

## Worcester Polytechnic Institute DigitalCommons@WPI

---

Computer Science Faculty Publications

Department of Computer Science

---

9-1-2006

# An Efficient Approximate Algorithm for the k-th Selection Problem

D. Cantone

*Università di Catania*, [cantone@cs.unict.it](mailto:cantone@cs.unict.it)

G. Cincotti

*Università di Catania*, [cincotti@cs.unict.it](mailto:cincotti@cs.unict.it)

M. Hofri

*Worcester Polytechnic Institute*, [hofri@wpi.edu](mailto:hofri@wpi.edu)

Follow this and additional works at: <http://digitalcommons.wpi.edu/computerscience-pubs>

 Part of the [Computer Sciences Commons](#)

---

### Suggested Citation

Cantone, D. , Cincotti, G. , Hofri, M. (2006). An Efficient Approximate Algorithm for the k-th Selection Problem. .

Retrieved from: <http://digitalcommons.wpi.edu/computerscience-pubs/54>

This Other is brought to you for free and open access by the Department of Computer Science at DigitalCommons@WPI. It has been accepted for inclusion in Computer Science Faculty Publications by an authorized administrator of DigitalCommons@WPI.

# An Efficient Approximate Algorithm for the $k$ -th Selection Problem

D. Cantone, G. Cincotti

Dipartimento di Matematica e Informatica, Università di Catania

Viale A. Doria 6, I-95125 Catania, Italy

email: {cantone,cincotti}@cs.unict.it

and M. Hofri

Department of Computer Science, WPI

100 Institute Road, Worcester MA 01609-2280, USA

email: hofri@wpi.edu

## Abstract

We present an efficient randomized algorithm for the approximate  $k$ -th selection problem. It works *in-place* and it is fast and easy to implement. The running time is linear in the length of the input.

For a large input set the algorithm returns, with high probability, an element which is very close to the exact  $k$ -th element. The quality of the approximation is analyzed theoretically and experimentally.

**Keywords:** Selection problem, Median finding, Approximation algorithms, In-place algorithms, Randomized algorithms, Analysis of algorithms.

## 1. Introduction

Given a multi-set  $S$  of  $n$  elements from some totally ordered universe and an integer  $k \in \{1, \dots, n\}$ , the  $k$ -th selection problem consists in finding the  $k$ -th smallest element in  $S$ . The median selection is a special case of this problem when  $k = \lceil \frac{n}{2} \rceil$ .

The standard cost metric used to compare selection algorithms is the *comparison cost* model, in which the computational complexity of an algorithm is determined by the number of element comparisons. This resulted in the invention of some algorithms that use a small number of comparisons, and “pay” for it by a relatively large overhead of element moves and exchanges. We shall consider the effect of the total number of operations.

Early in the sixties, Hoare [8] devised an efficient algorithm for the  $k$ -th selection problem, namely “Quickselect”, of great theoretical and practical importance. The algorithm takes the same approach as in “Quicksort,” making a sequence of random choices of pivot elements, executing about  $4n$  comparisons on the average. It has the problem of having a cost in  $\mathcal{O}(n^2)$  time in the worst-case.

Floyd and Rivest [7] improved the previous algorithm by choosing the pivot elements through an opportune sampling. They obtained a very practical algorithm requiring  $n + \min(k, n - k) + o(n)$ , i.e. at most  $1.5n + o(n)$ , comparisons on the average (cf. [5]).

The first discovered algorithm with linear worst-case cost was given by Blum et al. [2] in 1973, requiring only  $5.43n + o(n)$  comparisons in the worst-case. It was an important complexity result because till then the selection problem was assumed to be as difficult as sorting. Incidentally, a lower bound of  $2n$  comparisons (in the worst-case) for the median selection problem is due to Bent and John [3].

Three years later, Schönhage et al. [11] concentrated their attention on the median selection problem, devising an algorithm requiring at most  $3n + o(n)$  comparisons. Such result was slightly improved in 1995 by Dor and Zwick [6] that obtained a  $2.95n + o(n)$  upper bound modifying the algorithm by Schönhage et al. . These algorithms attempts to minimize the number of key comparisons but not the total number of operations, and hence do not lead to faster algorithms in practice.

All the worst-case linear time algorithms cited above require, apart from the space for the elements themselves, at least  $\Theta(\log n)$  additional space during the computation. From this point of view, the research community showed interest in finding optimal *in-place* solutions, i.e. algorithms that require only a constant amount of extra space. The first in-place algorithm, due to Lai and Wood [9] in 1988, needs at most  $6.83n + o(n)$  comparisons and it was obtained by an implicit emulation of the basic algorithm by Blum et al. .

Based on similar ideas as the median selection algorithm by Schönhage et al., Carlsson and Sundstrom [4] in 1995 achieved the  $(2.95 + \varepsilon)n + o(n)$  in-place upper bound, for any  $\varepsilon > 0$ . This is arbitrarily close to the Dor and Zwick's upper bound for the same problem, without space restrictions. It is to be noticed that these algorithms have mainly theoretical interest because their implementations are extremely complicated and not practical.

Selection arises in several applications others than order statistics. In most cases, the selected elements are used as threshold for filtering operations, because it is reasonable to concentrate attention only to the subset of the most important elements with respect to a suitable criterion.

For instance, a computer chess program could quickly evaluate all possible next moves, and then decide to further elaborate only the best  $p$ -percentage moves more carefully. Identifying such best moves can be done by a selection and a filtering operation.

Another explicit use of selection arises in filtering outlying elements. More specifically, in dealing with noisy statistical data samples, sometimes it is convenient to eliminate noise throwing out the largest and smallest  $p$ -percentage of them. Selection can be used to identify the items defining the  $p$ -th and  $(100 - p)$ -th percentiles and the outliers are then filtered out by comparing each item to the two selected elements.

A crucial observation for the current treatment is that many applications do not require the *exact* solution to an instance of the  $k$ -th selection problem, and an *approximate* solution would suffice. In such contexts it becomes important to deal with efficient approximate algorithms for the  $k$ -th selection problem.

An extremely efficient algorithm for the approximate median selection problem is described in [1]. The idea behind this algorithm is simple. To simplify the exposition, let the size of the input set be a power of 3, e.g.  $n = 3^r$  with  $r \in \mathbb{N}$ . The algorithm then proceeds in  $r$  stages. At each stage it divides the surviving input into subsets of three elements, and calculates the median of each such triplet. The "local medians" survive to the next stage. The algorithm continues recursively, using the local results to compute the approximate median of the initial set. In [1] is also described how to handle efficiently input sizes that are not powers of 3.

In this paper, we present an efficient algorithm for the approximate  $k$ -th selection problem. It can be viewed as a clever generalization of the algorithm described in [1]. In particular, it works in-place and its precision is strongly based on the good behavior of the approximate median selection algorithm.

The paper is organized as follows. In Section 2, we first show that the  $k$ -th and the median selection problems are equivalent. An algorithm for the approximate  $k$ -th selection problem is described in Section 3, and one of its possible implementation is discussed in Section 4. Moreover, in Section 5 we prove the linear computational cost of the algorithm, whereas in Section 6 we develop a probabilistic analysis aimed to establish the quality of the algorithm output. Section 7 collects some experimental results. The last section concludes the paper with some final remarks.

## 2. Reducing the $k$ -th selection to the median selection problem

In this section, we show that though the  $k$ -th selection is more general than the median selection problem they are equivalent. In particular, we describe how an instance of the  $k$ -th selection problem can be easily reduced to an instance of the median selection problem.

Let  $S^{(n)}$  be a multi-set<sup>1</sup> of  $n$  elements from a totally ordered universe. In the following, we denote by  $-\infty$  (resp.  $+\infty$ ) an element of the universe smaller than the minimum (resp. greater than the maximum) of  $S^{(n)}$ .

Solving an instance  $\langle S^{(n)} \rangle$  of the median selection problem means to find the median, i.e. the  $\lceil \frac{n}{2} \rceil$ -th, element out of the set  $S^{(n)}$ . In the more general  $k$ -th selection problem, solving an instance  $\langle S^{(n)}, k \rangle$ , with  $k \in \{1, \dots, n\}$ , requires finding the  $k$ -th element out of  $S^{(n)}$ .

**Reduction 1.** *Given an instance  $\langle S^{(n)}, k \rangle$  of the  $k$ -th selection problem, it can be reduced to an instance  $\langle S_{\infty}^{(n,m)} \rangle$  of the median selection problem, where the multi-set*

$$S_{\infty}^{(n,m)} = S^{(n)} \cup \bigcup_{i=1}^m \{\infty\}, \quad \text{with} \quad m = |n - 2k + 1|,$$

is obtained from  $S^{(n)}$  by adding  $m$  copies of the  $\infty$  value defined as follows:

$$\infty = \begin{cases} -\infty & \text{if } k \leq \lceil \frac{n}{2} \rceil, \\ +\infty & \text{otherwise.} \end{cases}$$

Obviously, we have  $0 \leq m < n$  with  $|S_{\infty}^{(n,m)}| = n + m$ , and the following:

**Theorem 1.** *The Reduction 1 from the  $k$ -th selection to the median selection problem is correct.*

*Proof.* When  $k \leq \lceil \frac{n}{2} \rceil$ , we construct the multi-set  $S_{\infty}^{(n,m)}$  obtained from  $S^{(n)}$  by adding  $m = n - 2k + 1 \geq 0$  copies of the element  $-\infty$ . It is immediate to recognize that the  $k$ -th element of  $S^{(n)}$  is exactly the  $m + k = (n - k + 1)$ -st element of  $S_{\infty}^{(n,m)}$ , namely the median element of  $S_{\infty}^{(n,m)}$ , because  $|S_{\infty}^{(n,m)}| = 2n - 2k + 1$  and the median is given by the  $\lceil \frac{|S_{\infty}^{(n,m)}|}{2} \rceil = \lceil n - k + \frac{1}{2} \rceil = (n - k + 1)$ -st element.

---

<sup>1</sup>A set that can contain multiple copies of the same element.

On the other hand, if  $k > \lceil \frac{n}{2} \rceil$ , then we construct the multi-set  $S_\infty^{(n,m)}$  by adding  $m = 2k - n - 1 > 0$  copies of the element  $+\infty$ . In this case, the  $k$ -th element of  $S^{(n)}$  is exactly the  $k$ -th element of  $S_\infty^{(n,m)}$  yet, namely the median element of  $S_\infty^{(n,m)}$ , because  $|S_\infty^{(n,m)}| = 2k - 1$  and the median is given by the  $\lceil \frac{|S_\infty^{(n,m)}|}{2} \rceil = \lceil k - \frac{1}{2} \rceil = k$ -th element. Thus, in both cases the problem of selecting the  $k$ -th element of  $S^{(n)}$  has been reduced to the one of selecting the median of  $S_\infty^{(n,m)}$ . ■

By virtue of the above reduction, any algorithm for the median selection can be used for solving the  $k$ -th selection problem. In particular, the reduction is useful for both the *exact* and the *approximate* variants of the  $k$ -th selection problem, and the latter case may be solved using the efficient and accurate algorithm for the approximate median selection presented in [1].

### 3. An probabilistic, approximate algorithm

The Reduction 1 described in Section 2 and the idea behind the approximate median selection algorithm are joined together, to get a randomized algorithm for the approximate  $k$ -th selection problem, that works *in-place*, without introducing any explicit infinite value.

In the following, we denote by  $\mathbb{N}$  the set of cardinals, with  $\mathbb{N}_0 = \mathbb{N} \cup \{0\}$ , and by  $\mathbb{Z}$  the set of relative integers. Moreover, we define  $[n]_k \stackrel{\text{def}}{=} n \bmod k$  for any  $n, k \in \mathbb{N}$ , thus  $[n]_k \in \{0, 1, \dots, k-1\}$ .

Let  $\langle S_\infty^{(n,m)} \rangle$  be the reduction of an instance  $\langle S^{(n)}, k \rangle$  of the  $k$ -th selection problem with  $n, m \in \mathbb{N}_0$  and  $m < n$ . For each  $i = 0..3$ , we calculate the probabilities:

$$p_i^{(n,m)} = \Pr \left\{ \begin{array}{l} \text{A random unordered triplet } t = \{x, y, z\} \in S_\infty^{(n,m)} \\ \text{contains } i \text{ infinite values} \end{array} \right\},$$

and find

$$p_0^{(n,m)} = \frac{\binom{n}{3}}{\binom{n+m}{3}}, \quad p_1^{(n,m)} = \frac{m \binom{n}{2}}{\binom{n+m}{3}}, \quad p_2^{(n,m)} = \frac{n \binom{m}{2}}{\binom{n+m}{3}}, \quad p_3^{(n,m)} = \frac{\binom{m}{3}}{\binom{n+m}{3}},$$

where  $\binom{n+m}{3}$  is the number of possible unordered triplets that can be selected from  $S_\infty^{(n,m)}$ . At this point, let the set  $S_\infty^{(n,m)}$  be partitioned into  $\frac{1}{3}(n+m)$  uniformly distributed triplets;<sup>2</sup> the expected number  $\mu_i(n, m)$  of unordered triplets in  $S_\infty^{(n,m)}$  containing  $i$  infinite values is given by:

$$\mu_i(n, m) = \frac{1}{3}(n+m)p_i^{(n,m)}, \quad \text{for each } i = 0..3.$$

---

<sup>2</sup>At this level of abstraction, it does not matter if  $(n+m)$  is not a multiple of 3.

In particular, we obtain:

$$\mu_0(n, m) = \frac{n(n-1)(n-2)}{3(n+m-1)(n+m-2)}, \quad (1)$$

$$\mu_1(n, m) = \frac{mn(n-1)}{(n+m-1)(n+m-2)}, \quad (2)$$

$$\mu_2(n, m) = \frac{m(m-1)n}{(n+m-1)(n+m-2)}, \quad (3)$$

$$\mu_3(n, m) = \frac{m(m-1)(m-2)}{3(n+m-1)(n+m-2)}. \quad (4)$$

Because  $\sum_{i=0}^3 p_i^{(n,m)} = 1$ , we get the equality  $\sum_{i=0}^3 \mu_i(n, m) = \frac{1}{3}(n+m)$  and, in particular, the fundamental system of identities:

$$\begin{cases} 3\mu_0(n, m) + 2\mu_1(n, m) + \mu_2(n, m) = n, \\ \mu_1(n, m) + 2\mu_2(n, m) + 3\mu_3(n, m) = m. \end{cases} \quad (5)$$

Equations (1) through 4 suggest a randomized algorithm for the approximate  $k$ -th selection problem where infinite values are only used for book-keeping, and for this reason we refer to them as *virtual infinities*. A detailed description of the algorithm follows.

1. Let  $\langle S, k \rangle$  be the problem instance in input, with  $n = |S|$  and  $m = |n - 2k + 1|$ . Furthermore, let the boolean variable *PositiveInfinities* be *True* whether  $k > \lceil \frac{n}{2} \rceil$  or *False* otherwise.
2. Repeat the following steps until the cardinality  $n$  becomes smaller than a fixed *Threshold*:
  - (a) Find integers  $r, r_\infty \in \mathbb{N}_0$ , by means of a predefined criterion, such that:

$$r + r_\infty = \lceil \frac{n+m}{3} \rceil.$$

These values represent the number of elements in  $S$  and the number of virtual infinities, respectively, that do not exactly fit in any triplet.

- (b) Compute values  $\bar{\mu}_i \in \mathbb{N}_0$ , with  $i = 0..3$ , such that  $\bar{\mu}_i$  is a “good” integer approximation for  $\mu_i(n, m)$  and the following integer system:

$$\begin{cases} 3\bar{\mu}_0 + 2\bar{\mu}_1 + \bar{\mu}_2 = n - r, \\ \bar{\mu}_1 + 2\bar{\mu}_2 + 3\bar{\mu}_3 = m - r_\infty, \end{cases} \quad (6)$$

is satisfied with such values.

- (c) Accordingly to the first equation of system (6), choose randomly in  $S$ :  $\bar{\mu}_0$  triplets of elements,  $\bar{\mu}_1$  couples of elements, and  $\bar{\mu}_2 + r$  single elements.
- (d) Perform the following operations on the selected elements of  $S$ :
  - remove from  $S$  the minimal and maximal elements of each of the  $\bar{\mu}_0$  triplets;
  - if *PositiveInfinities* is *True* then remove from  $S$  the minimal element of each of the  $\bar{\mu}_1$  couples, otherwise remove the maximal element of each of them;

- remove from  $S$  each of the  $\bar{\mu}_2$  singleton.
- (e) Retain or remove the  $r$  elements of  $S$  with a predefined policy, and set  $n$  to the cardinality of the remaining set. Treat analogously the  $r_\infty$  virtual infinities, saving or removing them, i.e., respectively, letting  $m = \bar{\mu}_2 + \bar{\mu}_3 + r_\infty$  or  $m = \bar{\mu}_2 + \bar{\mu}_3$  (see below).
3. Return the exact median of the multi-set obtained from the remaining elements in  $S$  joined with  $m$  copies of the value  $+\infty$  whether *PositiveInfinities* is *True*, or  $m$  copies of  $-\infty$  otherwise.

Observe that, any reasonable criterion introduced in Step (2a) should ensure  $r_\infty \leq r$  because we always have  $m < n$ . In particular, a satisfying criterion is obtained setting  $r = [n + m]_3 - r_\infty$  and

$$r_\infty = \begin{cases} 1 & \text{if } [n + m]_3 = 2 \text{ and } m \geq \frac{3}{4}n, \\ 0 & \text{otherwise,} \end{cases}$$

where the value  $\frac{3}{4}n$  is arbitrarily chosen and it means a large fraction of  $n$ .

Obviously, the system (6) has always integer solutions. Indeed,  $\bar{\mu}_2, \bar{\mu}_3$  can be considered as free variables, and  $\bar{\mu}_1$  can be obtained from the second equation; now, adding the two equations in (6), we obtain  $3 \sum_{i=0}^3 \bar{\mu}_i = n + m - [n + m]_3$ , and get  $\bar{\mu}_0$  because  $n + m - [n + m]_3$  is multiple of 3. Notice that the value  $\sum_{i=0}^3 \bar{\mu}_i$  is the number of *virtual triplets*, i.e. triplets in which virtual infinities can occur, processed in a single iteration.

Each operation listed in Step (2d) corresponds to the handling of virtual triplets containing, respectively, 0, 1 or 2 infinite values.

A good policy to treat the  $r$  elements in Step (2e) is addressed in [1]; these elements are processed joining them with one of the  $\bar{\mu}_0$  triplets and their local median (i.e., of 3, 4, or 5 elements) survives to the next iteration. Such technique could also include the  $r_\infty$  virtual infinities.

Finally, the value  $m$  computed in Step (2e) is the number of virtual infinities that survive to the next iteration; specifically, such infinities come out from virtual triplets containing 2 or 3 infinite values. Analogously, in Step (2e) we have  $n = \bar{\mu}_0 + \bar{\mu}_1 + r$  or  $n = \bar{\mu}_0 + \bar{\mu}_1$  elements coming out from virtual triplets containing 0 or 1 infinite values.

**Interval extremes adjustment.** The extreme cases of selection, such as the smallest (i.e.,  $k \approx 1$ ) or the largest (i.e.,  $k \approx n$ ) elements are better treated *ad hoc*, modifying the quantities given in equations (1) through (4).

## 4. An implementation of the algorithm

In this section we give a possible implementation of the algorithm described above. It is assumed that the selection is from values stored in an array in main memory. In this case, the computations disturb the order of the elements, but the contents of the array is unchanged.

Let  $n$  be the size of the input array  $A$ , and  $k$  be the index of the element that we want to select, and let *Threshold* be an integer value. The algorithm pseudo-code is shown in Fig. 1; it provided the experimental results reported in Section 7.

For the sake of brevity of this presentation, the value  $r_\infty$  introduced in Step (2a) is always set to zero, and the  $r$  elements treated in Step (2e) are simply thrown away. Asymptotically, these choices do not decrease the output quality in terms of precision, but it makes a noticeable difference in small arrays.

### Pseudo-code of the approximate selection function

```

FUNCTION ApproximateSelection ( $A, n, k, Threshold$ ) : integer
    This function returns the index of the approximate  $k$ -th element of the input array  $A[1, \dots, n]$ .

     $m = \lfloor n - 2k + 1 \rfloor$ ;
    if ( $k > \lceil n/2 \rceil$ ) then  $PositiveInfinities = True$ ;
        else  $PositiveInfinities = False$ ;
    while ( $n > Threshold$ ) do
         $r = \lfloor n + m \rfloor_3$ ;
         $\langle \bar{\mu}_0, \bar{\mu}_1, \bar{\mu}_2, \bar{\mu}_3 \rangle = \text{SystemSolve}(n, m)$ ;
        while ( $r > 0$ ) do
            Swap ( $A[n], A[\text{RandomInteger}(1 \dots n)]$ );
             $n = n - 1$ ;  $r = r - 1$ ;
            /* Eliminate  $r$  random */
            /* elements from consideration */
        end while;
         $i = 1$ ;  $j = 1$ ;  $m = \bar{\mu}_2 + \bar{\mu}_3$ ;
        while ( $j < n$ ) do
             $Choice = \text{RandomInteger}(1 \dots (\bar{\mu}_0 + \bar{\mu}_1 + \bar{\mu}_2))$ ;
            if ( $Choice \leq \bar{\mu}_2$ ) then
                /*Delete the singleton*/
                 $j = j + 1$ ;  $\bar{\mu}_2 = \bar{\mu}_2 - 1$ ;
            elseif ( $Choice \leq \bar{\mu}_1 + \bar{\mu}_2$ ) then
                /* Treat the couple */
                if ( $A[j] > A[j + 1]$ ) then Swap ( $A[j], A[j + 1]$ );
                if ( $PositiveInfinities = True$ ) then Swap ( $A[i], A[j + 1]$ );
                    else Swap ( $A[i], A[j]$ );
                 $j = j + 2$ ;  $i = i + 1$ ;  $\bar{\mu}_1 = \bar{\mu}_1 - 1$ ;
            else
                /* Treat the triplet */
                 $h = \text{TripletMedian}(A, j)$ ;
                Swap ( $A[i], A[h]$ );
                 $j = j + 3$ ;  $i = i + 1$ ;  $\bar{\mu}_0 = \bar{\mu}_0 - 1$ ;
            end while;
             $n = i - 1$ ;
        end while;
        SelectionSort ( $A, n$ );
        if ( $PositiveInfinities = True$ ) then
            /* We consider the  $m$  residues virtual infinities */
            if ( $m \geq n$ ) then return  $n$ ;
            else return  $\lceil (m + n) / 2 \rceil$ ;
        else
            if ( $m \geq n$ ) then return 1;
            else return  $\lceil (n - m) / 2 \rceil$ ;
    end while;

```

Figure 1: Pseudo-code for the approximate  $k$ -th selection algorithm.

The values  $\bar{\mu}_i$ , with  $i = 0..3$ , required in Step (2b) to solve the integer system (6), are computed by the function *SystemSolve*; its pseudo-code is given in Fig. 2. In this function, the real variables  $x, y$  are initialized by means of functions (3) and (4), using the value  $(n - r)$  because the  $r$  elements are removed at all. Furthermore, the best first choice for the couple in  $\Gamma$  is represented by  $\langle \text{Round}(x), \text{Round}(y) \rangle$ .



For convenience, the value  $m$  computed in Step (2e) is calculated before Step (2c).

Steps (2c) and (2d) are executed by scanning the array from left to right through the index  $j$ . For every single step of each scan, we flip a coin with three faces, i.e. the integer variable *Choice*, and depending on the result, we process one, two, or three consecutive elements of the array starting at position  $j$ . Notice that, so many calls to the random number generator can be avoided introducing further efficient pseudo-randomization methods as in the style of [1]. All the elements that survive to the next step are gathered on the leftmost portion of the input array through the index  $i$ . We show in Fig. 2 the function *TripletMedian* that finds the index of the median of the triplet formed by consecutive elements starting at position  $i$  of the array  $A$ .

The algorithm continues iteratively using the results of each stage as input for a new one. This is done until the number of the surviving elements falls below the *Threshold* value. Then, the procedure *SelectionSort*, shown in Fig. 2, sorts the remaining elements and the index of exact median of the set described in Step (3) is returned. In particular, if an infinite should be returned, i.e.  $m \geq n$ , it is replaced with the remaining element closest to it.

## 5. Run-time analysis

To evaluate the computational cost of our algorithm, we not only consider the usual metric of key comparisons, but also look briefly at the number of element swaps in this particular implementation, given in Section 4.

Let  $\delta_t = \frac{8}{3}$  be the average number of key comparisons needed to determine the median element of a triplet (cf. [1]) and let  $\delta_c = 1$  be the average number of key comparisons needed to determine the smallest or the largest element of a couple.

The following recurrence, where  $n$  and  $m$  stand for the number of real keys and virtual infinities respectively, describes the evolution of the expected number of key comparisons from stage to stage:

$$C(n, m) = \begin{cases} 0 & \text{if } n = 1, \\ c(n, m) + C(\mu_0(n, m) + \mu_1(n, m), \mu_2(n, m) + \mu_3(n, m)) & \text{if } n > 1. \end{cases} \quad (7)$$

Here  $c(n, m) = \delta_t \mu_0(n, m) + \delta_c \mu_1(n, m)$  is the average number of comparisons needed to process the  $\mu_0(n, m)$  triplets and the  $\mu_1(n, m)$  couples in a single stage of the algorithm; the recursive term adds the average number of comparisons needed to handle the real and infinite keys surviving to the next stage. The average number of key comparisons performed by the algorithm when the input size is  $n$  and we are searching for the  $k$ -th element is given by  $C(n, |n - 2k + 1|)$ .

It is easy to see that our algorithm has running-time which is linear in the input size  $n$ . Indeed, let us consider the following two bounding cases:  $k = \lceil \frac{n}{2} \rceil$  and  $k = 1$  (or symmetrically  $k = n$ ).

In the first case, when  $n$  is odd, we have  $m = 0$  and equation (7) can be rewritten as follows:

$$C(n, 0) = \frac{8}{9}n + C\left(\frac{n}{3}, 0\right).$$

This recurrence depends of a single variable and its exact solution, when  $n$  is a power of 3, is  $C(n, 0) = \frac{4}{3}(n - 1)$ .

In the second case, we have  $m = n - 1$  but, for the purpose of the analysis, we make the following approximations  $m \approx n$ ,  $\mu_0(n, m) \approx \mu_3(n, m) \approx \frac{n}{12}$ ,  $\mu_1(n, m) \approx \mu_2(n, m) \approx \frac{n}{4}$  and equation (7) can then be

### Pseudo-code for SystemSolve, TripletMedian, and SelectionSort

**FUNCTION SystemSolve** ( $n, m$ ) : integers quadruple

*This function returns four integer values satisfying conditions given in Step (2b).*

```

 $r = \lfloor n + m \rfloor_3;$ 
 $x = \mu_2(n - r, m);$ 
 $y = \mu_3(n - r, m);$ 
 $\Gamma = \{ \langle \lfloor x \rfloor, \lfloor y \rfloor \rangle, \langle \lfloor x \rfloor, \lceil y \rceil \rangle, \langle \lceil x \rceil, \lfloor y \rfloor \rangle, \langle \lceil x \rceil, \lceil y \rceil \rangle \};$ 
- Choose an ordered couple  $\langle \bar{\mu}_2, \bar{\mu}_3 \rangle \in \Gamma$  such that the following integer values:
-  $\bar{\mu}_1 = m - 2\bar{\mu}_2 - 3\bar{\mu}_3,$ 
-  $\bar{\mu}_0 = (n - r - 2\bar{\mu}_1 - \bar{\mu}_2)/3,$ 
- are non-negatives, i.e.  $\bar{\mu}_i \in \mathbb{N}_0$  with  $i = 0..3;$ 
return  $\langle \bar{\mu}_0, \bar{\mu}_1, \bar{\mu}_2, \bar{\mu}_3 \rangle;$ 

```

**FUNCTION TripletMedian** ( $A, i$ ) : integer

*This function returns the index of the median of the triplet  $(A[i], A[i + 1], A[i + 2])$ .*

```

if ( $A[i] < A[i + 1]$ ) then
  if ( $A[i + 2] < A[i]$ ) then return  $i;$ 
  elsif ( $A[i + 2] < A[i + 1]$ ) then return  $(i + 2);$ 
  else return  $(i + 1);$ 
else
  if ( $A[i] < A[i + 2]$ ) then return  $i;$ 
  elsif ( $A[i + 2] > A[i + 1]$ ) then return  $(i + 2);$ 
  else return  $(i + 1);$ 

```

**PROCEDURE SelectionSort** ( $A, Size$ )

*This procedure rearranges the elements of the array  $A[1, \dots, Size]$  in ascending sorted order.*

```

 $i = 1;$ 
while ( $i < Size$ ) do
   $min = i;$ 
   $j = i + 1;$ 
  while ( $j \leq Size$ ) do
    if ( $A[j] < A[min]$ ) then  $min = j;$ 
     $j = j + 1;$ 
  end while;
  Swap ( $A[i], A[min]$ );
   $i = i + 1;$ 
end while;

```

Figure 2: Pseudo-code for the procedures used by the approximate  $k$ -th selection algorithm.

rewritten as follows:

$$C(n, n) = \frac{17}{36}n + C\left(\frac{n}{3}, \frac{n}{3}\right),$$

with the exact solution, when  $n$  is a power of 3, given by  $C(n, n) = \frac{17}{24}(n - 1)$ .

Asymptotically, the leading coefficients of the above solutions does not change when  $n$  is not a power of 3. On the other hand, it is easy to verify that, among all the possible  $k$  values, the first case above is the most expensive in terms of key comparisons—since comparisons are used optimally when triplets have infinities. Hence can state the following:

**Theorem 2.** *The average number of key comparisons performed by the algorithm when the input size is  $n$ , and the  $k$ -th element is searched for, is no more than  $\frac{4}{3}n$ .*

The above reason underpins the way  $k$  influences the number of comparisons required by the algorithm. Unfortunately, there is little hope of securing a general, analytic solution of the recurrence 7, but one can numerically verify that the function  $C(n, m)$  is monotonically decreasing in  $m$ , with  $0 \leq m < n$  and  $n$  fixed; in particular, the minimum bound is reached in the above second case.

The worst-case analysis can be analogously handled considering the maximum number of key comparisons  $\delta_t = 3$  required to find the median of three elements. We get:

**Theorem 3.** *The number of key comparisons performed by the algorithm when the input size is  $n$ , and the  $k$ -th element is searched for, is no more than  $\frac{3}{2}n$ .*

Let us consider the number of element moves. Unlike the procedure developed in [1], it is not easy here to avoid doing as we do in Section 4, and move each surviving element (or swap it with another element in the array, if we need to preserve its content). Therefore the number of elements swaps, both in the average- and in the worst-case, is given by equation (7) setting  $\delta_t = \delta_c = 1$ . Again, not much hope for an exact solution, but it is a immediate to verify that  $C(n, m) < n$ . Moreover, one can numerically verify that  $C(n, m) < \frac{3}{4}n$  for all  $0 \leq m < n$ , so we can conclude that the number of element swaps performed by the algorithm given in Fig. 1, when the input size is  $n$  and we are searching for the  $k$ -th element, is no more than  $\frac{3}{4}n$ .

## 6. Probabilistic analysis

In this section we give an exact probabilistic analysis of the  $k$ -th selection algorithm implemented in Section 4, which returns the value *Output*. Specifically, for any  $n, k \in \mathbb{N}$ , with  $1 \leq k \leq n$ , we calculate the following probability distribution:

$$\mathcal{P}_{n,k}(z) = \Pr \{ |\text{rank}(\text{Output}) - k| \leq \lfloor zn \rfloor \}, \quad (8)$$

of the relative selection error, measured by  $z \in [0, 1]$ .

W.l.o.g., we can assume that the input array of size  $n$  is a permutation of the integers  $1, \dots, n$ . The set of the  $n!$  permutations is denoted by  $S_n$ , sometimes called the free group (it is a group, under composition, but we shall not use this fact here).

We make two reasonable uniformity assumptions. One concerns the distribution of the input set: all permutations are assumed equally likely. The second refers to the random choice of triplets, couples and singletons during the execution of the algorithm, which we assume are all independent of the array content and of each other.

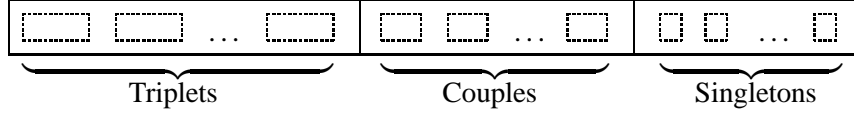


Figure 3: Looking at the array elements as a sequence of sub-sequences.

While the algorithm creates the successive triplets probabilistically, for the following exposition it is convenient to view the input array as a sequence of consecutive triplets followed by consecutive couples and by consecutive singletons as shown in Fig. 3, which does not show of course the virtual infinities.

W.l.o.g. in the analysis, we consider the case  $1 \leq k \leq \lceil \frac{n}{2} \rceil$ , with the symmetric one,  $\lceil \frac{n}{2} \rceil < k \leq n$ , being completely analogous. This means that the infinities are negative, and the surviving element of any couple is its smaller element.

We use of the following notation and identity for the multinomial coefficient:

$$\binom{N}{\{k\}^m} \stackrel{\text{def}}{=} \binom{N}{\underbrace{k, k, \dots, k}_{m \text{ times}}} = \prod_{i=0}^{m-1} \binom{N-ki}{k} = \frac{N!}{(k!)^m (N-km)!}, \quad (9)$$

for all  $N, k, m \in \mathbb{N}_0$ , with  $N \geq km$ .

We construct our solution for the probability distribution (8), from the answer to the question “what is the probability that when the array has the  $n$  values  $1, \dots, n$  in random order, the value  $a$  is selected, and becomes the  $b^{\text{th}}$  smallest among the surviving elements.” We decompose this calculation to a sequence of combinatorial problems.

A sequence of  $l$  two-position boxes is given. We place in them a random permutation of distinct integers,  $\pi \in S_{2l}$ . The smallest element in each box is selected, giving us  $l$  surviving elements.

**Definition 1.** Given  $l \in \mathbb{N}$  and  $a, b \in \mathbb{N}_0$ , we denote by  $Q_l^{[2]}(a, b)$  the number of permutations  $\pi$ , out of the  $(2l)!$  possible ones, in which exactly  $b$  elements survive among the  $a$  smallest elements in  $\pi$ .

For convenience, we extend Def. 1 for  $l = 0$  too, setting:

$$Q_0^{[2]}(a, b) = \begin{cases} 1 & \text{if } a = b = 0, \\ 0 & \text{otherwise.} \end{cases}$$

The following notation will simplify our discussion. For any permutation  $\pi$  and an integer  $a$ , we denote by “•” each element among the  $a$  smallest ones in  $\pi$ , and denote by “o” the others.

**Lemma 4.** Let  $a, b, l \in \mathbb{N}_0$ . The number of permutations  $Q_l^{[2]}(a, b)$  is given by

$$Q_l^{[2]}(a, b) = \begin{cases} 2^{(2b-a)} l! b! (l-b)! \binom{2l-a}{l-b} \binom{l-a+b}{2b-a} \binom{a}{a-b} & \text{if } b \leq l, \\ 0 & \text{otherwise.} \end{cases}$$

*Proof.* It is immediate to verify that any contribution to  $Q_l^{[2]}(a, b)$  can be given only if the condition  $b \leq a \leq 2b \leq 2l$  is satisfied.

In a permutation  $\pi$ , among the  $2l$  elements, there are  $a$  “•” and  $2l - a$  “o”. We count permutations that result in selections of  $b$  “•” and  $l - b$  “o”. All the surviving “o” must come from the  $l - b$  couples

of type  $\{\circ, \circ\}$  in  $\pi$ , and thus  $(2l - a) - 2(l - b) = 2b - a$  unordered couples of type  $\{\circ, \bullet\}$  occur in  $\pi$ . Consequently,  $[a - (2b - a)]/2 = a - b$  couples of type  $\{\bullet, \bullet\}$  are present in  $\pi$ .

We need to count the permutations  $\pi$  that satisfy the above constraints. First we count the number of ways to select occupants for the unmixed boxes: We have  $\frac{1}{(a-b)!} \binom{a}{2}^{a-b}$  ways to select  $a - b$  pairs of “ $\bullet$ ” from the  $a$  “ $\bullet$ ” in  $\pi$ , and  $\frac{1}{(l-b)!} \binom{2l-a}{2}^{l-b}$  ways to select  $(l - b)$  pairs of “ $\circ$ ” from the  $(2l - a)$  available “ $\circ$ ” in  $\pi$ . Then, we have  $\frac{1}{(2b-a)!} \prod_{i=0}^{2b-a-1} (2b - a - i)^2 = (2b - a)!$  ways to match the remaining  $(2b - a)$  “ $\circ$ ” and  $(2b - a)$  “ $\bullet$ ” with the  $(2b - a)$  boxes containing mixed couples.

Finally, to complete the specification of  $\pi$ , a factor  $2^l$  accounts for the order inside each box, and an  $l!$  accounts for the order of the boxes. Multiplying these factors and using the identity (9), we get:

$$Q_l^{[2]}(a, b) = \begin{cases} 2^{(2b-a)} l! \frac{a!(2l-a)!}{(a-b)!(l-b)!(2b-a)!} & \text{if } b \leq a \leq 2b \leq 2l, \\ 0 & \text{otherwise,} \end{cases} \quad (10)$$

and rewriting it using the binomial coefficients we get the desired result.  $\blacksquare$

Notice that, for any  $l \in \mathbb{N}_0$ , we have  $Q_l^{[2]}(0, 0) = (2l)!$  and  $Q_l^{[2]}(a, 0) = 0$  for all  $a \neq 0$ .

Now we consider the following consequence of the above result:

**Definition 2.** Given  $a, b, l \in \mathbb{N}$ , we denote by  $\overline{Q}_l^{[2]}(a, b)$  the number of permutations  $\pi$ , out of the  $(2l)!$  possible ones, in which the  $a^{\text{th}}$  smallest element in  $\pi$  is selected and becomes the  $b^{\text{th}}$  smallest among the surviving elements.

**Lemma 5.** Let  $a, b, l \in \mathbb{N}$ . The number of permutations  $\overline{Q}_l^{[2]}(a, b)$  is given by

$$\overline{Q}_l^{[2]}(a, b) = 2l(2l - a)Q_{l-1}^{[2]}(a - 1, b - 1).$$

*Proof.* Recycling the notation of the previous proof, we observe that the  $a^{\text{th}}$  element in  $\pi$  is the largest “ $\bullet$ ” and given that it is selected, it came from a mixed couple. Any permutation that contributes here was counted in  $Q_l^{[2]}(a, b)$ , if  $a$  was among those selected.

Since the candidate permutations have  $2b - a$  mixed couples, then among the permutations we counted above only in  $(2b - a)/a$  of them we find the largest “ $\bullet$ ” in a mixed couple, that is,

$$\overline{Q}_l^{[2]}(a, b) = \frac{2b - a}{a} Q_l^{[2]}(a, b) = 2l(2l - a)Q_{l-1}^{[2]}(a - 1, b - 1),$$

where the second equality follows from the explicit form above in Eq. (10).  $\blacksquare$

By convention, we extend Def. 2 to the case  $l = 0$  and we let  $\overline{Q}_0^{[2]}(a, b) = 0$  for any  $a, b \in \mathbb{N}$ .

We deal now with the the same combinatorial problem with respect to triplets of actual values, instead of couples. A random permutation  $\pi$  from  $S_{3l}$  is placed in a sequence of  $l$  three-position boxes, as before. We select the median element in each box, and get  $l$  surviving elements.

**Definition 3.** Given  $l \in \mathbb{N}$  and  $a, b \in \mathbb{N}_0$ , we denote by  $Q_l^{[3]}(a, b)$  the number of permutations  $\pi$ , out of the  $(3l)!$  possible ones, in which exactly  $b$  elements survive among the  $a$  smallest elements in  $\pi$ .

As above, it will be convenient to extend Def. 3 to  $l = 0$ , setting:

$$Q_0^{[3]}(a, b) = \begin{cases} 1 & \text{if } a = b = 0, \\ 0 & \text{otherwise.} \end{cases}$$

**Lemma 6.** Let  $a, b, l \in \mathbb{N}_0$ . The number of permutations  $Q_l^{[3]}(a, b)$  is given by

$$Q_l^{[3]}(a, b) = \begin{cases} 3^{a-b} a! (3l - a)! \binom{l}{b} \sum_{i=0}^b \frac{1}{9^i} \binom{b}{i} \binom{l-b}{a-2b-i} & \text{if } a \leq 3l, \\ 0 & \text{otherwise.} \end{cases} \quad (11)$$

*Proof.* Clearly this configuration requires  $0 \leq b \leq a \leq 3l$ . As above, we denote the  $a$  smallest elements in the array by “•” and the other  $(3l - a)$  entries by “o”. With the notation introduced before, four types of triplets can occur in  $\pi$ , namely  $t_0 : \{\bullet, \bullet, \bullet\}$ ,  $t_1 : \{\bullet, \bullet, o\}$ ,  $t_2 : \{\bullet, o, o\}$  and  $t_3 : \{o, o, o\}$ , where the order of the elements inside each triplet is not relevant at the moment.

Let denote by  $\langle \pi \rangle_i$  the number of triplets in  $\pi$  of type  $t_i$ , with  $0 \leq i \leq 3$ . All the permutations  $\pi$  contributing to  $Q_l^{[3]}(a, b)$  need to satisfy the following constraints:

$$\begin{cases} 3\langle \pi \rangle_0 + 2\langle \pi \rangle_1 + \langle \pi \rangle_2 = a, \\ \langle \pi \rangle_1 + 2\langle \pi \rangle_2 + 3\langle \pi \rangle_3 = 3l - a, \\ \langle \pi \rangle_0 + \langle \pi \rangle_1 = b. \end{cases} \quad (12)$$

The four numbers are not independent, of course, since they must all sum to  $l$ ; however, since adding the first two rows yields this sum,  $\sum_{i=0}^3 \langle \pi \rangle_i = l$ , the system (12) is undetermined. Specifically, we can describe the space of solutions as spanned by fixing  $\langle \pi \rangle_0$  at any value in the interval  $[0, b]$ .

Let us define the following equivalence relation<sup>3</sup> on the set  $S_{3l}$ , of the  $(3l)!$  permutations:

$$\pi \sim \pi' \iff \langle \pi \rangle_i = \langle \pi' \rangle_i, \quad \text{for each } i \in [0, 3],$$

for any  $\pi, \pi' \in S_{3l}$ . We denote by  $\tilde{\pi} \in S_{3l}/\sim$  the  $\sim$ -equivalence class of a permutation  $\pi \in S_{3l}$ . Moreover, by solving the system (12) in terms of  $\langle \pi \rangle_0$  we see that all the permutations  $\pi \in S_{3l}$  contributing to  $Q_l^{[3]}(a, b)$  can be gathered together in the following set of equivalence classes

$$\tilde{\Pi}_{a,b}^{(l)} = \left\{ \tilde{\pi} \in S_{3l}/\sim \mid 0 \leq \langle \pi \rangle_0 \leq b \quad \text{and} \quad \begin{pmatrix} \langle \pi \rangle_1 = b - \langle \pi \rangle_0 \\ \langle \pi \rangle_2 = a - 2b - \langle \pi \rangle_0 \\ \langle \pi \rangle_3 = l - a + b + \langle \pi \rangle_0 \end{pmatrix} \right\}. \quad (13)$$

The note about the solutions above implies that  $|\tilde{\Pi}_{a,b}^{(l)}| = b + 1$ . We start by computing the size  $|\tilde{\pi}|$  for each equivalence class  $\tilde{\pi} \in S_{3l}/\sim$ :

**Lemma 7.**

$$|\tilde{\pi}| = l! \times 6^l \times \prod_{i=0}^3 C_i(\tilde{\pi}),$$

<sup>3</sup>We do not use any of the standard properties of such a relation, except that its classes partition  $S_{3l}$ .

where

$$\begin{aligned}
C_0(\tilde{\pi}) &= \frac{1}{\langle \pi \rangle_0!} \binom{a}{\{3\}^{\langle \pi \rangle_0}} = \frac{1}{\langle \pi \rangle_0!} \prod_{i=0}^{\langle \pi \rangle_0 - 1} \binom{3\langle \pi \rangle_0 + 2\langle \pi \rangle_1 + \langle \pi \rangle_2 - 3i}{3}, \\
C_1(\tilde{\pi}) &= \frac{1}{\langle \pi \rangle_1!} \binom{2\langle \pi \rangle_1 + \langle \pi \rangle_2}{\{2\}^{\langle \pi \rangle_1}} \times \langle \pi \rangle_1! = \frac{1}{\langle \pi \rangle_1!} \prod_{i=0}^{\langle \pi \rangle_1 - 1} \binom{2\langle \pi \rangle_1 + \langle \pi \rangle_2 - 2i}{2} \langle \langle \pi \rangle_1 - i \rangle, \\
C_2(\tilde{\pi}) &= \frac{1}{\langle \pi \rangle_2!} \binom{\langle \pi \rangle_1 + 2\langle \pi \rangle_2}{\{2\}^{\langle \pi \rangle_2}} \times \langle \pi \rangle_2! = \frac{1}{\langle \pi \rangle_2!} \prod_{i=0}^{\langle \pi \rangle_2 - 1} \binom{\langle \pi \rangle_1 + 2\langle \pi \rangle_2 - 2i}{2} \langle \langle \pi \rangle_2 - i \rangle, \\
C_3(\tilde{\pi}) &= \frac{1}{\langle \pi \rangle_3!} \binom{3l - a}{\{3\}^{\langle \pi \rangle_3}} = \frac{1}{\langle \pi \rangle_3!} \prod_{i=0}^{\langle \pi \rangle_3 - 1} \binom{\langle \pi \rangle_1 + 2\langle \pi \rangle_2 + 3\langle \pi \rangle_3 - 3i}{3}.
\end{aligned}$$

*Proof.* The  $C_i(\tilde{\pi})$  have the following interpretation:

$C_0(\tilde{\pi})$  counts the ways to distribute  $3\langle \pi \rangle_0$  “•” elements out of the  $a$  available into  $\langle \pi \rangle_0$  boxes containing triplets of type  $t_0$ ; the factor  $\frac{1}{\langle \pi \rangle_0!}$  is inserted because we are not considering now the order among boxes. Similarly,  $C_3(\tilde{\pi})$  enumerates the ways to distribute  $3\langle \pi \rangle_3$  “o” elements out of  $(3l - a)$  into  $\langle \pi \rangle_3$  unordered boxes containing triplets of type  $t_3$ .

In  $C_1(\tilde{\pi})$  we count the ways to populate  $\langle \pi \rangle_1$  unordered boxes, each with 2 of the  $(a - 3\langle \pi \rangle_0)$  remaining “•” per box and one of the  $\langle \pi \rangle_1$  remaining “o” elements.

$C_2(\tilde{\pi})$  is defined in a completely analogous manner.

A factor  $6^l$  takes into account the order of the elements inside each box, and a factor  $l!$  is finally introduced to complete the specification of the permutation by accounting for the order of the boxes. ■

To return to proving lemma 6 we apply identity (9) and get after several cancellations:

$$\prod_{i=0}^3 C_i(\tilde{\pi}) = \frac{a!(3l - a)!}{2^{\langle \pi \rangle_1 + \langle \pi \rangle_2} 3^{\langle \pi \rangle_0 + \langle \pi \rangle_3} \prod_{i=0}^3 \langle \langle \pi \rangle_i \rangle!}$$

Now we can calculate the quantity of our interest; by the definition of  $\tilde{\Pi}_{a,b}^{(l)}$  we get:

$$\begin{aligned}
Q_l^{[3]}(a, b) &= |\{\pi \in S_{3l} \mid \tilde{\pi} \in \tilde{\Pi}_{a,b}^{(l)}\}| = \sum_{\tilde{\pi} \in \tilde{\Pi}_{a,b}^{(l)}} |\tilde{\pi}| \\
&= l! a! (3l - a)! \sum_{\tilde{\pi} \in \tilde{\Pi}_{a,b}^{(l)}} \frac{3^{\langle \pi \rangle_1 + \langle \pi \rangle_2}}{\prod_{i=0}^3 \langle \langle \pi \rangle_i \rangle!},
\end{aligned}$$

We use the relation in Eq. (13) to rewrite the terms the last row, and since each  $\tilde{\pi} \in \tilde{\Pi}_{a,b}^{(l)}$  can be seen as corresponding to a value of  $\langle \pi \rangle_0$  in  $[0, b]$ , the sum there can be changed to one on this index, and we get the desired result. ■

Notice that, for any  $l \in \mathbb{N}_0$ , we have  $Q_l^{[3]}(0, 0) = (3l)!$  and  $Q_l^{[3]}(a, 0)$  is not necessarily null when  $a \neq 0$ .

Again we derive a useful consequence of the last result:

**Definition 4.** Given  $a, b, l \in \mathbb{N}$ , we denote by  $\overline{Q}_l^{[3]}(a, b)$  the number of permutations  $\pi$ , out of the  $(3l)!$  possible ones, in which the  $a^{\text{th}}$  smallest element in  $\pi$  is selected and becomes the  $b^{\text{th}}$  smallest among the surviving elements.

**Lemma 8.** Let  $a, b, l \in \mathbb{N}$ . The number of permutations  $\overline{Q}_l^{[3]}(a, b)$  is given by

$$\overline{Q}_l^{[3]}(a, b) = 6l(a-1)(3l-a)Q_{l-1}^{[3]}(a-2, b-1).$$

*Proof.* All the permutations contributing to  $\overline{Q}_l^{[3]}(a, b)$  are among those we counted in  $Q_l^{[3]}(a, b)$ , but we only need those in which  $a$  ends up being selected. as before, since  $a$  element is the largest “•” in  $\pi$  it gets selected if  $\pi$  puts it in a triplet of type  $t_1$ . Such a triplet has 2 “•” elements, for a total of  $2\langle\pi\rangle_1$  out of  $a$  such elements. In the formula (11) we assign  $\langle\pi\rangle_0 = i$ , hence  $\langle\pi\rangle_1 = b - i$ , and therefore

$$\overline{Q}_l^{[3]}(a, b) = 3^{a-b} a! (3l-a)! \binom{l}{b} \sum_{i=0}^b \frac{1}{9^i} 2 \frac{b-i}{a} \binom{b}{i} \binom{l-b}{a-2b-i} \quad (14)$$

Rearranging the terms above we obtain the needed result. ■

By convention, we extend Def. 4 letting  $\overline{Q}_0^{[3]}(a, b) = 0$  for any  $a, b \in \mathbb{N}$ .

We are ready to define the next step in our analysis, combining all the results obtained so far.

An ordered sequence of consecutive  $l_0, l_1, l_2$  boxes is given, where each box contains exactly either three, two or one ordered elements, respectively (see Fig. 3). Initially,  $n = 3l_0 + 2l_1 + l_2$  totally ordered elements are uniformly distributed in the boxes, creating a random permutation  $\pi$ . Following the process of our algorithm, we select the median element in each triplet (cf. Def. 3) and the smallest element in each box containing a couple (cf. Def. 1). Thus,  $(l_0 + l_1)$  elements survive at this stage, including none of the singletons contained in the last  $l_2$  boxes, as shown in Fig. 4.

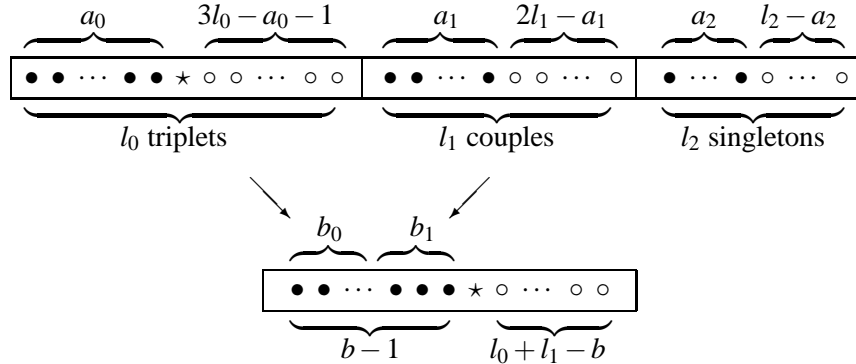


Figure 4: Looking at the elements array as a sequence of triplets, couples and singletons.

**Definition 5.** Given  $l_0, l_1, l_2 \in \mathbb{N}_0$  and  $a, b \in \mathbb{N}$ , we denote by  $\overline{Q}_{l_0, l_1, l_2}(a, b)$  the number of permutations  $\pi$ , out of the  $n!$  possible ones, in which the  $a^{\text{th}}$  smallest element in  $\pi$  is selected and becomes the  $b^{\text{th}}$  smallest among the surviving elements.

**Lemma 9.** Let  $l_0, l_1, l_2 \in \mathbb{N}_0$  and  $a, b \in \mathbb{N}$ . The number of permutations  $\overline{Q}_{l_0, l_1, l_2}(a, b)$  is given by

$$\overline{Q}_{l_0, l_1, l_2}(a, b) = l_2! \sum_{\substack{a_i \in [0, \dots, (3-i)l_i], i=0,2 \\ a_0 + a_1 + a_2 = a-1 \\ b_i \in [0, \dots, \min(a_i, l_i)], i=0,1 \\ b_0 + b_1 = b-1}} \binom{a-1}{a_0} \binom{a-a_0-1}{a_1} (S_1 + S_2),$$



where,

$$\begin{aligned} S_1 &= \binom{n-a}{3l_0-a_0-1} \binom{n-a-3l_0+a_0+1}{2l_1-a_1} \cdot \overline{Q}_{l_0}^{[3]}(a_0+1, b_0+1) \cdot Q_{l_1}^{[2]}(a_1, b_1), \\ S_2 &= \binom{n-a}{3l_0-a_0} \binom{n-a-3l_0+a_0}{2l_1-a_1-1} \cdot Q_{l_0}^{[3]}(a_0, b_0) \cdot \overline{Q}_{l_1}^{[2]}(a_1+1, b_1+1). \end{aligned}$$

*Proof.* In a generic permutation  $\pi$ , among the  $n$  elements, there are  $a$  “•” distributed into the boxes containing triplets, couples and singletons, and  $(n-a)$  “◦”. In particular, excluding the  $a^{\text{th}}$  element “★”,  $a_0$  “•” occurs in the  $l_0$  triplets (with  $0 \leq a_0 \leq 3l_0$ ),  $a_1$  “•” occurs in the  $l_1$  couples (with  $0 \leq a_1 \leq 2l_1$ ) and  $a_2$  “•” forms the  $l_2$  singletons (with  $0 \leq a_2 \leq l_2$ ), such that  $a_0 + a_1 + a_2 = a - 1$ . This situation is depicted in Fig. 4, where each rectangle does not preserve the order of the elements inside. In order to get the element “★” to be selected and to be the  $b^{\text{th}}$  among the surviving elements, we need  $b_0$  “•” survived from triplets (with  $0 \leq b_0 \leq \min(a_0, l_0)$ ) and  $b_1$  “•” survived from couples (with  $0 \leq b_1 \leq \min(a_1, l_1)$ ), such that  $b_0 + b_1 = b - 1$ .

To evaluate  $\overline{Q}_{l_0, l_1, l_2}(a, b)$  we just need to count the permutations that satisfy the above constraints. First we count the ways to fill up the  $a_0$  positions in the triplets with  $(a-1)$  “•” and the  $a_1$  positions in the couples with the remaining  $(a-a_0-1)$  “•”. Secondly, we evaluate the quantities  $S_1$  and  $S_2$  corresponding to the case of “★” belonging to a triplet or to a couple, respectively. In  $S_1$  we count how to fill up the  $(3l_0 - a_0 - 1)$  positions in the triplets with  $(n-a)$  “◦” and the  $(2l_1 - a_1)$  positions in the couples with the remaining  $(n-a-3l_0+a_0+1)$  “◦”; “★” is the  $(a_0+1)^{\text{th}}$  smallest element in the triplets and we need it becomes the  $(b_0+1)^{\text{th}}$  after the triplets are processed; at the same time, we need  $b_1$  “•” selected among the  $a_1$  smallest elements in the couples. The quantity  $S_2$  is derived analogously.

Finally, the factor  $l_2!$  takes into account the remaining elements forming the  $l_2$  singletons.  $\blacksquare$

The performance of the algorithm is characterized by the probability  $P_a^{(n,m)}$ , that the  $a^{\text{th}}$  smallest element of the input array of size  $n$  is returned in output as final winner, when the  $k$ -th element is searched for, and then  $m = |n - 2k + 1|$ .

Let  $n_i, m_i$ , for each  $i \in [0, f]$  be the sequence of values computed during the  $(f+1)$  iterations of the “while” loop of the algorithm, on an input array of size  $n = n_0$  with  $m_0 = |n_0 - 2k + 1|$ ; the value  $f \in \mathbb{N}_0$  is the smallest index such that  $n_{f+1} \leq \text{Threshold}$ .

The previous values are computed by

$$n_{i+1} = \overline{\mu}_0^{(i)} + \overline{\mu}_1^{(i)}, \quad m_{i+1} = \overline{\mu}_2^{(i)} + \overline{\mu}_3^{(i)}, \quad \text{for each } i = 0..f,$$

where  $\langle \overline{\mu}_0^{(i)}, \overline{\mu}_1^{(i)}, \overline{\mu}_2^{(i)}, \overline{\mu}_3^{(i)} \rangle$  are the integer values returned by SystemSolve  $(n_i, m_i)$ .

By Lemma 9, we are able to analyze a single iteration of the algorithm. Indeed, for any values of the integer variables  $n_i, m_i$ , we can derive the probability  $p_{a,b}^{(n_i, m_i)}$  that the  $a^{\text{th}}$  smallest element at  $i$ -th step becomes the  $b^{\text{th}}$  smallest element in the next  $(i+1)$ -st step obtained applying rules in Step (2d) of Section 3:

$$p_{a,b}^{(n_i, m_i)} = \frac{1}{n_i!} \cdot \overline{Q}_{\overline{\mu}_0^{(i)}, \overline{\mu}_1^{(i)}, \overline{\mu}_2^{(i)} + [n_i + m_i]_3}(a, b).$$

A recurrence equation can be established to evaluate the probability of our interest:

$$P_a^{(n_i, m_i)} = \sum_{b_i \in [1..a]} P_{a, b_i}^{(n_i, m_i)} \cdot P_{b_i}^{(n_{i+1}, m_{i+1})}, \quad \text{for all } i = 0..f,$$

and expanding, we get:

$$P_a^{(n,m)} = \sum_{b_0, b_1, \dots, b_{f-1}, b_f} P_{a, b_0}^{(n_0, m_0)} \cdot P_{b_0, b_1}^{(n_1, m_1)} \cdots P_{b_{f-1}, b_f}^{(n_f, m_f)} \cdot P_{b_f}^{(n_{f+1}, m_{f+1})}, \quad (15)$$

where,

$$P_a^{(n_{f+1}, m_{f+1})} = \begin{cases} 1 & \text{if } \left( \begin{array}{l} n_{f+1} \leq m_{f+1} \text{ and } a = 1 \\ \text{or} \\ n_{f+1} > m_{f+1} \text{ and } a = \left\lceil \frac{n_{f+1} - m_{f+1}}{2} \right\rceil \end{array} \right) \\ 0 & \text{otherwise.} \end{cases}, \quad (16)$$

It is straightforward that probabilities given by (15) and (16) can be symmetrically defined for  $\lceil \frac{n}{2} \rceil < k \leq n$ , so that we can define the probability distribution given by (8), for each  $1 \leq k \leq n$ :

$$\mathcal{P}_{n,k}(z) = \sum_{\substack{1 \leq a \leq n \\ |a-k| \leq \lfloor zn \rfloor}} P_a^{(n, |n-2k+1|)}, \quad \text{with } z \in [0, 1]. \quad (17)$$

While these expressions for the distribution of the returned value are far too complex to allow for an analytic characterization, generating numbers from them is possible, as we show in the next section, which also demonstrates how closely they track the actual performance of the algorithm.

## 7. Experimental results

In this section we present empirical results, demonstrating the effectiveness of the algorithm described in Section 4. Our implementation is in standard C (GNU C compiler v2.7). All the experiments were carried out on a PC Pentium II 350Mhz with the Linux (Red Hat distribution) operating system. The arrays were permuted using the pseudo-random number generator suggested by Park and Miller in 1988 and updated in 1993 [10]. The algorithm was run on random arrays of sizes  $n = 10^i$ , with  $i = 2..5$ , whose entry keys were always the integers  $1, \dots, n$ .

Our first experiment aims to evaluate how the theoretical probability distribution (15) fits with the empirical one. In order to do this, we chose a sample of 100,000 random input arrays of size  $n = 100$  and run the algorithm on it for some  $k$  values chosen between 1 and  $\lceil \frac{n}{2} \rceil$ , with a fixed *Threshold* equal to  $\lfloor \sqrt{n} \rfloor$ . We compute the histograms of the relative frequencies of the algorithm outputs so that theoretical and empirical distributions are compared in Fig. 5, where, for each chosen  $k$ , we reported the average  $\mu_k$  and the standard deviation  $\sigma_k$  of both the distributions.

A more interesting experiment was done with variable input sizes  $n = 10^i$ ,  $i = 2..5$ , with *Threshold* equal to  $\lfloor \sqrt{n} \rfloor$ . Also in this case we fixed some  $k$  values among  $\lfloor \frac{n}{8} \rfloor$ ,  $\lfloor \frac{3n}{8} \rfloor$ ,  $\lfloor \frac{n}{2} \rfloor$  and  $\lfloor \frac{3n}{4} \rfloor$ , with the only exception for  $n = 100$ , and we built the histograms of the relative frequencies of the algorithm outputs. Each single histogram, shown in Fig. 6, has been obtained on a different sample of size 100,000.

The *Threshold* value  $\lfloor \sqrt{n} \rfloor$  until now used was completely arbitrary. In Fig. 7 we report some histograms referred to the relative frequencies of algorithm outputs for three different *Threshold* values, i.e. 3,  $\lfloor \sqrt{n} \rfloor$ ,  $2 \lfloor \sqrt{n} \rfloor$ . Data are obtained from a fixed sample of 100,000 arrays of size  $n = 1,000$  and for some fixed  $k$  values. By examining the standard deviation  $\sigma_k$  values, it is clear that larger is the *Threshold*

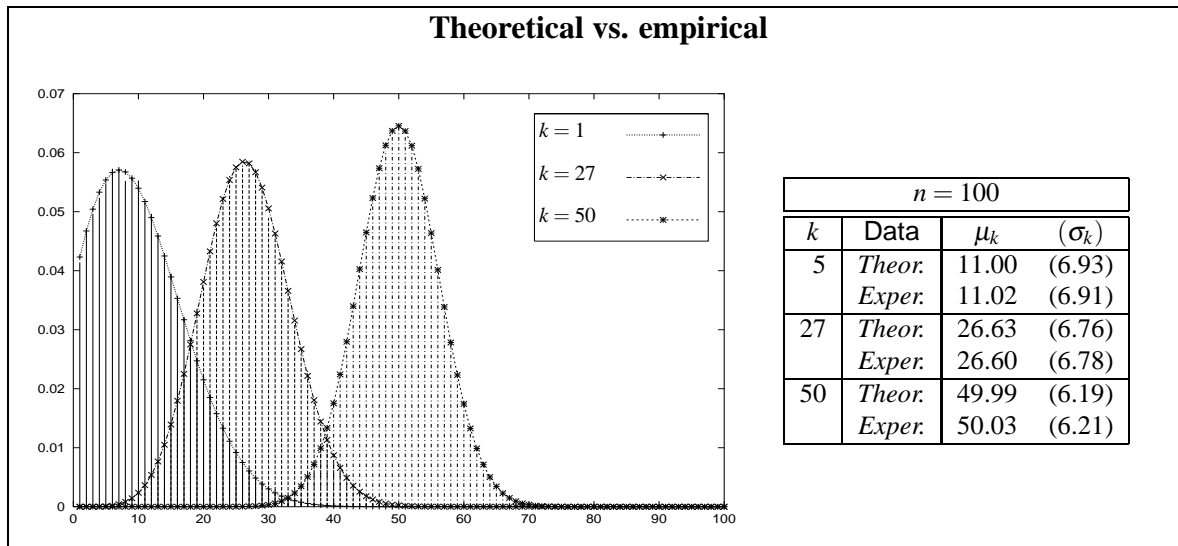


Figure 5: Comparing the theoretical probability distributions with the experimental results histograms.

value, larger is the computational cost of the selection and better is the quality of the selection. For this reason, depending on the specific application, the *Threshold* value has to be appropriately tuned.

In the previous experiments, we worked on some representative  $k$  values arbitrarily chosen. Now, we concentrate our attention on the percentage error made by the algorithm, independently by  $k$ . For each  $i = 2..5$ , we fixed a sample of 1,000 arrays of size  $n = 10^i$  and we used a *Threshold* value  $\lfloor \sqrt{n} \rfloor$ . We run the algorithm on the input sample for each  $k$  value, such that  $1 \leq k \leq n$ , with the only exception for  $i = 4, 5$  for which only 1,000 uniformly distributed  $k$  values were examined, i.e. values  $k = 1 + j \cdot 10^{i-3}$  with  $j = 0..999$ . For each single algorithm output, corresponding to a value  $k$  and to an array of the sample, we calculate the percentage error distance

$$d_k = \frac{|k - \text{rank}(\text{Output})|}{n-1} \cdot 100, \quad \text{with } 1 \leq k \leq n,$$

from the correct  $k$ -th element. The extreme values assumed by  $d_k$  can be 0%, when the  $k$ -th element is returned by the algorithm, and 100% when the algorithm returns the largest element in the array whereas the smallest one is searched for, or viceversa. In Fig. 8, for each single  $k$ , we plotted the maximum error  $M_k = \max(d_k)$ , the average error  $\epsilon_k = E[d_k]$  and the standard deviation  $\sigma_k = \sigma[d_k]$  of the 1,000  $d_k$  values.

Finally, in Fig. 9 we give a plot that summarizes some of the statistical indexes obtained from the previous data-sets. It well illustrates the good quality of the algorithm output when the input size  $n$  is sufficiently large. More specifically, the trend of the standard deviation  $E[\sigma_k]$  may be viewed as a measure of the improvement of the selection effectiveness for increasing  $n$ .

## 8. Conclusion

We have presented an approximate algorithm for the  $k$ -th selection problem, based on a statistical approach. It is a generalization of the algorithm for the approximate median selection problem described in

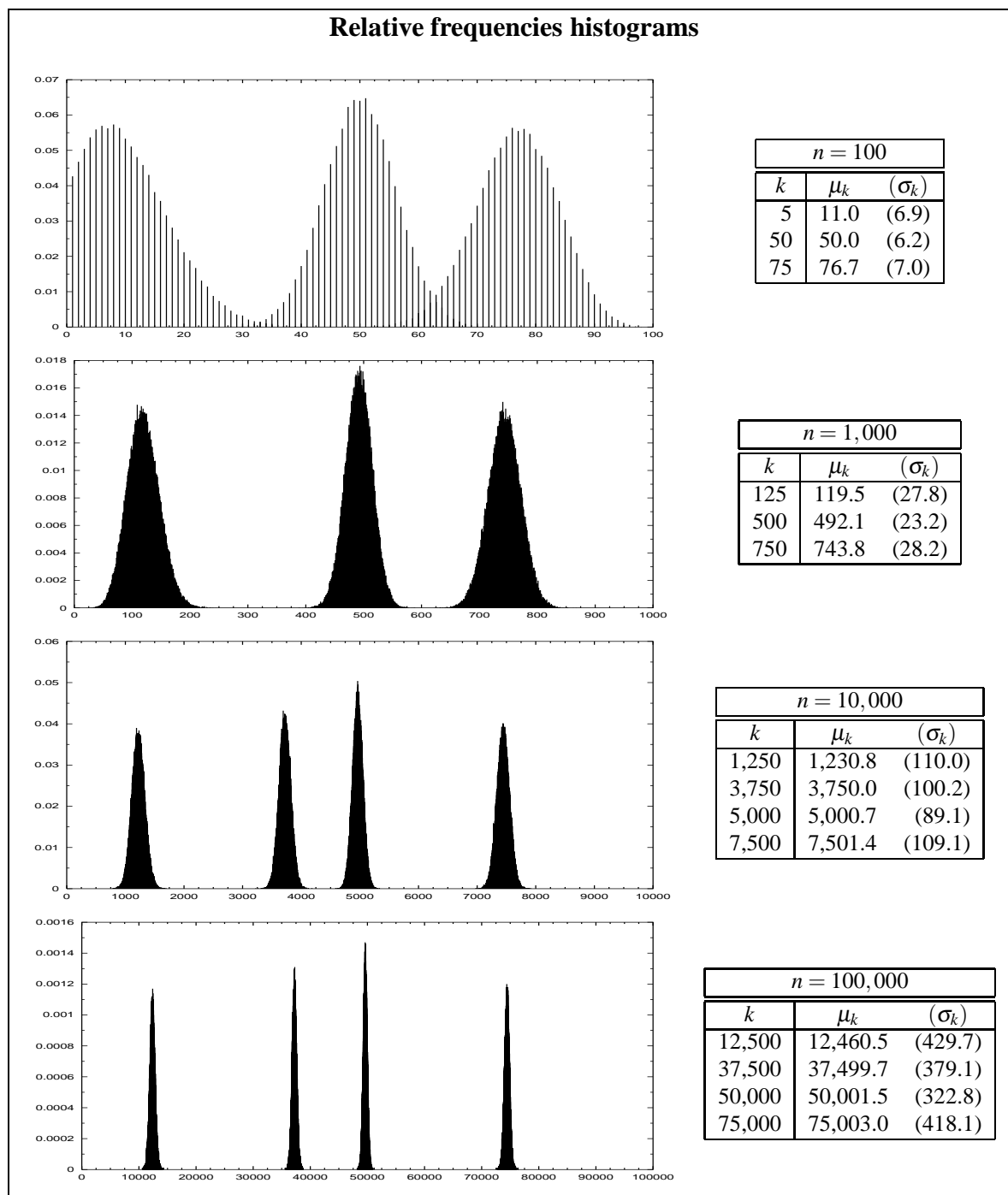


Figure 6: Relative frequencies for  $n = 10^i, i = 2..5$  and some fixed  $k$  values.

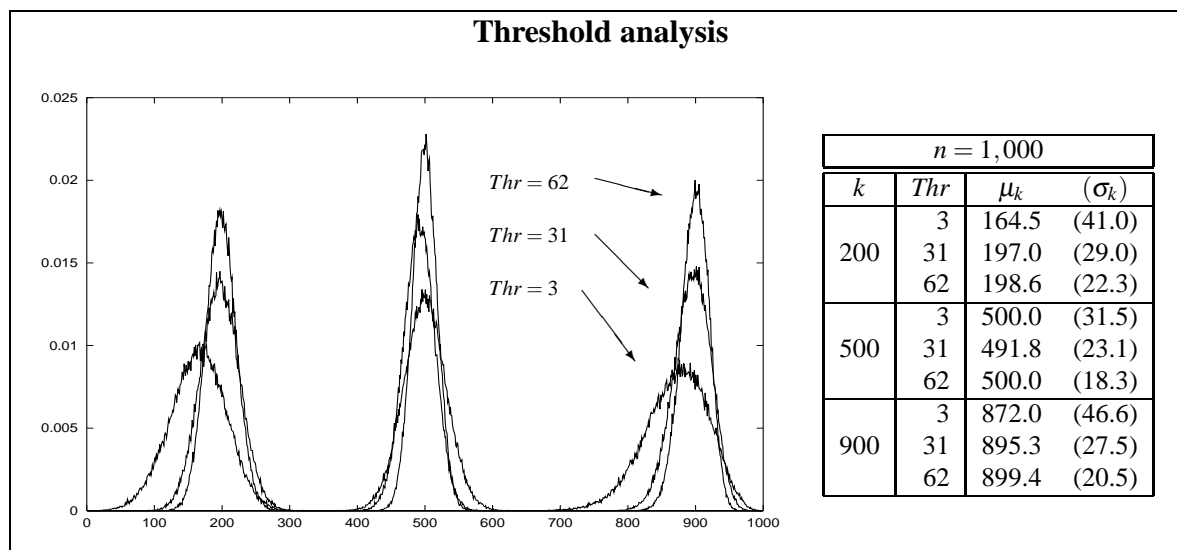


Figure 7: Quality of selection as a function of the *Threshold* value.

[1]. Its precision is strongly based on the good behavior of the approximate median selection algorithm. As the experimental section shows, the quality of the algorithm output becomes better and better for larger input sizes.

In order to get an improving in the quality of results, our future research is oriented to find efficient ways to manipulate larger blocks rather than triplets when selecting for the  $k$ -th element. In particular, an alternative method to treat the extremes of the range of selection to get higher precision is under investigation. We are also studying how to apply the technique presented here for the linear case to other higher dimension domains.

## References

- [1] S. Battiato, D. Cantone, D. Catalano, G. Cincotti, and M. Hofri. An efficient algorithm for the approximate median selection problem. Proceedings of the 4th Italian Conference on Algorithms and Complexity (CIAC 2000, Rome, Italy). *Lecture Notes in Computer Science*, Vol. 1767, pp. 226–238, 2000.
- [2] M. Blum, R.W. Floyd, V. Pratt, R.L. Rivest, and R. Tarjan. Time bounds for selection. *Journal of Computer and System Sciences*, Vol. 7, pp. 448–461, 1973.
- [3] S. Bent and J. John. Finding the median requires  $2n$  comparisons. Proceedings of the 17th ACM Symposium on Theory of Computing, pp. 213–216, 1985.
- [4] S. Carlsson and M. Sundstrom. Linear-time in-place selection in less than  $3n$  comparisons. Proceedings of the 6th International Symposium on Algorithms and Computation (ISAAC 95, Cairns, Australia). *Lecture Notes in Computer Science*, Vol. 1004, pp. 245–253, 1995.
- [5] W. Cunto and J.I. Munro. Average case selection. *Journal of the ACM*, Vol. 36, pp. 270–279, 1989.

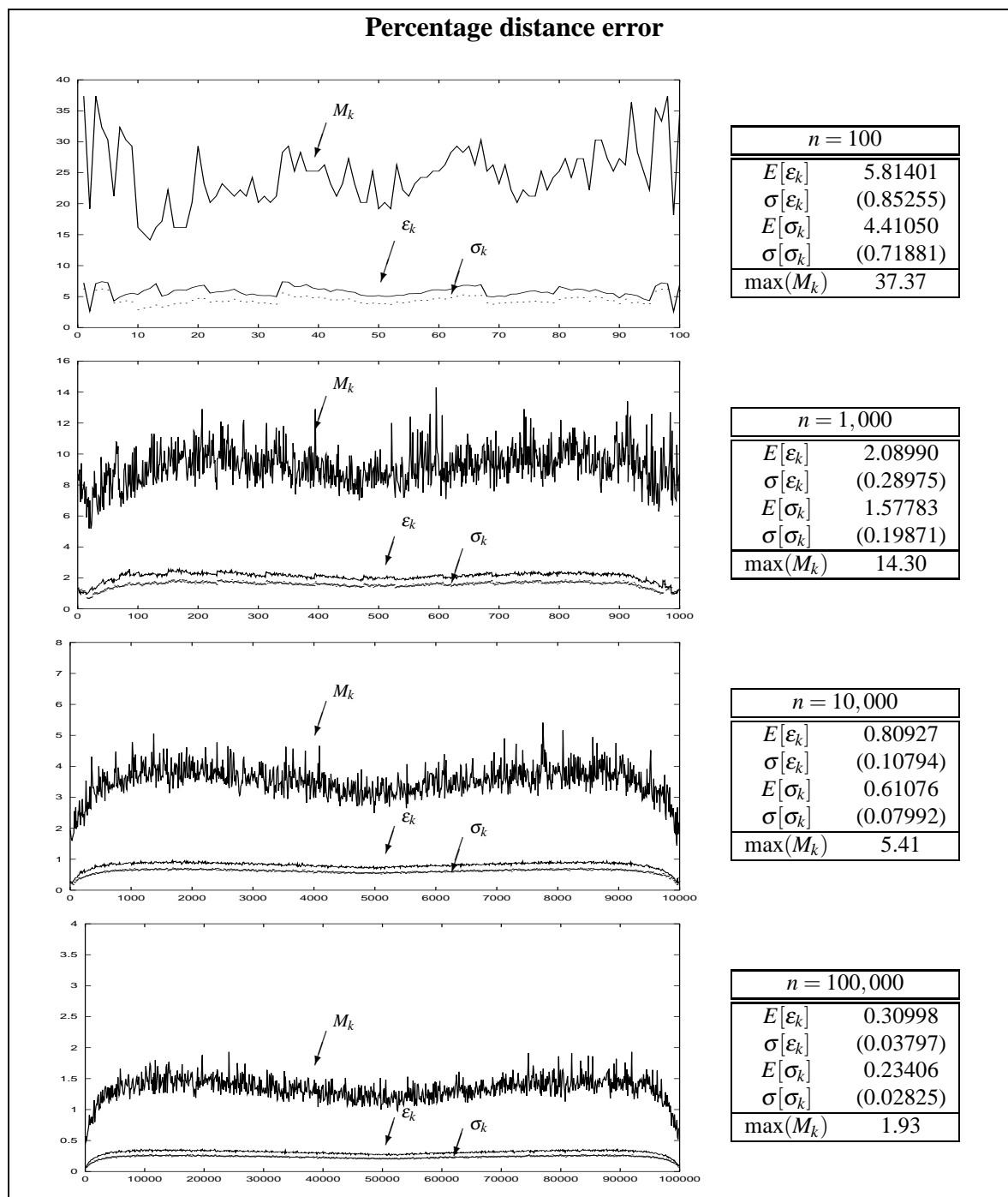


Figure 8: Percentage distance error for  $n = 10^i$ , with  $i = 2..5$ .

[6] D. Dor and U. Zwick. Selecting the median. Proceedings of the 6th ACM-SIAM Symposium on

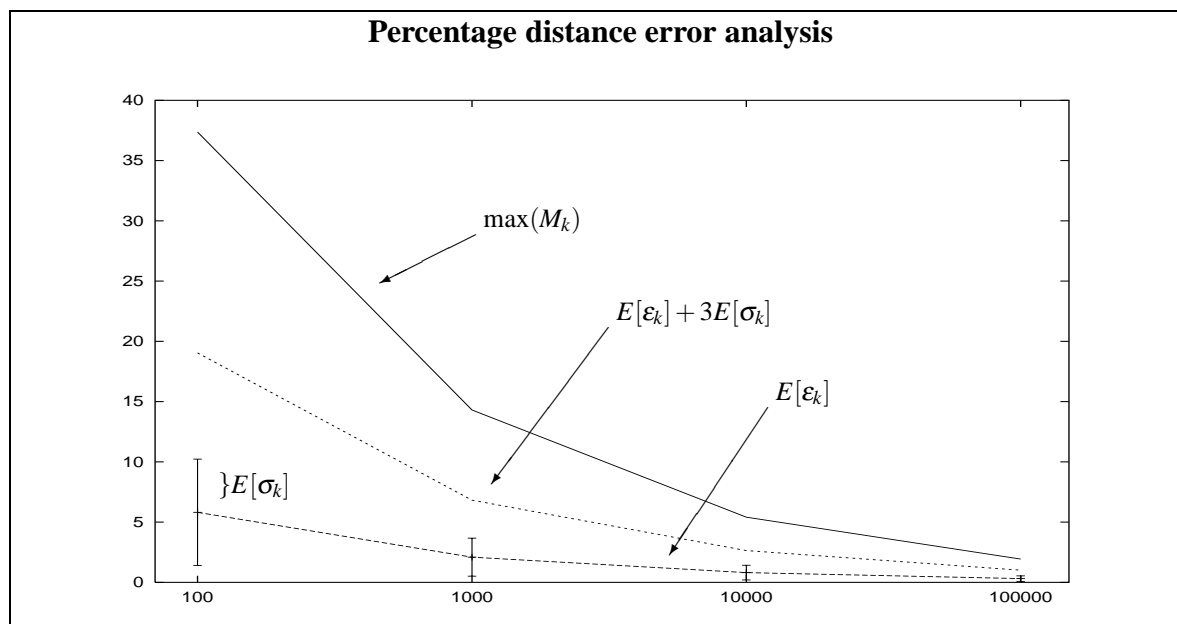


Figure 9: Asymptotically behaviour of percentage distance error w.r.t.  $n$ .

Discrete Algorithms (SODA 95, San Francisco, California). *SIAM Journal on Computing*, Vol. 28, pp. 1722–1758, 1999.

- [7] R.W. Floyd and R.L. Rivest. Expected time bounds for selection. *Communications of the ACM*, Vol. 18, pp. 165–172, 1975.
- [8] C.A.R. Hoare. Algorithm 63 (partition) and algorithm 65 (find). *Communications of the ACM*, Vol. 4, pp. 321–322, 1961.
- [9] T.W. Lai and D. Wood. Implicit Selection. Proceedings of the Scandinavian Workshop on Algorithm Theory (SWAT 88, Halmstad, Sweden). *Lecture Notes in Computer Science*, Vol. 318, pp. 18–23, 1988.
- [10] S. K. Park and K. W. Miller. Random number generators: good ones are hard to find. *Communications of the ACM*, Vol. 31, pp. 1192–1201, 1988. Updated in *Communications of the ACM*, Vol. 36, pp. 108–110, 1993.
- [11] A. Schönhage, M. Paterson, and N. Pippenger. Finding the median. *Journal of Computer and System Sciences*, Vol. 13, pp. 184–199, 1976.