiserson, Prokop and Ramachandran [13], a hierarchical-memory model with an infinite external memory and an internal memory of capacity M elements, where data is transferred between internal and external memory in blocks of B consecutive elements. Algorithms are compared by their I/O cost, i.e., the number of block transfers or *I/Os* (input/output operations). This is similar to the external-memory model by Aggarwal and Vitter [1]. Crucially, in the cache-oblivious model, no variables are allowed to depend on the model parameters M and B and *I/Os* are assumed to be performed automatically by an optimal (offline) paging algorithm. The algorithms are analyzed with respect to the parameters M and B . Cache-oblivious algorithms hence work for any parameters M and B , and they even adapt to multi-level memory hierarchies (under certain conditions [13]).

The multiple-selection problem was first formally addressed by Chambers [8], who considered it a generalization of quickselect [14]. Prodingler [17] proved that Chambers' algorithm achieves an optimal *expected* running time up to constant factors: $O(\mathcal{B} + N)$, where $\mathcal{B} = \sum_{i=1}^{q+1} \Delta_i \lg \frac{N}{\Delta_i}$ with $\Delta_i = r_i - r_{i-1}$, for $1 \leq i \leq q+1$, assuming $r_0 = 0$ and $r_{q+1} = N+1$, and \lg denoting the binary logarithm. We call \mathcal{B} the (*query-rank*) *entropy* of the sequence of query ranks [2]. It should be noted that $\mathcal{B} + N = O(N(1 + \lg q))$, but the latter bound does not take the location of query ranks into account; for example, if $q = \Theta(\sqrt{n})$ queries are in a range of size $O(N/\lg N)$, i.e., $r_q - r_1 = O(N/\lg N)$, then the entropy bound is $O(N)$ whereas the latter $N(1 + \lg q) = \Theta(N \lg N)$.

■ **Table 1** Algorithms for selection and multiple selection. CO = cache-oblivious, \mathbb{E} = expected, wc = worst-case bounds. Note that Barbay *et al.* assume a tall cache $M \geq B^{1+\varepsilon}$, whereas Hu *et al.* do not.

Reference		Comparisons	I/Os	Comments
Single selection				
Hoare [14]	\mathbb{E}	$2 \ln 2\mathcal{B} + 2N + o(N)$	$O(N/B)$	CO, randomized
Floyd & Rivest [11]	\mathbb{E}	$N + \min\{r, N-r\} + o(N)$	$O(N/B)$	CO, randomized
Blum <i>et al.</i> [4]	wc	$5.4305N$	$O(N/B)$	CO, deterministic
Schönhage <i>et al.</i> [18]	wc	$3N + o(N)$?	CO, median, deterministic
Dor & Zwick [10]	wc	$2.95N + o(N)$?	CO, median, deterministic
Multiple selection				
Chambers [8, 17]	\mathbb{E}	$2 \ln 2\mathcal{B} + O(N)$	$O((\mathcal{B} + N)/B)$	CO, randomized
Dobkin & Munro [9]	wc	$3\mathcal{B} + O(N)$	$O((\mathcal{B} + N)/B)$	CO, deterministic
Kaligosi <i>et al.</i> [16]	wc	$\mathcal{B} + o(\mathcal{B}) + O(N)$	$O((\mathcal{B} + N)/B)$	CO, deterministic
Hu <i>et al.</i> [15]	wc	$O(N \lg(q))$	$O(N/B \cdot \log_{M/B}(q/B))$	deterministic
	wc	$O(\mathcal{B} + N)$	$O(\mathcal{B}_{I/O} + N/B)$	(from closer analysis)
Barbay <i>et al.</i> [2]	wc	$\mathcal{B} + o(\mathcal{B}) + O(N)$	$O(\mathcal{B}_{I/O} + N/B)$	online, determ., $M \geq B^{1+\varepsilon}$
Brodal & Wild [7]	\mathbb{E}	$O(\mathcal{B} + N)$	$O(\mathcal{B}_{I/O} + N/B)$	CO, randomized, $M \geq B^{1+\varepsilon}$
<i>This paper</i>	wc	$O(\mathcal{B} + N)$	$O(\mathcal{B}_{I/O} + N/B)$	CO, deterministic, $M \geq B^{1+\varepsilon}$

Dobkin and Munro [9] showed that $\mathcal{B} - O(N)$ comparisons are necessary to find all ranks r_1, \dots, r_q (in the worst case). Deterministic algorithms with that same $O(\mathcal{B} + N)$ running time are also known [9, 16], but as for single selection, the deterministic algorithms were presented later than the randomized algorithms and require more sophistication. Multiple selection in external-memory was studied by Hu *et al.* [15] and Barbay *et al.* [2]. Their algorithms have an I/O cost of $O(\mathcal{B}_{I/O} + \frac{N}{B})$, where the “I/O entropy” $\mathcal{B}_{I/O} = \frac{\mathcal{B}}{B \lg(M/B)}$. An I/O cost of $\Omega(\mathcal{B}_{I/O}) - O(\frac{N}{B})$ is known to be necessary [2, 7]. A more comprehensive history of the multiple-selection problem appears in [7]; Table 1 gives an overview.

Algorithms designed for a single level memory model (RAM model) can be analyzed in the cache-oblivious model. Analyzing existing time- and comparison-optimal multiple-selection algorithms [8, 9, 16, 17] in the cache-oblivious model, shows that they are not optimal with respect to the number of I/Os performed. The obtained I/O bounds are a factor $\lg(M/B)$ away from being optimal. To get a feeling for the I/O bounds in Table 1, consider for example the case where $B = \lg N$, $M = B^2$, and there are $q = \lg N$ evenly distributed queries ($r_i = \lfloor \frac{iN}{q+1} \rfloor$). In this case $\mathcal{B} = \Theta(N \lg \lg N)$. The cache-oblivious sorting bound (see Lemma 6) is $O(\frac{N}{B} \log_M N) = O(\frac{N}{B} \cdot \frac{\lg N}{\lg \lg N})$, $O((\mathcal{B} + N)/B) = O(\frac{N}{B} \lg \lg N)$ and $O(\mathcal{B}_{I/O} + \frac{N}{B}) = O(\frac{N}{B})$, where $O(\frac{N}{B})$ is the number of I/Os required to scan the input.

1.2 Result

Our main result is the cache-oblivious algorithm *deterministic funnelselect* achieving the following efficiency (see Theorem 10 for the full statement and proof).

► **Theorem 1.** *There exists a deterministic cache-oblivious algorithm solving the multiple-selection problem using $O(\mathcal{B} + N)$ comparisons and $O(\mathcal{B}_{I/O} + \frac{N}{B})$ I/Os in the worst case, assuming a tall cache $M \geq B^{1+\varepsilon}$.*

At the high level, our algorithm uses the standard overall idea of a recursive partitioning algorithm and pruning recursive calls containing no rank queries, an idea dating back to the first algorithm by Chambers [8]. In the cache-aware external-memory model, I/O efficient

algorithms are essentially obtained by replacing binary partitioning (as used in [8]) by an external-memory $\Theta(M/B)$ -way partitioning [2, 15]. Unfortunately, in the cache-oblivious model this is not possible, since the parameters M and B are unknown to the algorithm. To be I/O efficient in the cache oblivious model, both our previous algorithm randomized funnelselect [7] and our new algorithm deterministic funnelselect apply a cache-oblivious multi-way k -partitioner to distribute elements into k buckets given a set of $k - 1$ pivot elements, essentially reversing the computation done by the k -merger used by funnelsort [12]. The k -partitioner is a balanced binary tree of $k - 1$ pipelined binary partitioners.

The key difference between our randomized and deterministic algorithms is that in our randomized algorithm we use a single $N^{\Theta(\epsilon)}$ -way partitioner using randomly selected pivots and truncate work inside the partitioner for subproblems that (with high probability) will not contain any rank queries. This is done by estimating the ranks of the pivots through sampling and pruning subproblems estimated to be sufficiently far from any query ranks. In our deterministic version, we choose k smaller and deterministically compute pivots, such that all elements are pushed all the way down through a k -partitioner without truncation (eliminating the need to know the (approximate) ranks of the pivots before the k -partitioning is finished), while we choose k such that the buckets with unresolved rank queries (that we have to recurse on) in total contain *at most half* of the elements. To compute k , we apply a linear-time *weighted*-median finding algorithm on $\Delta_1, \dots, \Delta_{q+1}$. While randomized funnelselect can handle buckets with unresolved rank queries directly using sorting, deterministic funnelselect needs to recursively perform multiple-selection on the buckets to achieve the desired I/O performance.

2 Preliminaries

Throughout the paper we assume that the input to a multiple-selection algorithm is given as two arrays $S[1..N]$ and $R[1..q]$, where S is an unsorted array of N elements from a totally ordered universe, and R is a sorted array r_1, \dots, r_q of q distinct query ranks, where $1 \leq r_1 < \dots < r_q \leq N$. The array S is allowed to contain duplicate elements. Our task is to produce/report an array of the q order statistics $S_{(r_1)}, \dots, S_{(r_q)}$, where $S_{(r)}$ is the r th smallest element in S , i.e., the element at index r in an array storing S after sorting it.

Our new deterministic cache-oblivious multiple-selection algorithm makes use of the following three existing cache-oblivious results for single selection, weighted selection, sorting, and multi-way partitioning.

► **Lemma 2** (Blum, Floyd, Pratt, Rivest, Tarjan [4, Theorem 1]). *Selecting the k -th smallest element in an unsorted array of N elements can be done with $O(N)$ comparisons and $O(1 + \frac{N}{B})$ I/Os in the cache-oblivious model.*

► **Remark 3** (Median of medians: I/O cost). Although the original paper by Blum *et al.* [4] predates the cache-oblivious model [12] by decades, analyzing the algorithm in the cache-oblivious model with a stack-oriented memory allocator gives a linear I/O cost, since the algorithm is based on repeatedly scanning geometrically decreasing subproblems.

► **Remark 4** (Median of medians: duplicates). The original algorithm in [4] assumes that all elements are distinct. The algorithm can be extended to handle duplicates (by performing a three-way partition of the elements into those less-than, equal-to, and greater-than a pivot, respectively), and to return a triple S_{\leq}, p, S_{\geq} , that is a partition of S , where p is the element of rank k , S_{\leq} are the elements of rank $1, \dots, k - 1$ in arbitrary order, and S_{\geq} are the elements of rank $k + 1, \dots, |S|$ in arbitrary order (where duplicate elements are assigned consecutive ranks in an arbitrary order).

In the *weighted selection* problem we are given an array of N elements, each with an associated non-negative weight, and a target weight W . The goal is find the smallest k , where the sum of the weights of the k smallest elements is at least W , and to return the k -th smallest element. A linear-time weighted-selection algorithm can be derived from the unweighted selection algorithm by Blum *et al.* [4] (Lemma 2) – as hinted by Shamos in [19] and spelled out in detail by Bleich and Overton [3] – by computing the weighted rank of the pivot. The weighted selection algorithm follows essentially the same recursion as [4], and it similarly follows that it is cache oblivious and performs $O(1 + \frac{N}{B})$ I/Os.

► **Lemma 5** (Bleich, Overton [3]). *Weighted selection in an unsorted array of N weighted elements can be done with $O(N)$ comparisons and $O(1 + \frac{N}{B})$ I/Os in the cache-oblivious model.*

► **Lemma 6** (Frigo, Leiserson, Prokop, Ramachandran [13, Theorem 7], Brodal, Fagerberg [5, Theorem 2]). *Funnelsort sorts an array of N elements using $O(N \lg N)$ comparisons and $O(\frac{N}{B}(1 + \log_M N))$ I/Os in a cache-oblivious model with a tall-cache assumption $M \geq B^{1+\varepsilon}$, for constant $\varepsilon > 0$.*

► **Remark 7** (Tall and taller). The original description of funnelsort by Frigo *et al.* [12] assumed the tall cache assumption $M = \Omega(B^2)$, whereas [5] observed that this could be relaxed to the weaker tall cache assumption $M = \Omega(B^{1+\varepsilon})$. I/O optimality of funnelsort follows from a matching external-memory lower bound by Aggarwal and Vitter [1, Theorem 3.1].

The key innovation in our previous randomized algorithm funnelselect [7] is the *k-partitioner* (Figure 2), a cache-oblivious and I/O-efficient multi-way partitioning algorithm to distribute a batch of elements around $k - 1$ given pivots into k buckets; the precise characteristics are summarized in the following lemma.

► **Lemma 8** (Brodal and Wild [7, Lemma 3]). *Given an unsorted array of $N \geq k^d$ elements and $k - 1$ pivots $P_1 \leq \dots \leq P_{k-1}$, a *k-partitioner* can partition the elements into k buckets S_1, \dots, S_k , such all elements x in bucket S_i satisfy $P_{i-1} \leq x \leq P_i$. The algorithm is cache-oblivious and performs $O(N \lg k)$ comparisons and $O(k + \frac{N}{B}(1 + \log_M k))$ I/Os, provided a tall-cache assumption $M \geq B^{1+\varepsilon}$ and $d \geq \max\{1 + 2/\varepsilon, 2\}$. The working space for the *k-partitioner* (ignoring input and output buffers) is $O(k^{(d+1)/2})$. This is also the time required to construct a *k-partitioner* (again ignoring input and output buffers).*

The *k-partitioners* are structurally similar to the *k-mergers* from funnelsort for merging k runs cache obliviously. In [7] we pipeline the partitioning by essentially reversing the computations done by funnelsort, and replace each binary merging node by a binary partitioning node.

3 Deterministic multiple-selection

In this section we present our deterministic cache-oblivious multiple-selection algorithm that performs optimal $O(B + N)$ comparisons and $O(B_{I/O} + \frac{N}{B})$ I/Os, under a tall-cache assumption $M \geq B^{1+\varepsilon}$. Detailed pseudo-code is given in Algorithm 1 and Algorithm 2, and the basic idea is illustrated in Figure 3.

Given a tall-cache assumption $M \geq B^{1+\varepsilon}$, we let $d = \max\{1 + 2/\varepsilon, 2\}$. The algorithm follows the general idea of making a recursive multi-way partition of the array of elements and to only recurse on subproblems with unresolved rank queries. For two consecutive query ranks r_{i-1} and r_i , we say that the $\Delta_i = r_i - r_{i-1}$ elements of rank $r_{i-1} + 1, \dots, r_i$ are in a

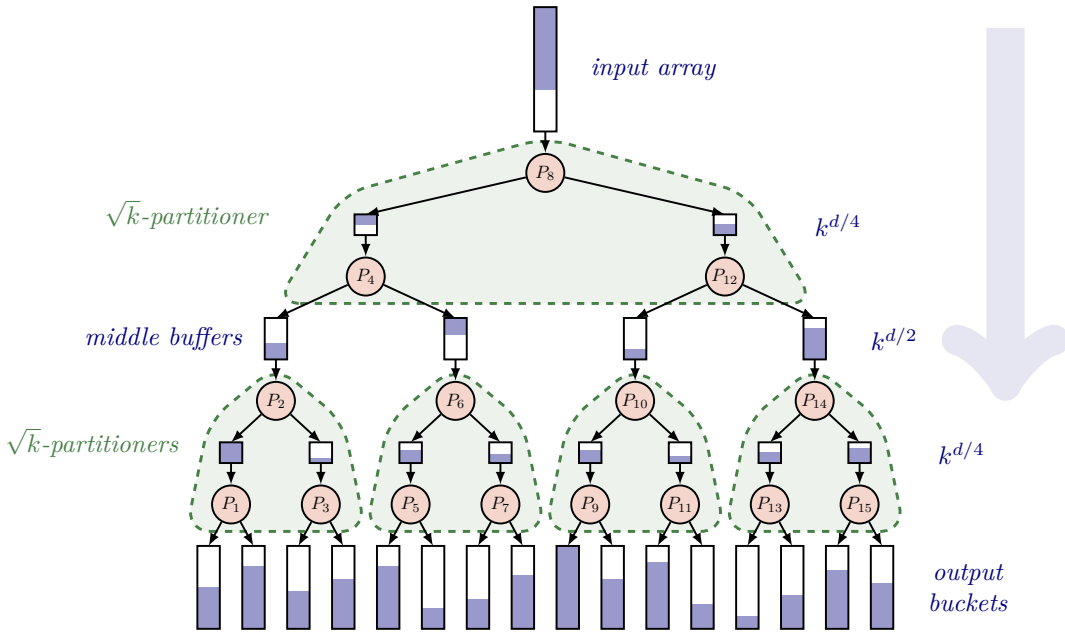


Figure 2 A k -partitioner for $k = 16$ buckets. Content in the buffers is shaded; buffers are filled bottom-to-top; when full, they are flushed and then consumed from the bottom. The figure shows the situation where the input buffer for P_6 is being flushed down to its children (by partitioning elements around pivot P_6). The flush at P_6 was triggered during flushing P_4 's input buffer, which in turn has been called while flushing P_8 (the input). Buffer sizes for the three internal levels are shown next to the buffers. k -partitioners are defined recursively from a \sqrt{k} -partitioner at the top, a collection of \sqrt{k} middle buffers, and \sqrt{k} further \sqrt{k} -partitioners, each partitioning from one middle buffer to \sqrt{k} output buffers. (All sizes here ignore floors and ceilings; for the precise definition valid for all k , see [7].)

gap of size Δ_i . We choose a parameter Δ , such that at least half of the elements are in gaps of size $\leq \Delta$ and simultaneously at least half (rounded down) of the elements are in gaps of size $\geq \Delta$. To compute Δ (Algorithm 1, line 4), we compute $\Delta_i = r_i - r_{i-1}$ by a scan over the query ranks r_1, \dots, r_q (and $r_0 = 0$ and $r_{q+1} = N + 1$), and perform weighted selection (Lemma 5) among $\Delta_1, \dots, \Delta_{q+1}$, where Δ_i has weight $w_i = \Delta_i$, and return the smallest Δ where $\sum_{\Delta_i \leq \Delta} w_i \geq N/2 + 1$.

For the case when Δ is small compared to N (formally, $(2N)^d \geq \Delta^{d+1}$ or $N^{1+\frac{1}{1+\varepsilon}} \geq \Delta^2$), we simply solve the multiple-selection problem by sorting the elements (cache-obliviously using funnelsort [13]), and report the elements with ranks r_1, \dots, r_q by a single scan over the sorted elements. The condition on Δ implies $\mathcal{B}_{I/O} = \Omega(\text{Sort}_{M,B}(N))$, where $\text{Sort}_{M,B}(N) = \Theta\left(\frac{N}{B} \left(1 + \log_{M/B} \frac{N}{B}\right)\right)$ is the number of I/Os required to sort N elements in external memory [1], so this is within a constant factor of the I/O lower bound (detailed analysis in Section 4).

Otherwise, we create a k -partition, where $k = \Theta\left(\frac{N}{\Delta}\right)$ as follows (MULTIPARTITION in Algorithm 2): We repeatedly distribute batches of Δ elements into a set of buckets separated by pivot elements. Initially we have one empty bucket and no pivot. Whenever a bucket reaches size $> \Delta$, the bucket is split into two buckets of size $\leq \Delta$ separated by a new pivot using the (cache-oblivious) linear-time median selection algorithm (Lemma 2). To distribute a batch of elements into the current set of buckets we use a cache-oblivious k -partitioner (Lemma 8, which depends on the tall-cache assumption parameter d) built using the current set of pivots. Note that we need to construct a new k -partitioner after each batch of Δ elements has been distributed, since the number of buckets and pivots can increase. For the

■ **Algorithm 1** Deterministic cache-oblivious multiple-selection.

```

1: procedure DETERMINISTICFUNNELSELECT( $S[1..N]$ ,  $R[1..q]$ )
2:   if  $q > 0$  then
3:      $\Delta_i \leftarrow R[i] - R[i-1]$  for  $i = 1, \dots, q+1$ , assuming  $R[0] = 0$  and  $R[q+1] = N+1$ 
4:      $\Delta \leftarrow \min\{\Delta_i \in \{\Delta_1, \dots, \Delta_{q+1}\} \mid \sum_{j \in \{1, \dots, q+1\}: \Delta_j \leq \Delta_i} \Delta_j \geq N/2 + 1\}$ 
5:     if  $(2N)^d \geq \Delta^{d+1}$  or  $N^{1+\frac{1}{1+\varepsilon}} \geq \Delta^2$  then ▷  $\mathcal{B}_{I/O} = \Omega(\text{Sort}_{M,B}(N))$ 
6:        $S \leftarrow \text{FUNNELSORT}(S)$ 
7:       Report  $S[R[1]], \dots, S[R[q]]$ 
8:     else
9:        $(P_1, \dots, P_{k-1}), (S_1, \dots, S_k) \leftarrow \text{MULTIPARTITION}(S, \Delta)$ 
10:       $\bar{r}_0 \leftarrow 0$ 
11:      for  $i \leftarrow 1, \dots, k$  do
12:         $\bar{r}_i \leftarrow \bar{r}_{i-1} + |S_i| + 1$  ▷  $\bar{r}_i$  is rank of  $P_i$ 
13:         $R_i \leftarrow \{r \mid r \in R \wedge \bar{r}_{i-1} < r < \bar{r}_i\}$  ▷ Rank queries to bucket  $S_i$ 
14:        if  $|R_i| > 0$  then
15:           $r_i^{\max} \leftarrow \max(R_i)$ 
16:           $\bar{S}_i, p_{\max}, S_{\geq} \leftarrow \text{SELECT}(S_i, r_i^{\max} - \bar{r}_{i-1})$ 
17:          if  $|R_i| > 1$  then
18:             $r_i^{\min} \leftarrow \min(R_i)$ 
19:             $S_{\leq}, p_{\min}, \bar{S}_i \leftarrow \text{SELECT}(\bar{S}_i, r_i^{\min} - \bar{r}_{i-1})$ 
20:            Report  $p_{\min}$ 
21:            if  $|R_i| > 2$  then
22:               $\bar{R}_i \leftarrow \{r - r_i^{\min} \mid r \in R_i \setminus \{r_i^{\min}, r_i^{\max}\}\}$ 
23:              DETERMINISTICFUNNELSELECT( $\bar{S}_i, \bar{R}_i$ )
24:            Report  $p_{\max}$ 
25:          if  $\bar{r}_i \in R$  then
26:            Report  $P_i$ 

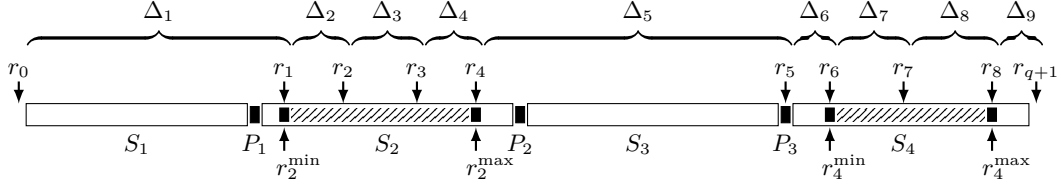
```

■ **Algorithm 2** Given an array S with N elements and a bucket capacity Δ , where $(2N)^{\frac{d}{d+1}} \leq \Delta \leq N$, partition S into k buckets S_1, \dots, S_k separated by $k-1$ pivots P_1, \dots, P_{k-1} , where $\lfloor \frac{\Delta}{2} \rfloor \leq |S_i| \leq \Delta$.

```

1: procedure MULTIPARTITION( $S[1..N]$ ,  $\Delta$ )
2:   Requires  $(2N)^{\frac{d}{d+1}} \leq \Delta \leq N$ 
3:    $k \leftarrow 1$ ,  $S_1 \leftarrow \{\}$  ▷ Initially only one empty bucket and no pivots
4:   for  $i \leftarrow 1$  to  $N$  step  $\Delta$  do
5:      $\bar{S} \leftarrow S[i.. \min(i + \Delta - 1, N)]$  ▷ Next batch to distribute to buckets
6:     Distribute  $\bar{S}$  to buckets  $S_1, \dots, S_k$  using pivots  $P_1, \dots, P_{k-1}$  with a  $k$ -partitioner
7:     while there exists a bucket  $S_j$  with  $|S_j| > \Delta$  do ▷ Split bucket  $S_j$ 
8:        $S_{\leq}, p, S_{\geq} \leftarrow \text{SELECT}(S_j, \lceil |S_j|/2 \rceil)$ 
9:       Rename  $S_{j+1}, \dots, S_k$  to  $S_{j+2}, \dots, S_{k+1}$  and  $P_j, \dots, P_{k-1}$  to  $P_{j+1}, \dots, P_k$ 
10:       $S_j \leftarrow S_{\leq}$ ,  $P_j \leftarrow p$ ,  $S_{j+1} \leftarrow S_{\geq}$ 
11:       $k \leftarrow k + 1$ 
12:   return  $(P_1, \dots, P_{k-1}), (S_1, \dots, S_k)$ 

```



■ **Figure 3** Deterministic multiple selection. The partition of an array S into buckets S_1, \dots, S_4 separated by pivots P_1, \dots, P_3 , and query ranks r_1, \dots, r_8 . In the example the maximum allowed bucket size is $\Delta = \Delta_1$, since $\Delta_1 + \Delta_2 + \Delta_3 + \Delta_4 + \Delta_6 + \Delta_7 + \Delta_8 + \Delta_9 \geq |S|/2 + 1$ and $\Delta_2 + \Delta_3 + \Delta_4 + \Delta_6 + \Delta_7 + \Delta_8 + \Delta_9 < |S|/2 + 1$. Black squares are pivots and the shaded regions in buckets are the subproblems to recurse on.

computation to be I/O efficient, we allocate in memory space for a $\lfloor \frac{2N}{\Delta} \rfloor$ -partitioner followed by space for $\lfloor \frac{2N}{\Delta} \rfloor$ buckets of capacity 2Δ (in the proof of Lemma 9 we argue that the number of buckets created is at most $\frac{2N}{\Delta}$ and each bucket will never exceed 2Δ elements). The space for the partitioner is reused for each new batch, and whenever a bucket is split into two new buckets, one bucket remains in the old bucket's allocated space and the other bucket is placed in next available slot for a bucket. This ensures all buckets are stored consecutively in memory, albeit in arbitrary order.

After having constructed the buckets we compute the ranks of the pivots from the bucket sizes, and consider the buckets with at least one unresolved rank query. If the rank of a pivot coincides with a query rank, we report this pivot just after having considered the preceding bucket. Before recursing on the elements in a bucket, we first find the minimum and maximum query ranks r^{\min} and r^{\max} in the bucket by a scan over the bucket's query ranks, and find and report the corresponding elements in the bucket using linear-time selection (Lemma 2). Finally, we only recurse on the elements between ranks r^{\min} and r^{\max} , provided there are any unresolved rank queries to the bucket. This ensures that when recursing on a subproblem of size \bar{N} , all elements in the subproblem are in gaps of size $< \bar{N}$ in the original input. By reporting the elements at the appropriate times during the recursion, elements will be reported in increasing order.

The partitioning of an array S into buckets is illustrated in Figure 3. The crucial property is that for a gap $\Delta_i \geq \Delta$, the two query ranks r_{i-1} and r_i defining the gap *cannot* be in the same bucket, implying that no element in this gap will be part of a recursive subproblem (see, e.g., gaps Δ_1 and Δ_5 in Figure 3).

Pseudocode for our algorithm is shown in Algorithm 1 and Algorithm 2. We assume $\text{SELECT}(S, k)$ is the deterministic linear-time selection algorithm from Lemma 2, and that it returns a triple S_{\leq}, p, S_{\geq} , that is a partition of S , where p is the element of rank k , S_{\leq} are the elements of rank $1, \dots, k-1$ in arbitrary order, and S_{\geq} the elements of rank $k+1, \dots, |S|$ in arbitrary order.

4 Analysis

We first analyze the number of comparisons and I/Os performed by MULTIPARTITION in Algorithm 2, that deterministically performs a k -way partition of N elements into $k = O(\frac{N}{\Delta})$ buckets separated by $k-1$ pivots, where each bucket has size at most Δ . The following lemma summarizes the precise properties of MULTIPARTITION .

► **Lemma 9.** *For $N \geq \Delta$ and $\Delta^{d+1} \geq (2N)^d$, MULTIPARTITION creates $k \leq \frac{2N}{\Delta}$ buckets and $k-1$ pivots, each bucket has size at most Δ , and performs $O(N \lg k)$ comparisons and $O(k^2 + \frac{N}{B}(1 + \log_M k))$ I/Os.*

Proof. We first bound the sizes of the buckets created by MULTIPARTITION. The algorithm repeatedly distributes batches of at most Δ elements to buckets and splits all overflowing buckets of size $> \Delta$ before considering the next batch. It is an invariant that before distributing a batch, all buckets have size at most Δ . Furthermore, as soon as the first bucket is split, all buckets have size at least $\lfloor \frac{\Delta}{2} \rfloor$, since whenever an overflowing bucket of size $s > \Delta$ is split the new buckets have initial sizes $\lfloor \frac{s-1}{2} \rfloor$ and $\lceil \frac{s-1}{2} \rceil$. Here “ -1 ” is due to one element becomes a pivot. The smallest bucket size is achieved when $s = \Delta + 1$, where the smallest bucket size is $\lfloor \frac{\Delta+1-1}{2} \rfloor = \lfloor \frac{\Delta}{2} \rfloor$. Note that the buckets after the split have size at most Δ , since all buckets had at most Δ elements before the distribution of a batch of at most Δ elements to the buckets, i.e., $s \leq 2\Delta$. To bound the total number of buckets k created, observe that if $\Delta = N$ then no bucket will be split and $k = 1$. Otherwise, $\Delta < N$ and at least two buckets are created, and $k \lfloor \frac{\Delta}{2} \rfloor + k - 1 \leq N$, since all buckets have size at least $\lfloor \frac{\Delta}{2} \rfloor$ and there are $k - 1$ pivots. We have $N \geq k(\frac{\Delta}{2} - \frac{1}{2}) + k - 1 = \frac{k\Delta}{2} + \frac{k}{2} - 1 \geq \frac{k\Delta}{2}$, since $k \geq 2$, i.e., the total number of buckets created $k \leq \frac{2N}{\Delta}$.

To analyze the number of comparisons and I/Os performed, we need to consider the $\lceil \frac{N}{\Delta} \rceil$ distribution steps and at most $\frac{2N}{\Delta} - 1$ bucket splittings. Since each bucket splitting involves at most 2Δ elements, each bucket splitting can be performed cache-obliviously by a linear-time selection algorithm (Lemma 2) using $O(\Delta)$ comparisons and $O(1 + \frac{\Delta}{B})$ I/Os, assuming each bucket is stored in a buffer of 2Δ consecutive memory cells. In total the $k - 1 = \Theta(\frac{N}{\Delta})$ bucket splittings require $O(N)$ comparisons and $O(k + \frac{N}{B})$ I/Os. A k -partitioner for partitioning Δ elements uses $O(\Delta \lg k)$ comparisons and $O(k + \frac{\Delta}{B}(1 + \log_M k))$ I/Os (Lemma 8), assuming k is sufficiently small according to the tall-cache assumption (see below). This includes the cost of constructing the k -partitioner. The total cost for all $\lceil \frac{N}{\Delta} \rceil$ distribution steps becomes $O(N \lg k)$ comparisons and $O(k \frac{N}{\Delta} + \frac{N}{B}(1 + \log_M k)) = O(k^2 + \frac{N}{B}(1 + \log_M k))$ I/Os.

By Lemma 8, the tall-cache assumption $M \geq B^{1+\varepsilon}$ implies that for a k -partitioner and an input of size Δ , it is required that $\Delta \geq k^d$ for the I/O bounds to hold (recall $d = \max\{1 + 2/\varepsilon, 2\}$). The input assumption $\Delta \geq (\frac{2N}{\Delta})^d$ together with $k \leq \frac{2N}{\Delta}$ ensure that $\Delta \geq k^d$. \blacktriangleleft

We now prove our main result that DETERMINISTICFUNNELSELECT in Algorithm 1 is an optimal deterministic cache-oblivious multiple-selection algorithm. Crucial to the analysis is to show that the choice of Δ balances early pruning of buckets without queries with simultaneously achieving efficient I/O bounds.

► Theorem 10. DETERMINISTICFUNNELSELECT performs $O(\mathcal{B} + N)$ comparisons and $O(\mathcal{B}_{I/O} + \frac{N}{B})$ I/Os cache-obliviously in a cache model with tall assumption $M \geq B^{1+\varepsilon}$, for some constant $\varepsilon > 0$.

Proof. We first consider the consequences of the choice of Δ . By the choice of Δ , we have $\sum_{\Delta_i < \Delta} \Delta_i < N/2 + 1$. Since each bucket S_i has size at most Δ , and we only recurse on subsets that are (the union of) gaps where the two bounding rank queries of the gaps are *both* in the same bucket, we only recurse on gaps with $\Delta_i < \Delta$ elements (see Figure 3). A recursive subproblem between query ranks r_s and r_t , where $1 \leq s < t \leq q$, contains $r_t - r_s - 1 = (\sum_{i=s+1}^t \Delta_i) - 1$ elements. It follows that

- (A) all recursive subproblems in total contain at most $\sum_{\Delta_i < \Delta} \Delta_i - 1 < N/2$ elements and each subproblem has size $\leq \Delta - 2$.
- (B) $\sum_{\Delta_i \leq \Delta} \Delta_i \geq N/2 + 1$, i.e., at least $N/2$ elements are in gaps of size at most Δ .

17:10 Deterministic Cache-Oblivious Funnelselect

To analyze the number of comparisons performed, we use a potential argument where one unit of potential can pay for $O(1)$ comparisons, and all comparisons performed can be charged to the released potential. We define the potential of an element x in a gap of size Δ_i to be $1 + \lg \frac{N}{\Delta_i}$, where N is the size of the current recursive subproblem x resides in. The total initial potential is at most $N + \sum_{i=1}^{q+1} \Delta_i \lg \frac{N}{\Delta_i} = O(\mathcal{B} + N)$.

We first consider the number of comparisons for the non-sorting case (Algorithm 1, lines 9–26). If an element x in a gap of size $\Delta_i \leq \Delta$ participates in a recursive call of size $< \Delta$, the potential released for x is at least $(1 + \lg \frac{N}{\Delta_i}) - (1 + \lg \frac{\Delta}{\Delta_i}) = \lg \frac{N}{\Delta}$. If an element x in a gap of size $\Delta_i \leq \Delta$ does not participate in a recursive call, the potential released for x is $1 + \lg \frac{N}{\Delta_i} \geq 1 + \lg \frac{N}{\Delta}$. Finally, elements in gaps of size $> \Delta$ will not participate in recursive calls, and will each release at least potential 1. It follows that the released potential is at least $\frac{N}{2} + \frac{N}{2} \lg \frac{N}{\Delta}$, since at least $N/2$ elements are in gaps of size $\leq \Delta$ (property (B), contributing the second summand) and at most $N/2$ elements are in gaps of size $< \Delta$ and participate in recursive calls (property (A)), i.e., at least $N/2$ elements are in gaps of size $\geq \Delta$ (contributing the first summand). By Lemma 9, MULTIPARTITION requires $O(N \lg k)$ comparisons. Since $k = O(N/\Delta)$ (Lemma 9), these comparisons can be covered by the released potential. The additional comparisons required for computing Δ with a linear-time weighted section algorithm (Lemma 5) and performing SELECT (Lemma 2) at most twice on each bucket require in total at most $O(N)$ comparisons, and can also be charged to the released potential. It follows that for the non-sorting case the released potential can cover for all comparisons performed.

In the sorting case, a single call to FUNNELSORT is performed causing $O(N \lg N)$ comparisons (Lemma 6). No further recursive calls are made and the potential of all elements is released. At least $N + \frac{N}{2} \lg \frac{N}{\Delta}$ potential is released, since at least $N/2$ elements are in gaps of size $\leq \Delta$ (property (B)). In the sorting case, either $(2N)^d \geq \Delta^{d+1}$ or $N^{1+\frac{1}{1+\varepsilon}} \geq \Delta^2$. If $(2N)^d \geq \Delta^{d+1}$, we have $\Delta \leq (2N)^{\frac{d}{d+1}}$ and $\frac{N}{\Delta} \geq N/(2N)^{\frac{d}{d+1}} \geq \frac{1}{2} N^{\frac{1}{d+1}}$. It follows that the released potential is at least $N + \frac{N}{2} \lg(\frac{1}{2} N^{\frac{1}{d+1}}) \geq \frac{1}{2(d+1)} N \lg N$, covering the cost for the comparisons. Otherwise, $N^{1+\frac{1}{1+\varepsilon}} \geq \Delta^2$, i.e., $\Delta \leq N^{\frac{1}{2}(1+\frac{1}{1+\varepsilon})}$ and we have $\frac{N}{\Delta} \geq N/N^{\frac{1}{2}(1+\frac{1}{1+\varepsilon})} = N^{\frac{\varepsilon}{2(1+\varepsilon)}}$ and the potential released is at least $N + \frac{N}{2} \lg \frac{N}{\Delta} \geq N + \frac{\varepsilon}{4(1+\varepsilon)} N \lg N$ and can cover the cost for the comparisons. Note that the comparison bound depends on the tall-cache parameters ε and d .

To analyze the I/O cost we assign an I/O potential to an element x in gap of size Δ_i of $\frac{1}{B}(1 + \log_M \frac{N}{\Delta_i})$, where N is the size of the current subproblem x resides in. Similar to the comparison potential, it follows that the non-sorting case releases I/O potential $\frac{1}{2}(\frac{N}{B} + \frac{N}{B} \log_M \frac{N}{\Delta})$. The number of I/Os required is $O(1 + \frac{q}{B}) = O(1 + \frac{N}{B})$ I/Os for scanning R and computing Δ using weighted selection (Lemma 5), $O(k + \frac{N}{B})$ I/Os for selecting the minimum and maximum rank elements in each bucket (Lemma 2), and $O(k^2 + \frac{N}{B}(1 + \log_M k))$ I/Os for the k -partitioning (Lemma 8), i.e., in total $O(k^2 + \frac{N}{B}(1 + \log_M k))$ I/Os. It follows that the I/O cost can be charged to the released potential, provided $k^2 = O(\frac{N}{B})$. To address this, we need to consider two cases depending on the size N of a subproblem. If the problem completely fits in internal memory together with all the geometric decreasing recursive subproblems, assuming a stack-oriented memory allocation, then considering this problem will in total cost $O(1 + \frac{N}{B})$ I/Os, including all recursive subproblems. That means, there exists a constant $c > 0$ such that for $N \leq cM$, the I/O cost for handling such problems can be charged to the parent subproblem creating the subproblem. It follows that we only need to consider the I/O cost for subproblems of size $N \geq cM$. Since $M \geq B^{1+\varepsilon}$, we have $N \geq cM \geq cB^{1+\varepsilon}$, i.e., $B \leq (\frac{N}{c})^{\frac{1}{1+\varepsilon}}$. Since $k = O(\frac{N}{\Delta})$, to prove $k^2 = O(\frac{N}{B})$ it is sufficient

to prove $(\frac{N}{\Delta})^2 = O\left(\frac{N}{(N/c)^{1/(1+\varepsilon)}}\right)$. This holds, e.g., when $N^{1+\frac{1}{1+\varepsilon}} \leq \Delta^2$, which is always fulfilled in the non-sorting case. For the sorting case, we have similarly to the comparison potential that $\Omega(\frac{N}{B} \log_M N)$ I/O potential is released, which can cover the I/O cost for cache-oblivious sorting (Lemma 6). ◀

5 Conclusion

With deterministic funnelselect, we close the gap left in previous work and obtain an I/O-optimal cache-oblivious multiple-selection algorithm that does not need to resort to randomization to achieve its performance. This settles the complexity of the multiple-selection problem in the cache-oblivious model (including the fine-grained analysis based on the query-rank entropy \mathcal{B}).

There are open questions left in other variants of the problem. Like randomized funnel-select [7], deterministic funnelselect cannot deal with queries arriving in an online fashion, one after the other. This problem has been addressed in the external-memory model [2], but no cache-oblivious I/O-optimal solution is known.

Concerning the transition from single selection by rank to sorting, which multiple selection allows us to study, some questions remain unanswered. For example, in the cache-oblivious model, it is known that sorting with optimal I/O-complexity is only possible under a tall-cache assumption [6] (such as the one made in this work); for single selection, however, such a restriction is not necessary. It would be interesting to study the transition between the problems and find out, how “sorting-like” a multiple-selection instance has to be to likewise require a tall cache for I/O-optimal cache-oblivious algorithms.

Another direction for future work are parallel algorithms for multiple selection that are also cache-oblivious and I/O efficient.

References

- 1 Alok Aggarwal and Jeffrey Scott Vitter. The input/output complexity of sorting and related problems. *Commun. ACM*, 31(9):1116–1127, 1988. doi:10.1145/48529.48535.
- 2 Jérémy Barbay, Ankur Gupta, Srinivasa Rao Satti, and Jon Sorenson. Near-optimal online multiselection in internal and external memory. *Journal of Discrete Algorithms*, 36:3–17, January 2016. doi:10.1016/j.jda.2015.11.001.
- 3 Chaya Bleich and Michael L. Overton. A linear-time algorithm for the weighted median problem. Technical Report 75, New York University, Department of Computer Science, April 1983. URL: <https://archive.org/details/lineartimealgori00blei/>.
- 4 Manuel Blum, Robert W. Floyd, Vaughan R. Pratt, Ronald L. Rivest, and Robert Endre Tarjan. Time bounds for selection. *J. Comput. Syst. Sci.*, 7(4):448–461, 1973. doi:10.1016/S0022-0000(73)80033-9.
- 5 Gerth Stølting Brodal and Rolf Fagerberg. Cache oblivious distribution sweeping. In Peter Widmayer, Francisco Triguero Ruiz, Rafael Morales Bueno, Matthew Hennessy, Stephan J. Eidenbenz, and Ricardo Conejo, editors, *Automata, Languages and Programming, 29th International Colloquium, ICALP 2002, Malaga, Spain, July 8-13, 2002, Proceedings*, volume 2380 of *Lecture Notes in Computer Science*, pages 426–438. Springer, 2002. doi:10.1007/3-540-45465-9_37.
- 6 Gerth Stølting Brodal and Rolf Fagerberg. On the limits of cache-obliviousness. In Lawrence L. Larmore and Michel X. Goemans, editors, *Proceedings of the 35th Annual ACM Symposium on Theory of Computing, June 9-11, 2003, San Diego, CA, USA*, pages 307–315. ACM, 2003. doi:10.1145/780542.780589.

- 7 Gerth Stølting Brodal and Sebastian Wild. Funnelselect: Cache-oblivious multiple selection. In Inge Li Gørtz, Martin Farach-Colton, Simon J. Puglisi, and Grzegorz Herman, editors, *31st Annual European Symposium on Algorithms, ESA 2023, September 4-6, 2023, Amsterdam, The Netherlands*, volume 274 of *LIPICs*, pages 25:1–25:17. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2023. doi:10.4230/LIPICs.ESA.2023.25.
- 8 J. M. Chambers. Partial sorting [M1] (algorithm 410). *Commun. ACM*, 14(5):357–358, 1971. doi:10.1145/362588.362602.
- 9 David P. Dobkin and J. Ian Munro. Optimal time minimal space selection algorithms. *J. ACM*, 28(3):454–461, 1981. doi:10.1145/322261.322264.
- 10 Dorit Dor and Uri Zwick. Selecting the median. *SIAM Journal on Computing*, 28(5):1722–1758, 1999. doi:10.1137/s0097539795288611.
- 11 Robert W. Floyd and Ronald L. Rivest. Expected time bounds for selection. *Communications of the ACM*, 18(3):165–172, March 1975. doi:10.1145/360680.360691.
- 12 Matteo Frigo, Charles E. Leiserson, Harald Prokop, and Sridhar Ramachandran. Cache-oblivious algorithms. In *40th Annual Symposium on Foundations of Computer Science, FOCS '99, 17-18 October, 1999, New York, NY, USA*, pages 285–298. IEEE Computer Society, 1999. doi:10.1109/SFFCS.1999.814600.
- 13 Matteo Frigo, Charles E. Leiserson, Harald Prokop, and Sridhar Ramachandran. Cache-oblivious algorithms. *ACM Trans. Algorithms*, 8(1):4:1–4:22, 2012. doi:10.1145/2071379.2071383.
- 14 C. A. R. Hoare. Algorithm 65: find. *Commun. ACM*, 4(7):321–322, 1961. doi:10.1145/366622.366647.
- 15 Xiaocheng Hu, Yufei Tao, Yi Yang, and Shuigeng Zhou. Finding approximate partitions and splitters in external memory. In *Proceedings of the 26th ACM symposium on Parallelism in algorithms and architectures*. ACM, June 2014. doi:10.1145/2612669.2612691.
- 16 Kanella Kaligosi, Kurt Mehlhorn, J. Ian Munro, and Peter Sanders. Towards optimal multiple selection. In Luís Caires, Giuseppe F. Italiano, Luís Monteiro, Catuscia Palamidessi, and Moti Yung, editors, *Automata, Languages and Programming, 32nd International Colloquium, ICALP 2005, Lisbon, Portugal, July 11-15, 2005, Proceedings*, volume 3580 of *Lecture Notes in Computer Science*, pages 103–114. Springer, 2005. doi:10.1007/11523468_9.
- 17 Helmut Prodinger. Multiple Quickselect – Hoare’s Find algorithm for several elements. *Information Processing Letters*, 56(3):123–129, November 1995. doi:10.1016/0020-0190(95)00150-b.
- 18 Arnold Schönhage, Mike Paterson, and Nicholas Pippenger. Finding the median. *J. Comput. Syst. Sci.*, 13(2):184–199, 1976. doi:10.1016/S0022-0000(76)80029-3.
- 19 Michael Ian Shamos. Geometry and statistics: Problems at the interface. In Joseph Frederick Traub, editor, *Algorithms and Complexity: New Directions and Recent Results*, pages 251–280. Academic Press, 1976. URL: <http://euro.ecom.cmu.edu/people/faculty/mshamos/1976Stat.pdf>.