

36/2002

Raport Badawczy
Research Report

RB/80/2002

**Randomized Selection
with Tripartitioning**

K.C. Kiwiel

Instytut Badań Systemowych
Polska Akademia Nauk

Systems Research Institute
Polish Academy of Sciences



POLSKA AKADEMIA NAUK

Instytut Badań Systemowych

ul. Newelska 6

01-447 Warszawa

tel.: (+48) (22) 8373578

fax: (+48) (22) 8372772

Kierownik Pracowni zgłaszający pracę:
Prof. dr hab. inż. Krzysztof C. Kiwiel

Warszawa 2002

Randomized selection with tripartitioning

Krzysztof C. Kiwiel*

December 15, 2002

Abstract

We show that several versions of Floyd and Rivest's algorithm SELECT [Comm. ACM 18 (1975) 173] for finding the k th smallest of n elements require at most $n + \min\{k, n - k\} + o(n)$ comparisons on average, even when equal elements occur. This parallels our recent analysis of another variant due to Floyd and Rivest [Comm. ACM 18 (1975) 165–172]. Our computational results suggest that both variants perform well in practice, and may compete with other selection methods, such as Hoare's FIND or quickselect with median-of-three partitioning.

Key words. Selection, medians, partitioning, computational complexity.

1 Introduction

The *selection problem* is defined as follows: Given a set $X := \{x_j\}_{j=1}^n$ of n elements, a total order $<$ on X , and an integer $1 \leq k \leq n$, find the k th *smallest* element of X , i.e., an element x of X for which there are at most $k - 1$ elements $x_j < x$ and at least k elements $x_j \leq x$. The *median* of X is the $\lceil n/2 \rceil$ th smallest element of X .

Selection is one of the fundamental problems in computer science; see, e.g., the references in [DHUZ01, DoZ99, DoZ01] and [Knu98, §5.3.3]. Most references concentrate on the number of comparisons between pairs of elements made in selection algorithms. In the worst case, selection needs at least $(2 + \epsilon)n$ comparisons [DoZ01], whereas the algorithm of [BFP⁺72] makes at most $5.43n$, that of [SPP76] needs $3n + o(n)$, and that in [DoZ99] takes $2.95n + o(n)$. In the average case, for $k \leq \lceil n/2 \rceil$, at least $n + k - O(1)$ comparisons are necessary [CuM89], whereas the best upper bound is $n + k + O(n^{1/2} \ln^{1/2} n)$ [Knu98, Eq. (5.3.3.16)]. The classical algorithm FIND of [Hoa61], also known as quickselect, has an upper bound of $3.39n + o(n)$ for $k = \lceil n/2 \rceil$ in the average case [Knu98, Ex. 5.2.2–32], which improves to $2.75n + o(n)$ for median-of-three partitioning [Grü99, KMP97].

In practice FIND is most popular. One reason is that the algorithms of [BFP⁺72, SPP76] are much slower on the average [Mus97, Val00], whereas [KMP97] adds that other methods proposed so far, although better than FIND in theory, are not practical because they are difficult to implement, their constant factors and hidden lower order terms are

*Systems Research Institute, Newelska 6, 01-447 Warsaw, Poland (kiwiel@ibspan.waw.pl)

too large, etc. It is quite surprising that these references [KMP97, Mus97, Val00] ignore the algorithm SELECT of [FIR75b], since most textbooks mention that SELECT is asymptotically faster than FIND. In contrast, this paper shows that SELECT can compete with FIND in both theory and practice, even for moderate values of the input size n .

We now outline our contributions in more detail. The initial two versions of SELECT [FIR75b] had gaps in their analysis (cf. [Bro76, PRKT83], [Knu98, Ex. 5.3.3–24]); the first version was validated in [Kiw02], and the second one will be addressed elsewhere. This paper deals with the third version of SELECT from [FIR75a], which operates as follows. Using a small random sample, it finds an element v almost sure to be just above the k th if $k < n/2$, or below the k th if $k \geq n/2$. Partitioning X about v leaves $\min\{k, n - k\} + o(n)$ elements on average for the next recursive call, in which k is near 1 or n with high probability, so this second call eliminates almost all the remaining elements.

Apparently this version of SELECT has not been analyzed in the literature, even in the case of distinct elements. We first revise it slightly to simplify our analysis. Then, without assuming that the elements are distinct, we show that SELECT needs at most $n + \min\{k, n - k\} + O(n^{2/3} \ln^{1/3} n)$ comparisons on average, with $\ln^{1/3} n$ replaced by $\ln^{1/2} n$ for the original samples of [FIR75a]. Thus the average cost of SELECT reaches the lower bounds of $1.5n + o(n)$ for median selection and $1.25n + o(n)$ for selecting an element of random rank. For the latter task, FIND has the bound $2n + o(n)$ when its pivot is set to the median of a random sample of s elements, with $s \rightarrow \infty$, $s/n \rightarrow \infty$ as $n \rightarrow \infty$ [Mar01]; thus SELECT improves upon FIND mostly by using k , the rank of the element to be found, for selecting the pivot v in each recursive call.

In principle, SELECT can be implemented like FIND by using any well-known bipartitioning scheme [Sed77] (an enhancement of the scheme of [FIR75a] is given in §6.3). However, such schemes can perform quite poorly when equal elements occur, in which case the ternary scheme of [BeM93, BeS97] may be preferred. This scheme works rather well in practice, but we present a faster ternary scheme that obviates subscript range checking. Our scheme is only slightly slower than binary schemes when the elements are distinct; it thus combines reliability and efficiency. We add that the implementation of [FIR75a], like several popular implementations of FIND, avoids random number generation by assuming that the input file is in random order, but this results in poor performance on some inputs of [Val00]; hence our implementation of SELECT employs random sampling.

Our computational experience shows that SELECT outperforms even quite sophisticated implementations of FIND in both comparison counts and computing times. To save space, only selected results are reported for the version of [Val00], but our experience with other versions on many different inputs was similar. SELECT turned out to be more stable than FIND, having much smaller variations of solution times and numbers of comparisons. Quite surprisingly, contrary to the folklore saying that SELECT is only asymptotically faster than FIND, SELECT makes significantly fewer comparisons even for small inputs.

To relate our results with those of [Kiw02], let's call QSELECT the quintary method of [Kiw02] stemming from [FIR75b, §2.1]. QSELECT eliminates almost all elements on its first call by using two pivots, almost sure to be just below and above the k th element, in a quintary partitioning scheme. Thus most work occurs on the first call of QSELECT, which corresponds to the first two calls of SELECT. Hence SELECT and QSELECT share

the same efficiency estimates, and in practice make similarly many comparisons. However, QSELECT tends to be slightly faster on median finding: although its quintary scheme is more complex, most of its work is spent on the first pass through X , whereas SELECT first partitions X and then the remaining part (about half) of X on its second call to achieve a similar problem reduction. On the other hand, SELECT makes fewer comparisons on small inputs. Of course, future work should assess more fully the relative merits of SELECT and QSELECT. For now, the tests reported in [Kiw02] and in §7 suggest that both SELECT and QSELECT can compete successfully with refined implementations of FIND [Val00].

The paper is organized as follows. A general version of SELECT is introduced in §2, and its basic features are analyzed in §3. The average performance of SELECT is studied in §4. A modification that improves practical performance is introduced in §5. Partitioning schemes are discussed in §6. Finally, our computational results are reported in §7.

Our notation is fairly standard. $|A|$ denotes the cardinality of a set A . In a given probability space, P is the probability measure, E is the mean-value operator and $P[\cdot|\mathcal{E}]$ is the probability conditioned on an event \mathcal{E} ; the complement of \mathcal{E} is denoted by \mathcal{E}' .

2 The algorithm SELECT

In this section we describe a general version of SELECT in terms of two auxiliary functions $s(n)$ and $g(n)$ (the sample size and rank gap), which will be chosen later. We omit their arguments in general, as no confusion can arise.

Algorithm 2.1.

SELECT(X, k) (Selects the k th smallest element of X , with $1 \leq k \leq n := |X|$)

Step 1 (*Initiation*). If $n = 1$, return x_1 . Choose the sample size $s \leq n - 1$ and gap $g > 0$.

Step 2 (*Sample selection*). Pick randomly a sample $S := \{y_1, \dots, y_s\}$ from X .

Step 3 (*Pivot selection*). Let v be the output of SELECT(S, i_v), where

$$i_v := \begin{cases} \min \{ \lceil ks/n + g \rceil, s \} & \text{if } k < n/2, \\ \max \{ \lceil ks/n - g \rceil, 1 \} & \text{if } k \geq n/2. \end{cases} \quad (2.1)$$

Step 4 (*Partitioning*). By comparing each element x of $X \setminus S$ to v , partition X into the three sets $L := \{x \in X : x < v\}$, $E := \{x \in X : x = v\}$ and $U := \{x \in X : v < x\}$.

Step 5 (*Stopping test*). If $|L| < k \leq |L \cup E|$, return v .

Step 6 (*Reduction*). If $k \leq |L|$, set $\hat{X} := L$, $\hat{n} := |\hat{X}|$ and $\hat{k} := k$; else set $\hat{X} := U$, $\hat{n} := |\hat{X}|$ and $\hat{k} := k - |L \cup E|$.

Step 7 (*Recursion*). Return SELECT(\hat{X}, \hat{k}).

A few remarks on the algorithm are in order.

Remarks 2.2. (a) The correctness and finiteness of SELECT stem by induction from the following observations. The returns of Steps 1 and 5 deliver the desired element. At Step

6, \hat{X} and \hat{k} are chosen so that the k th smallest element of X is the \hat{k} th smallest element of \hat{X} , and $\hat{n} < n$ (since $v \notin \hat{X}$). Also $|S| < n$ for the recursive call at Step 3.

(b) When Step 5 returns v , SELECT may also return information about the positions of the elements of X relative to v . For instance, if X is stored as an array, its k smallest elements may be placed first via interchanges at Step 4 (cf. §6). Hence Step 4 need only compare v with the elements of $X \setminus S$.

(c) The following elementary property is needed in §4. Let c_n denote the maximum number of comparisons taken by SELECT on any input of size n . Since Step 3 makes at most c_s comparisons with $s < n$, Step 4 needs at most $n - s$, and Step 7 takes at most c_n with $\hat{n} < n$, by induction $c_n < \infty$ for all n .

3 Sampling deviations

In this section we analyze general features of sampling used by SELECT. Our analysis hinges on the following bound on the tail of the hypergeometric distribution established in [Hoe63] and rederived shortly in [Chv79].

Fact 3.1. *Let s balls be chosen uniformly at random from a set of n balls, of which r are red, and r' be the random variable representing the number of red balls drawn. Let $p := r/n$. Then*

$$P[r' \geq ps + g] \leq e^{-2g^2/s} \quad \forall g \geq 0. \quad (3.1)$$

Denote by $x_1^* \leq \dots \leq x_n^*$ and $y_1^* \leq \dots \leq y_s^*$ the sorted elements of the input set X and the sample set S , respectively, so that $v = y_{i_v}^*$. The following result will give bounds on the position of v in the sorted input sequence.

Lemma 3.2. *Suppose $\bar{i} := \max\{1, \min(\lceil \kappa s \rceil, s)\}$, $\bar{j}_l := \max\{\lceil \kappa n - gn/s \rceil, 1\}$, and $\bar{j}_r := \min\{\lceil \kappa n + gn/s \rceil, n\}$, where $-g < \kappa s \leq s + g$, $1 \leq s \leq n$ and $g \geq 0$. Then:*

- (a) $P[y_{\bar{i}}^* < x_{\bar{j}_l}^*] \leq e^{-2g^2/s}$ if $\bar{i} \geq \lceil \kappa s \rceil$.
- (b) $P[x_{\bar{j}_r}^* < y_{\bar{i}}^*] \leq e^{-2g^2/s}$ if $\bar{i} \leq \lceil \kappa s \rceil$.

Proof. Note that $-g < \kappa s \leq s + g$ implies that $\bar{j}_l \leq n$ and $\bar{j}_r \geq 1$ are well-defined.

(a) If $y_{\bar{i}}^* < x_{\bar{j}_l}^*$, at least \bar{i} samples satisfy $y_i \leq x_{\bar{j}_l}^*$, where $r := \max_{x_j^* < x_{\bar{j}_l}^*} j$. In the setting of Fact 3.1, we have r red elements $x_j \leq x_{\bar{j}_l}^*$, $ps = rs/n$ and $r' \geq \bar{i}$. Now, $1 \leq r \leq \bar{j}_l - 1$ implies $2 \leq \bar{j}_l = \lceil \kappa n - gn/s \rceil < \kappa n - gn/s + 1$, so $-rs/n > -\kappa s + g$. Hence $\bar{i} - ps - g > \kappa s - \kappa s + g - g = 0$, i.e., $r' > ps + g$. Thus $P[y_{\bar{i}}^* < x_{\bar{j}_l}^*] \leq e^{-2g^2/s}$ by (3.1).

(b) If $x_{\bar{j}_r}^* < y_{\bar{i}}^*$, $s - \bar{i} + 1$ samples are at least $x_{\bar{j}_r}^*$ with $\bar{j} := \max_{x_j^* > x_{\bar{j}_r}^*} j$. Thus we have $r := n - \bar{j}$ red elements $x_j \geq x_{\bar{j}_r}^*$, $ps = s - \bar{j}s/n$ and $r' \geq s - \bar{i} + 1$. Since $\bar{i} < \kappa s + 1$ and $n > \bar{j} \geq \bar{j}_r \geq \kappa n + gn/s$, we get $s - \bar{i} + 1 - ps - g > \bar{j}s/n - \kappa s - g \geq \kappa s + g - \kappa s - g = 0$. Hence $r' > ps + g$ and $P[x_{\bar{j}_r}^* < y_{\bar{i}}^*] \leq P[r' \geq ps + g] \leq e^{-2g^2/s}$ by (3.1). \square

We now bound the position of v relative to x_k^* , $x_{k_l}^*$ and $x_{k_r}^*$, where

$$k_l := \max\{\lceil k - 2gn/s \rceil, 1\} \quad \text{and} \quad k_r := \min\{\lceil k + 2gn/s \rceil, n\}. \quad (3.2)$$

Table 4.1: Sample size $f(n) := n^{2/3} \ln^{1/3} n$ and relative sample size $\phi(n) := f(n)/n$.

n	10^3	10^4	10^5	10^6	$5 \cdot 10^6$	10^7	$5 \cdot 10^7$	10^8
$f(n)$	190.449	972.953	4864.76	23995.0	72287.1	117248	353885	568986
$\phi(n)$.190449	.097295	.048648	.023995	.014557	.011725	.007078	.005690

Corollary 3.3. (a) $P[v < x_k^*] \leq e^{-2g^2/s}$ if $i_v = \lceil ks/n + g \rceil$ and $k < n/2$.

(b) $P[x_{k_r}^* < v] \leq e^{-2g^2/s}$ if $k < n/2$.

(c) $P[x_k^* < v] \leq e^{-2g^2/s}$ if $i_v = \lceil ks/n - g \rceil$ and $k \geq n/2$.

(d) $P[v < x_{k_l}^*] \leq e^{-2g^2/s}$ if $k \geq n/2$.

(e) If $k < n/2$, then $i_v \neq \lceil ks/n + g \rceil$ iff $n < k + gn/s$; similarly, if $k \geq n/2$, then $i_v \neq \lceil ks/n - g \rceil$ iff $k \leq gn/s$.

Proof. Use Lem. 3.2 with $\kappa s = ks/n + g$ for (a,b), and $\kappa s = ks/n - g$ for (c,d). \square

4 Average case performance

In this section we analyze the average performance of SELECT for various sample sizes.

4.1 Floyd-Rivest's samples

For positive constants α and β , consider choosing $s = s(n)$ and $g = g(n)$ as

$$s := \min \{ \lceil \alpha f(n) \rceil, n - 1 \} \text{ and } g := (\beta s \ln n)^{1/2} \text{ with } f(n) := n^{2/3} \ln^{1/3} n. \quad (4.1)$$

This form of g gives a probability bound $e^{-2g^2/s} = n^{-2\beta}$ for Cor. 3.3. To get more feeling, suppose $\alpha = \beta = 1$ and $s = f(n)$. Let $\phi(n) := f(n)/n$. Then $s/n = g/s = \phi(n)$ and it will be seen that the recursive call reduces n at least by the factor $4\phi(n)$ on average, i.e., $\phi(n)$ is a contraction factor; note that $\phi(n) \approx 2.4\%$ for $n = 10^6$ (cf. Tab. 4.1).

Theorem 4.1. Let C_{nk} denote the expected number of comparisons made by SELECT for s and g chosen as in (4.1) with $\beta \geq 1/6$. There exists a positive constant γ such that

$$C_{nk} \leq n + \min\{k, n - k\} + \gamma f(n) \quad \forall 1 \leq k \leq n. \quad (4.2)$$

Proof. We need a few preliminary facts. The function $\phi(t) := f(t)/t = (\ln t/t)^{1/3}$ decreases to 0 on $[e, \infty)$, whereas $f(t)$ grows to infinity on $[2, \infty)$. Let $\delta := 4(\beta/\alpha)^{1/2}$. Pick $\bar{n} \geq 3$ large enough so that $e - 1 \leq \alpha f(\bar{n}) \leq \bar{n} - 1$ and $e \leq \delta f(\bar{n})$. Let $\bar{\alpha} := \alpha + 1/f(\bar{n})$. Then, by (4.1) and the monotonicity of f and ϕ , we have for $n \geq \bar{n}$

$$s \leq \bar{\alpha} f(n) \text{ and } f(s) \leq \bar{\alpha} \phi(\bar{\alpha} f(\bar{n})) f(n), \quad (4.3)$$

$$f(\lfloor \delta f(n) \rfloor) \leq f(\delta f(n)) \leq \delta \phi(\delta f(\bar{n})) f(n). \quad (4.4)$$

For instance, the first inequality of (4.3) yields $f(s) \leq f(\bar{\alpha}f(n))$, whereas

$$f(\bar{\alpha}f(n)) = \bar{\alpha}\phi(\bar{\alpha}f(n))f(n) \leq \bar{\alpha}\phi(\bar{\alpha}f(\bar{n}))f(n).$$

Also for $n \geq \bar{n}$, we have $s = \lceil \alpha f(n) \rceil = \alpha f(n) + \epsilon$ with $\epsilon \in [0, 1)$ in (4.1). Writing $s = \bar{\alpha}f(n)$ with $\bar{\alpha} := \alpha + \epsilon/f(n) \in [\alpha, \bar{\alpha})$, we deduce from (4.1) that

$$gn/s = (\beta/\bar{\alpha})^{1/2}f(n) \leq (\beta/\alpha)^{1/2}f(n). \quad (4.5)$$

In particular, $4gn/s \leq \delta f(n)$, since $\delta := 4(\beta/\alpha)^{1/2}$. Next, (4.1) implies

$$ne^{-2g^2/s} \leq n^{1-2\beta} = f(n)n^{1/3-2\beta} \ln^{-1/3} n. \quad (4.6)$$

Using the monotonicity of f and ϕ , increase \bar{n} if necessary to get for all $n \geq \bar{n}$

$$2\bar{\alpha}\phi(\bar{\alpha}f(\bar{n})) + \delta\phi(\delta f(\bar{n})) + 2n^{-2\beta} + 2 \max \left\{ [\delta f(n)]^{2/3-2\beta} n^{-2/3}, n^{-2\beta} \right\} \leq 0.95. \quad (4.7)$$

By Rem. 2.2(c), there is γ such that (4.2) holds for all $n \leq \bar{n}$; increasing γ if necessary, and using the monotonicity of f and the assumption $\beta \geq 1/6$, we have for all $n \geq \bar{n}$

$$2\bar{\alpha} + 2\delta + 5n^{1/3-2\beta} \ln^{-1/3} n + 3 \max \left\{ \delta^{1-2\beta} f(n)^{-2\beta}, n^{1/3-2\beta} \ln^{-1/3} n \right\} \leq 0.05\gamma. \quad (4.8)$$

Let $n' \geq \bar{n}$. Assuming (4.2) holds for all $n \leq n'$, for induction let $n = n' + 1$.

We need to consider the following two cases in the first call of SELECT.

Left case: $k < n/2$. First, suppose the event $\mathcal{E}_l := \{x_k^* \leq v \leq x_{k_r}^*\}$ occurs. By the rules of Steps 4–6, we have $\hat{X} = L$ (from $x_k^* \leq v$), $\hat{k} = k$ and $\hat{n} := |\hat{X}| \leq k_r - 1$ (from $v \leq x_{k_r}^*$); since $k_r < k + 2gn/s + 1$ by (3.2), we get the two (equivalent) bounds

$$\hat{n} < k + 2gn/s \quad \text{and} \quad \hat{n} - \hat{k} < 2gn/s. \quad (4.9)$$

Note that if $i_v = \lceil ks/n + g \rceil$ then, by Cor. 3.3(a,b), the Boole-Benferroni inequality and the choice (4.1), the complement \mathcal{E}'_l of \mathcal{E}_l has $P[\mathcal{E}'_l] \leq 2e^{-2g^2/s} = 2n^{-2\beta}$. Second, if $i_v \neq \lceil ks/n + g \rceil$, then $n < k + gn/s$ (Cor. 3.3(e)) combined with $k < n/2$ gives $n < 2gn/s$; hence $\hat{n} - \hat{k} < \hat{n} < n < 2gn/s$ implies (4.9). Since also \mathcal{E}_l implies (4.9), we have

$$P[\mathcal{A}'_l] \leq 2n^{-2\beta} \quad \text{for} \quad \mathcal{A}_l := \left\{ \hat{n} - \hat{k} < 2gn/s \right\}. \quad (4.10)$$

Right case: $k \geq n/2$. First, suppose the event $\mathcal{E}_r := \{x_{k_l}^* \leq v \leq x_k^*\}$ occurs. By the rules of Steps 4–6, we have $\hat{X} = U$ (from $v \leq x_k^*$), $\hat{n} - \hat{k} = n - k$ and $\hat{n} := |\hat{X}| \leq n - k_l$ (from $x_{k_l}^* \leq v$); since $k_l \geq k - 2gn/s$ by (3.2), we get the two (equivalent) bounds

$$\hat{n} \leq n - k + 2gn/s \quad \text{and} \quad \hat{k} \leq 2gn/s, \quad (4.11)$$

using $\hat{n} - \hat{k} = n - k$. If $i_v = \lceil ks/n - g \rceil$ then, by Cor. 3.3(c,d), the complement \mathcal{E}'_r of \mathcal{E}_r has $P[\mathcal{E}'_r] \leq 2e^{-2g^2/s} = 2n^{-2\beta}$. Second, if $i_v \neq \lceil ks/n - g \rceil$, then $k \leq gn/s$ (Cor. 3.3(e)) combined with $k \geq n/2$ gives $n \leq 2gn/s$; hence $\hat{k} \leq \hat{n} < n \leq 2gn/s$ implies (4.11). Thus

$$P[\mathcal{A}'_r] \leq 2n^{-2\beta} \quad \text{for} \quad \mathcal{A}_r := \left\{ \hat{k} \leq 2gn/s \right\}. \quad (4.12)$$

Since $k < n - k$ if $k < n/2$, $n - k \leq k$ if $k \geq n/2$, (4.9) and (4.11) yield

$$P[\mathcal{B}'] \leq 2n^{-2\beta} \quad \text{for } \mathcal{B} := \{ \hat{n} \leq \min\{k, n - k\} + 2gn/s \}. \quad (4.13)$$

Note that $\min\{k, n - k\} \leq \lfloor n/2 \rfloor \leq n/2$; this relation will be used implicitly below.

For the recursive call of Step 7, let \hat{s} , \hat{g} and \hat{i}_v denote the quantities generated as in (4.1) and (2.1) with n and k replaced by \hat{n} and \hat{k} , let \hat{v} be the pivot found at Step 3, and let \hat{X} , \hat{n} and \hat{k} correspond to \hat{X} , \hat{n} and \hat{k} at Step 7, so that $\hat{n} := |\hat{X}| < \hat{n}$.

The cost of selecting v and \hat{v} at Step 3 may be estimated as

$$C_{siv} + C_{\hat{s}i_v} \leq 1.5s + \gamma f(s) + 1.5\hat{s} + \gamma f(\hat{s}) \leq 3s + 2\gamma f(s), \quad (4.14)$$

since f is increasing and (4.2) holds for $\hat{s} \leq s \leq n - 1 = n'$ (cf. (4.1)) from $\hat{n} < n$.

Let $c := n - s$ and $\hat{c} := \hat{n} - \hat{s}$ denote the costs of Step 4 for the two calls. Since $0 \leq \hat{c} < n$ and $E\hat{c} = E[\hat{c}|\mathcal{B}]P[\mathcal{B}] + E[\hat{c}|\mathcal{B}']P[\mathcal{B}'] \leq E[\hat{c}|\mathcal{B}] + nP[\mathcal{B}']$, by (4.13) we have

$$c + E\hat{c} \leq n - s + \min\{k, n - k\} + 2gn/s + 2n^{1-2\beta}. \quad (4.15)$$

Using (4.2) again with $\hat{n} < n$, the cost of finishing up at Step 7 is at most

$$EC_{\hat{n}\hat{k}} \leq E[1.5\hat{n} + \gamma f(\hat{n})] = 1.5E\hat{n} + \gamma Ef(\hat{n}). \quad (4.16)$$

Thus we need suitable bounds for $E\hat{n}$ and $Ef(\hat{n})$, which may be derived as follows.

To generalize (4.13) to the recursive call, consider the events

$$\hat{\mathcal{B}} := \{ \hat{n} \leq \min\{\hat{k}, \hat{n} - \hat{k}\} + 2\hat{g}\hat{n}/\hat{s} \} \quad \text{and} \quad \mathcal{C} := \{ \hat{n} \leq \lfloor \delta f(n) \rfloor \}. \quad (4.17)$$

By (4.10) and (4.12), $\hat{\mathcal{B}} \cap \mathcal{A}_l$ and $\hat{\mathcal{B}} \cap \mathcal{A}_r$ imply \mathcal{C} , since $2gn/s + 2\hat{g}\hat{n}/\hat{s} \leq \delta f(n)$ by (4.5) with $\hat{n} < n$ and $\delta := 4(\beta/\alpha)^{1/2}$. For the recursive call, proceeding as in the derivation of (4.13) with n replaced by $\hat{n} = i$, k by \hat{k} , etc., shows that, due to random sampling,

$$P[\hat{\mathcal{B}}|\mathcal{A}_l, \hat{n} = i] \leq 2i^{-2\beta} \quad \text{and} \quad P[\hat{\mathcal{B}}|\mathcal{A}_r, \hat{n} = i] \leq 2i^{-2\beta}. \quad (4.18)$$

In the left case of $k < n/2$, using $\hat{n} < n$ and $P[\mathcal{A}'_l] \leq 2n^{-2\beta}$ (cf. (4.10)), we get

$$E\hat{n} = E[\hat{n}|\mathcal{A}_l]P[\mathcal{A}_l] + E[\hat{n}|\mathcal{A}'_l]P[\mathcal{A}'_l] \leq E[\hat{n}|\mathcal{A}_l] + n2n^{-2\beta}.$$

Partitioning \mathcal{A}_l into the events $\mathcal{D}_i := \mathcal{A}_l \cap \{\hat{n} = i\}$, $i = 0: n - 1$ ($\hat{n} < n$ always), we have

$$E[\hat{n}|\mathcal{A}_l] = \sum_{i=0}^{n-1} E[\hat{n}|\mathcal{D}_i]P[\mathcal{D}_i|\mathcal{A}_l] \leq \max_{i=0:n-1} E[\hat{n}|\mathcal{D}_i],$$

where $E[\hat{n}|\mathcal{D}_i] \leq \lfloor \delta f(n) \rfloor$ if $i \leq \lfloor \delta f(n) \rfloor + 1$, because $\hat{n} < \hat{n}$ always. As for the remaining terms, $\hat{\mathcal{B}} \cap \mathcal{A}_l \subset \mathcal{C}$ implies $P[\mathcal{C}'|\mathcal{D}_i] \leq P[\hat{\mathcal{B}}|\mathcal{D}_i] \leq 2i^{-2\beta}$ by (4.18), where $\mathcal{C}' := \{\hat{n} \leq \lfloor \delta f(n) \rfloor\}$ and $\hat{n} < \hat{n} = i$ when the event \mathcal{D}_i occurs, so $E[\hat{n}|\mathcal{D}_i] \leq \lfloor \delta f(n) \rfloor + i2i^{-2\beta}$. Hence

$$\max_{i=0:n-1} E[\hat{n}|\mathcal{D}_i] \leq \lfloor \delta f(n) \rfloor + \max_{i=\lfloor \delta f(n) \rfloor + 2:n-1} 2i^{1-2\beta},$$

where the final term is omitted if $\lfloor \delta f(n) \rfloor > n - 3$; otherwise it is at most

$$2 \max \left\{ (\lfloor \delta f(n) \rfloor + 1)^{1-2\beta}, n^{1-2\beta} \right\} \leq 2 \max \left\{ \delta^{1-2\beta} f(n)^{-2\beta}, n^{1/3-2\beta} \ln^{-1/3} n \right\} f(n),$$

since $\max_{i=\lfloor \delta f(n) \rfloor+1..n} 2i^{1-2\beta}$ is bounded as above (consider $\beta \geq 1/2$, then $\beta < 1/2$ and use $\delta f(n) < \lfloor \delta f(n) \rfloor + 1$, the monotonicity of f and (4.6) for the final inequality). Collecting the preceding estimates, we obtain

$$E\bar{n} \leq \lfloor \delta f(n) \rfloor + 2n^{1-2\beta} + 2 \max \left\{ \delta^{1-2\beta} f(n)^{-2\beta}, n^{1/3-2\beta} \ln^{-1/3} n \right\} f(n). \quad (4.19)$$

Similarly, replacing \bar{n} by $f(\bar{n})$ in our derivations and using the monotonicity of f yields

$$Ef(\bar{n}) \leq f(\lfloor \delta f(n) \rfloor) + 2f(n)n^{-2\beta} + \max_{i=\lfloor \delta f(n) \rfloor+2..n-1} 2f(i)i^{-2\beta}, \quad (4.20a)$$

where the final term is omitted if $\lfloor \delta f(n) \rfloor > n - 3$; otherwise it is at most

$$2 \max \left\{ \frac{f(\lfloor \delta f(n) \rfloor + 1)}{(\lfloor \delta f(n) \rfloor + 1)^{2\beta}}, \frac{f(n)}{n^{2\beta}} \right\} \leq 2 \max \left\{ [\delta f(n)]^{2/3-2\beta} n^{-2/3}, n^{-2\beta} \right\} f(n). \quad (4.20b)$$

To see this, use the monotonicity of f and the fact that for $i \leq n$ (cf. (4.1))

$$f(i)i^{-2\beta}/f(n) = i^{2/3-2\beta} n^{-2/3} (\ln i / \ln n)^{1/3} \leq i^{2/3-2\beta} n^{-2/3}.$$

For the right case, replace \mathcal{A}_l by \mathcal{A}_r in the preceding paragraph to get (4.19)–(4.20). Add the costs (4.14), (4.15) and (4.16), using (4.19)–(4.20), to get

$$\begin{aligned} C_{nk} &\leq 3s + 2\gamma f(s) + n - s + \min\{k, n - k\} + 2gn/s + 2n^{1-2\beta} \\ &\quad + 1.5\lfloor \delta f(n) \rfloor + 3n^{1-2\beta} + 3 \max \left\{ \delta^{1-2\beta} f(n)^{-2\beta}, n^{1/3-2\beta} \ln^{-1/3} n \right\} f(n) \\ &\quad + \gamma f(\lfloor \delta f(n) \rfloor) + 2\gamma f(n)n^{-2\beta} + 2\gamma \max \left\{ [\delta f(n)]^{2/3-2\beta} n^{-2/3}, n^{-2\beta} \right\} f(n). \end{aligned}$$

Now, using the bounds (4.3)–(4.4), $2gn/s \leq \frac{1}{2}\delta f(n)$ (cf. (4.5)) and (4.6) gives

$$\begin{aligned} C_{nk} &\leq n + \min\{k, n - k\} \\ &\quad + \left[2\bar{\alpha} + 2\delta + 5n^{1/3-2\beta} \ln^{-1/3} n + 3 \max \left\{ \delta^{1-2\beta} f(n)^{-2\beta}, n^{1/3-2\beta} \ln^{-1/3} n \right\} \right] f(n) \\ &\quad + \left[2\bar{\alpha}\phi(\bar{\alpha}f(\bar{n})) + \delta\phi(\delta f(\bar{n})) + 2n^{-2\beta} + 2 \max \left\{ [\delta f(n)]^{2/3-2\beta} n^{-2/3}, n^{-2\beta} \right\} \right] \gamma f(n). \end{aligned}$$

By (4.7)–(4.8), the two bracketed terms above are at most $0.05\gamma f(n)$ and $0.95\gamma f(n)$, respectively; thus (4.2) holds as required. \square

4.2 Other sampling strategies

We now indicate briefly how to adapt the proof of Thm 4.1 to several variations on (4.1); a choice similar to (4.21) below was used in [FIR75a].

Remarks 4.2. (a) Theorem 4.1 remains true for $\beta \geq 1/6$ and (4.1) replaced by

$$s := \min \left\{ \lceil \alpha n^{2/3} \rceil, n-1 \right\}, \quad g := (\beta s \ln n)^{1/2} \text{ and } f(n) := n^{2/3} \ln^{1/2} n. \quad (4.21)$$

Indeed, using $e^{3/2} - 1 \leq \alpha \bar{n}^{2/3} \leq \bar{n} - 1$, $e^{3/2} \leq \delta f(\bar{n})$, $\bar{\alpha} := \alpha + \bar{n}^{-2/3}$ and $s = \bar{\alpha} \bar{n}^{2/3}$ with $\bar{\alpha} \in [\alpha, \bar{\alpha}]$ yields (4.3)–(4.5) as before, and $\ln^{-1/2}$ replaces $\ln^{-1/3}$ in (4.6), (4.8) and (4.19).

(b) Theorem 4.1 holds for the following modification of (4.1) with $\epsilon_t > 1$

$$s := \min \{ \lceil \alpha f(n) \rceil, n-1 \} \text{ and } g := (\beta s \ln^\epsilon n)^{1/2} \text{ with } f(n) := n^{2/3} \ln^{\epsilon/3} n. \quad (4.22)$$

First, using $e^{\epsilon_t} - 1 \leq \alpha f(\bar{n}) \leq \bar{n} - 1$ and $e^{\epsilon_t} \leq \delta f(\bar{n})$ gives (4.3)–(4.5) as before. Next, fix $\tilde{\beta} \geq 1/6$. Let $\beta_n := \beta \ln^{\epsilon_t-1} n$. Increase \bar{n} if necessary so that $\beta_i \geq \tilde{\beta}$ for all $i \leq \min\{\bar{n}, \lceil \delta f(\bar{n}) \rceil\}$; then replace β by $\tilde{\beta}$ and $\ln^{-1/3}$ by $\ln^{-\epsilon_t/3}$ in (4.6) and below.

(c) Several other replacements for (4.1) may be analyzed as in [Kiw02, §§4.1–4.2].

(d) None of these choices gives $f(n)$ better than that in (4.1) for the bound (4.2).

We now comment briefly on the possible use of sampling with replacement.

Remarks 4.3. (a) Suppose Step 2 of SELECT employs sampling with replacement. Since the tail bound (3.1) remains valid for the binomial distribution [Chv79, Hoe63], Lemma 3.2 is not affected. However, when Step 4 no longer skips comparisons with the elements of S , $-s$ in (4.15) is replaced by 0; the resulting change in the bound on C_{nk} only needs replacing $2\bar{\alpha}$ in (4.8) by $3\bar{\alpha}$. Hence the preceding results remain valid.

(b) Of course, sampling with replacement needs additional storage for S . However, the increase in both storage and the number of comparisons may be tolerated because the sample sizes are relatively small.

4.3 Handling small subfiles

Since the sampling efficiency decreases when X shrinks, consider the following modification. For a fixed cut-off parameter $n_{\text{cut}} \geq 1$, let $\text{sSelect}(X, k)$ be a “small-select” routine that finds the k th smallest element of X in at most $C_{\text{cut}} < \infty$ comparisons when $|X| \leq n_{\text{cut}}$ (even bubble sort will do). Then SELECT is modified to start with the following

Step 0 (*Small file case*). If $n := |X| \leq n_{\text{cut}}$, return $\text{sSelect}(X, k)$.

Our preceding results remain valid for this modification. In fact it suffices if C_{cut} bounds the *expected* number of comparisons of $\text{sSelect}(X, k)$ for $n \leq n_{\text{cut}}$. For instance, (4.2) holds for $n \leq n_{\text{cut}}$ and $\gamma \geq C_{\text{cut}}$, and by induction as in Rem. 2.2(c) we have $C_{nk} < \infty$ for all n , which suffices for the proof of Thm 4.1.

Another advantage is that even small n_{cut} (1000 say) limits nicely the stack space for recursion. Specifically, the tail recursion of Step 7 is easily eliminated (set $X := \bar{X}$, $k := \hat{k}$ and go to Step 0), and the calls of Step 3 deal with subsets whose sizes quickly reach n_{cut} . For example, for the choice of (4.1) with $\alpha = 1$ and $n_{\text{cut}} = 600$, at most four recursive levels occur for $n \leq 2^{31} \approx 2.15 \cdot 10^9$.

5 A modified version

We now consider a modification inspired by a remark of [Bro76]. For k close to $\lceil n/2 \rceil$, by symmetry it is best to choose v as the sample median with $i_v = \lceil s/2 \rceil$, thus attempting to get v close to x_k^* instead of $x_{\lceil k-gn/s \rceil}^*$ or $x_{\lceil k+gn/s \rceil}^*$; then more elements are eliminated. Hence we may let

$$i_v := \begin{cases} \lceil ks/n + g \rceil & \text{if } k < n/2 - gn/s, \\ \lceil s/2 \rceil & \text{if } n/2 - gn/s \leq k \leq n/2 + gn/s, \\ \lceil ks/n - g \rceil & \text{if } k > n/2 + gn/s. \end{cases} \quad (5.1)$$

Note that (5.1) coincides with (2.1) in the *left* case of $k < n/2 - gn/s$ and the *right* case of $k > n/2 + gn/s$, but the *middle* case of $n/2 - gn/s \leq k \leq n/2 + gn/s$ fixes i_v at the median position $\lceil s/2 \rceil$; in fact i_v is the median of the three values in (5.1):

$$i_v := \max \{ \min \{ \lceil ks/n + g \rceil, \lceil s/2 \rceil \}, \lceil ks/n - g \rceil \}. \quad (5.2)$$

Corollary 3.3 remains valid for the left and right cases. For the middle case, letting

$$j_l := \max \{ \lceil n/2 - gn/s \rceil, 1 \} \quad \text{and} \quad j_r := \min \{ \lceil n/2 + gn/s \rceil, n \}, \quad (5.3)$$

we obtain from Lemma 3.2 with $\kappa = 1/2$ the following complement of Corollary 3.3.

Corollary 5.1. $P[v < x_{j_l}^*] \leq e^{-2g^2/s}$ and $P[x_{j_r}^* < v] \leq e^{-2g^2/s}$ if $n/2 - gn/s \leq k \leq n/2 + gn/s$.

Theorem 5.2. *Theorem 4.1 holds for SELECT with Step 3 using (5.1).*

Proof. We only indicate how to adapt the proof of Thm 4.1 following (4.8). As noted after (5.1), the left case now has $k < n/2 - gn/s$ and the right case has $k > n/2 + gn/s$, so we only need to discuss the middle case.

Middle case: $n/2 - gn/s \leq k \leq n/2 + gn/s$. Suppose the event $\mathcal{E}_m := \{x_{j_l}^* \leq v \leq x_{j_r}^*\}$ occurs (note that $P[\mathcal{E}'_m] \leq 2e^{-2g^2/s} = 2n^{-2\beta}$ by Cor. 5.1). If $\hat{X} = L$ then, by the rules of Steps 4–6, we have $\hat{k} = k$ and $\hat{n} \leq j_r - 1$; since $j_r < n/2 + gn/s + 1$ by (5.3), we get $\hat{n} < n/2 + gn/s$. Hence $k \geq n/2 - gn/s$ yields $\hat{n} < k + 2gn/s$ and $\hat{n} - \hat{k} < 2gn/s$ as in (4.9). Next, if $\hat{X} = U$ then $\hat{n} - \hat{k} = n - k$ and $\hat{k} := k - |L \cup E|$, so $L \cup E = \{x \in X : x \leq v\} \ni x_{j_l}^*$ gives $\hat{k} \leq k - j_l$. Since $k \leq n/2 + gn/s$ and $j_l \geq n/2 - gn/s$ by (5.3), we get $\hat{k} \leq 2gn/s$ and $\hat{n} \leq \hat{n} - \hat{k} + 2gn/s$ as in (4.11); further, $\hat{n} \leq n - j_l$ yields $\hat{n} \leq n/2 + gn/s$. Noticing that $n/2 - gn/s \leq k \leq n/2 + gn/s$ implies $n/2 \leq \min\{k, n - k\} + gn/s$, we have $\hat{n} \leq \min\{k, n - k\} + 2gn/s$ in both cases.

Thus in the middle case we again have (4.13) and hence (4.15); further, by (4.10) and (4.12), the event $\mathcal{E}_m \subset \mathcal{A}_l \cup \mathcal{A}_r$ is partitioned into $\mathcal{E}_m \cap \mathcal{A}_l$ and $\mathcal{E}_m \cap \mathcal{A}'_l \cap \mathcal{A}_r$.

Next, reasoning as before, we see that (4.18) and hence (4.19)–(4.20) remain valid in the left and right cases, whereas in the middle case we have

$$P[\hat{\mathcal{B}}' | \mathcal{E}_m, \mathcal{A}_l, \hat{n} = i] \leq 2i^{-2\beta} \quad \text{and} \quad P[\hat{\mathcal{B}}' | \mathcal{E}_m, \mathcal{A}'_l, \mathcal{A}_r, \hat{n} = i] \leq 2i^{-2\beta}. \quad (5.4)$$

In the middle case, $E\tilde{n} = E[\tilde{n}|\mathcal{E}_m]P[\mathcal{E}_m] + E[\tilde{n}|\mathcal{E}'_m]P[\mathcal{E}'_m]$ is bounded by $E[\tilde{n}|\mathcal{E}_m] + 2n^{1-2\beta}$, since $P[\mathcal{E}'_m] \leq 2n^{-2\beta}$ and $\tilde{n} < n$ always. Next, partitioning \mathcal{E}_m into $\mathcal{E}_m \cap \mathcal{A}_l$ and $\mathcal{E}_m \cap \mathcal{A}'_l \cap \mathcal{A}_r$, we obtain $E[\tilde{n}|\mathcal{E}_m] \leq \max\{E[\tilde{n}|\mathcal{E}_m, \mathcal{A}_l], E[\tilde{n}|\mathcal{E}_m, \mathcal{A}'_l, \mathcal{A}_r]\}$, where $E[\tilde{n}|\mathcal{E}_m, \mathcal{A}_l]$ and $E[\tilde{n}|\mathcal{E}_m, \mathcal{A}'_l, \mathcal{A}_r]$ may be bounded like $E[\tilde{n}|\mathcal{A}_l]$ and $E[\tilde{n}|\mathcal{A}_r]$ in the left and right cases to get (4.19). Then (4.20) is obtained similarly, and the conclusion follows as before. \square

6 Ternary and binary partitioning

In this section we discuss ways of implementing `SELECT` when the input set is given as an array $x[1:n]$. In particular, we recall the ternary partitioning scheme of [BeM93, BeS97] and its modification of [Kiw02] that obviates subscript range checking.

The following notation is needed to describe the operations of `SELECT` in more detail.

Each stage works with a segment $x[l:r]$ of the input array $x[1:n]$, where $1 \leq l \leq r \leq n$ are such that $x_i < x_l$ for $i = 1:l-1$, $x_r < x_i$ for $i = r+1:n$, and the k th smallest element of $x[1:n]$ is the $(k-l+1)$ th smallest element of $x[l:r]$. The task of `SELECT` is *extended*: given $x[l:r]$ and $l \leq k \leq r$, `SELECT`(x, l, r, k, k_-, k_+) permutes $x[l:r]$ and finds $l \leq k_- \leq k \leq k_+ \leq r$ such that $x_i < x_k$ for all $l \leq i < k_-$, $x_i = x_k$ for all $k_- \leq i \leq k_+$, $x_i > x_k$ for all $k_+ < i \leq r$. The initial call is `SELECT`($x, 1, n, k, k_-, k_+$).

A vector swap denoted by $x[a:b] \leftrightarrow x[b+1:c]$ means that the first $d := \min(b+1-a, c-b)$ elements of array $x[a:c]$ are exchanged with its last d elements in arbitrary order if $d > 0$; e.g., we may exchange $x_{a+i} \leftrightarrow x_{c-i}$ for $0 \leq i < d$, or $x_{a+i} \leftrightarrow x_{c-d+1+i}$ for $0 \leq i < d$.

6.1 Ternary partitions

For a given pivot $v := x_k$ from the array $x[l:r]$, the following *ternary* scheme partitions the array into three blocks, with $x_m < v$ for $l \leq m < a$, $x_m = v$ for $a \leq m \leq b$, $x_m > v$ for $b < m \leq r$. The basic idea is to work with the five inner parts of the array

$$\begin{array}{|c|c|c|c|c|c|c|} \hline x < v & x = v & x < v & ? & x > v & x = v & x > v \\ \hline l & \bar{l} & p & i & j & q & \bar{r} & r \\ \hline \end{array} \quad (6.1)$$

until the middle part is empty or just contains an element equal to the pivot

$$\begin{array}{|c|c|c|c|c|} \hline x = v & x < v & x = v & x > v & x = v \\ \hline \bar{l} & p & j & i & q & \bar{r} \\ \hline \end{array} \quad (6.2)$$

(i.e., $j = i - 1$ or $j = i - 2$), then swap the ends into the middle for the final arrangement

$$\begin{array}{|c|c|c|} \hline x < v & x = v & x > v \\ \hline \bar{l} & a & b & \bar{r} \\ \hline \end{array} \quad (6.3)$$

- A1. [Initialize.] Set $v := x_k$ and exchange $x_l \leftrightarrow x_k$. Set $i := \bar{l} := l$, $p := l + 1$, $j := \bar{r} := r$, $q := r - 1$. If $v < x_r$, set $\bar{r} := r - 1$. If $v > x_r$, exchange $x_l \leftrightarrow x_r$ and set $\bar{l} := l + 1$.
- A2. [Increase i until $x_i \geq v$.] Increase i by 1; then if $x_i < v$, repeat this step.

- A3.** [Decrease j until $x_j \leq v$.] Decrease j by 1; then if $x_j > v$, repeat this step.
- A4.** [Exchange.] (Here $x_j \leq v \leq x_i$.) If $i < j$, exchange $x_i \leftrightarrow x_j$; then if $x_i = v$, exchange $x_i \leftrightarrow x_p$ and increase p by 1; if $x_j = v$, exchange $x_j \leftrightarrow x_q$ and decrease q by 1; return to A2. If $i = j$ (so that $x_i = x_j = v$), increase i by 1 and decrease j by 1.
- A5.** [Cleanup.] Exchange $x[\bar{l}:p-1] \leftrightarrow x[p:j]$ and $x[i:q] \leftrightarrow x[q+1:\bar{r}]$. Finally, set $a := \bar{l} + j - p + 1$ and $b := \bar{r} - q + i - 1$.

Step A1 ensures that $x_l \leq v \leq x_r$, so steps A2 and A3 don't need to test whether $i \leq j$. However, this scheme involves two extraneous comparisons (only one when $i = j$ at A4). Consider, therefore, the following scheme of [BeM93, BeS97] and [Knu97, Ex. 5.2.2–41], also based on the arrangements (6.1)–(6.3) with $\bar{l} := l$, $\bar{r} := r$ and final $j = i - 1$.

- B1.** [Initialize.] Set $v := x_k$ and exchange $x_l \leftrightarrow x_k$. Set $i := p := l + 1$, $\bar{l} := l$ and $j := q := \bar{r} := r$.
- B2.** [Increase i until $x_i > v$.] If $i \leq j$ and $x_i < v$, increase i by 1 and repeat this step. If $i \leq j$ and $x_i = v$, exchange $x_p \leftrightarrow x_i$, increase p and i by 1, and repeat this step.
- B3.** [Decrease j until $x_j < v$.] If $i \leq j$ and $x_j > v$, decrease j by 1 and repeat this step. If $i \leq j$ and $x_j = v$, exchange $x_j \leftrightarrow x_q$, decrease j and q by 1, and repeat this step.
- B4.** [Exchange.] If $i \leq j$, exchange $x_i \leftrightarrow x_j$, increase i by 1, decrease j by 1, and return to B2.
- B5.** [Cleanup.] Swap $x[\bar{l}:p-1] \leftrightarrow x[p:j]$ and $x[i:q] \leftrightarrow x[q+1:\bar{r}]$. Finally, set $a := \bar{l} + i - p$ and $b := \bar{r} - q + j$.

Relative to scheme A, scheme B saves one or two v -comparisons at the expense of $r - l + 2$ comparisons of i vs. j . Since for usual choices of n_{cut} , we have $r - l \gg 2$ in SELECT and relatively few small partitions in sSelect, scheme A is faster than B unless the cost of key comparisons is extremely large.

6.2 Preparing for ternary partitions

At Step 1, $r - l + 1$ replaces n in finding s and g . At Step 2, it is convenient to place the sample in the initial part of $x[l:r]$ by exchanging $x_i \leftrightarrow x_{i+\text{rand}(r-i)}$ for $l \leq i \leq r_s := l + s - 1$, where $\text{rand}(r - i)$ denotes a random integer, uniformly distributed between 0 and $r - i$.

Step 3 uses $i := k - l + 1$ and $m := r - l + 1$ instead of k and n to find the pivot position

$$k_v := \begin{cases} \min \{ [l - 1 + is/m + g], r_s \} & \text{if } i < m/2, \\ \max \{ [l - 1 + is/m - g], l \} & \text{if } i \geq m/2, \end{cases} \quad (6.4)$$

so that the recursive call of SELECT($x, l, r_s, k_v, k_v^-, k_v^+$) produces $v := x_{k_v}$.

After v has been found, our array looks as follows

$$\begin{array}{|c|c|c|c|} \hline x < v & x = v & x > v & ? \\ \hline l & k_v^- & k_v^+ & r_s \quad r \\ \hline \end{array} \quad (6.5)$$

Setting $\bar{l} := k_v^-$, $p := k_v^+ + 1$, $\bar{r} := r - r_s + k_v^+$ and $q := \bar{r}$, we swap $x[p:r_s] \leftrightarrow x[r_s + 1:r]$; if $k_v^+ = r_s$, we exchange $x_{k_v^+} \leftrightarrow x_q$ and decrease p and q by 1. This yields the arrangement of (6.1), so we may use either scheme A of §6.1 with initial $i := p - 1$, $j := q + 1$ and step A1 omitted, or scheme B with $i := p$, $j := q$ and step B1 omitted.

After partitioning l and r are updated by setting $l := b + 1$ if $a \leq k$, $r := a - 1$ if $k \leq b$. If $l \geq r$, SELECT may return $k_- := k_+$ if $l = r$, $k_- := r + 1$ and $k_+ := l - 1$ if $l > r$. Otherwise, instead of calling SELECT recursively, Step 6 may jump back to Step 1, or to Step 0 if sSelect is used (cf. §4.3).

A simple version of sSelect is obtained if Steps 2 and 3 choose $v := x_k$ when $r - l + 1 \leq n_{\text{cut}}$ (this choice of [FIR75a] works well in practice, but more sophisticated pivots could be tried); then the ternary partitioning code can be used by sSelect as well.

6.3 Binary partitions

We now consider a *binary* version of SELECT, called BSELECT, which employs less refined but potentially faster partitioning. This version works with $x[l:r]$ such that $x_i \leq x_l$ for $i = 1:l - 1$, $x_r \leq x_i$ for $i = r + 1:n$, and its task is standard: given $x[l:r]$ and $l \leq k \leq r$, BSELECT(x, l, r, k) permutes $x[l:r]$ so that $x_i \leq x_k$ for all $l \leq i < k$, and $x_k \leq x_i$ for all $k < i \leq r$; the initial call is BSELECT($x, 1, n, k$).

For a given pivot $v := x_k$ from the array $x[l:r]$, the following *binary* scheme partitions the array into three blocks, with $x_m \leq v$ for $l \leq m < a$, $x_m = v$ for $a \leq m \leq b$, $v \leq x_m$ for $b < m \leq r$; usually $a = b$ and the middle block is singleton.

- C1. [Initialize.] Set $v := x_k$ and exchange $x_l \leftrightarrow x_k$. Set $i := \hat{p} := l$ and $j := r$. If $v > x_r$, exchange $x_{\hat{p}} \leftrightarrow x_r$ and set $\hat{p} := r$. (Thus $x_i \leq v = x_{\hat{p}} \leq x_r$ always.)
- C2. [Increase i until $x_i \geq v$.] Increase i by 1; then if $x_i < v$, repeat this step.
- C3. [Decrease j until $x_j \leq v$.] Decrease j by 1; then if $x_j > v$, repeat this step.
- C4. [Exchange.] (Here $x_j \leq v \leq x_i$.) If $i < j$, exchange $x_i \leftrightarrow x_j$ and return to C2. If $i = j$ (so that $x_i = x_j = v$), increase i by 1 and decrease j by 1.
- C5. [Cleanup.] If $\hat{p} \neq r$, exchange $x_{\hat{p}} \leftrightarrow x_j$, set $a := j$ and $b := i - 1$; otherwise exchange $x_i \leftrightarrow x_{\hat{p}}$, set $a := j + 1$ and $b := i$.

The setup of §6.2 changes as follows. Step 3 calls BSELECT(x, l, r_s, k_v) to find $v := x_{k_v}$; then (6.5) changes to

$$\begin{array}{cccc} \boxed{x \leq v} & \boxed{v} & \boxed{x \geq v} & \boxed{?} \\ l & k_v & r_s & r \end{array} \quad (6.6)$$

Setting $i := \hat{p} := k_v$ and $j := r - r_s + k_v$, we swap $x[i + 1:r_s] \leftrightarrow x[r_s + 1:r - 1]$. If $v > x_r$, we exchange $x_{\hat{p}} \leftrightarrow x_r$ and set $\hat{p} := r$. Then we may use scheme C (with C1 omitted) for updating l and r , also in sSelect as in §6.2.

The inner loops of schemes A and C (i.e., A2, A3, C2, C3) coincide, whereas C4 is like A4 without its equality tests and associated updates. When equal elements are absent, scheme C (although not equivalent to A) yields correct partitions for Step 4. When equal

elements occur, its partitions needn't meet the requirements of Step 4, but are still usable, because $r-l$ shrinks. In effect, BSELECT works like SELECT in the case of distinct elements, but may require more comparisons otherwise.

7 Experimental results

7.1 Implemented algorithms

An implementation of SELECT was programmed in Fortran 77 and run on a notebook PC (Pentium II 400 MHz, 256 MB RAM) under MS Windows 98. The input set X was specified as a double precision array. For efficiency, the recursion was removed and small arrays with $n \leq n_{\text{cut}}$ were handled as if Steps 2 and 3 chose $v := x_k$; the resulting version of sSelect (cf. §§4.3 and 6.2) typically required less than $3.5n$ comparisons. The choice of (4.21) was employed, with the parameters $\alpha = 0.5$, $\beta = 0.25$ and $n_{\text{cut}} = 600$ as proposed in [FIR75a]; future work should test other sample sizes and parameters.

A similar implementation of BSELECT was programmed as described in §6.3.

7.2 Testing examples

As in [Kiw02], we used minor modifications of the input sequences of [Val00]:

random A random permutation of the integers 1 through n .

onezero A random permutation of $\lceil n/2 \rceil$ ones and $\lfloor n/2 \rfloor$ zeroes.

sorted The integers 1 through n in increasing order.

rotated A sorted sequence rotated left once; i.e., $(2, 3, \dots, n, 1)$.

organpipe The integers 1 through $n/2$ in increasing order, followed by $n/2$ through 1 in decreasing order.

m3killer Musser's "median-of-3 killer" sequence with $n = 4j$ and $k = n/2$:

$$\begin{pmatrix} 1 & 2 & 3 & 4 & \dots & k-2 & k-1 & k & k+1 & \dots & 2k-2 & 2k-1 & 2k \\ 1 & k+1 & 3 & k+3 & \dots & 2k-3 & k-1 & 2 & 4 & \dots & 2k-2 & 2k-1 & 2k \end{pmatrix}.$$

twofaced Obtained by randomly permuting the elements of an m3killer sequence in positions $4\lceil \log_2 n \rceil$ through $n/2 - 1$ and $n/2 + 4\lceil \log_2 n \rceil - 1$ through $n - 2$.

For each input sequence, its (lower) median element was found.

Table 7.1: Performance of SELECT on randomly generated inputs.

Sequence	Size n	Time [sec]			Comparisons [n]			γ_{avg}	L_{avg} [n]	P_{avg} [$\ln n$]	N_{avg} [$\ln n$]	p_{avg}	s_{avg} [% n]
		avg	max	min	avg	max	min						
random	50K	0.01	0.06	0.01	1.67	1.76	1.60	1.88	1.67	0.45	0.55	8.03	2.63
	100K	0.02	0.06	0.01	1.62	1.67	1.58	1.68	1.62	0.61	0.69	7.81	2.12
	500K	0.06	0.11	0.05	1.56	1.63	1.53	1.41	1.56	0.67	0.75	8.45	1.19
	1M	0.13	0.17	0.11	1.55	1.59	1.53	1.23	1.55	0.70	0.77	8.38	0.92
	2M	0.23	0.28	0.22	1.54	1.57	1.52	1.19	1.54	0.77	0.84	8.48	0.72
	4M	0.47	0.50	0.44	1.53	1.55	1.52	1.13	1.53	0.86	0.92	8.35	0.57
	8M	0.91	0.99	0.82	1.52	1.54	1.51	1.05	1.52	0.88	0.94	8.27	0.44
	16M	1.82	1.93	1.70	1.52	1.53	1.51	1.11	1.52	0.94	1.00	8.33	0.35
onezero	50K	0.01	0.05	0.01	1.33	1.51	1.00	0.00	1.33	0.24	0.18	1.23	1.99
	100K	0.01	0.06	0.01	1.10	1.51	1.00	0.00	1.10	0.21	0.16	1.22	1.27
	500K	0.08	0.11	0.05	1.25	1.50	1.00	0.00	1.25	0.26	0.14	1.24	0.86
	1M	0.16	0.22	0.11	1.20	1.50	1.00	0.00	1.20	0.23	0.13	1.23	0.65
	2M	0.33	0.39	0.22	1.30	1.50	1.00	0.00	1.30	0.26	0.14	1.19	0.56
	4M	0.62	0.72	0.49	1.25	1.50	1.00	0.00	1.25	0.22	0.12	1.24	0.43
	8M	1.24	1.48	0.99	1.25	1.50	1.00	0.00	1.25	0.22	0.12	1.28	0.34
	16M	2.62	2.97	2.03	1.33	1.50	1.00	0.00	1.33	0.22	0.12	1.12	0.29
twofaced	50K	0.01	0.06	0.01	1.67	1.79	1.60	1.93	1.67	0.49	0.58	8.21	2.66
	100K	0.02	0.06	0.01	1.62	1.67	1.56	1.58	1.62	0.59	0.68	7.15	2.10
	500K	0.06	0.11	0.05	1.57	1.60	1.54	1.51	1.57	0.67	0.75	8.52	1.20
	1M	0.12	0.17	0.11	1.55	1.59	1.53	1.36	1.55	0.72	0.80	8.55	0.93
	2M	0.24	0.28	0.21	1.53	1.55	1.52	1.04	1.53	0.72	0.79	8.52	0.72
	4M	0.45	0.50	0.44	1.53	1.55	1.51	1.06	1.53	0.82	0.88	8.29	0.56
	8M	0.92	0.99	0.87	1.52	1.54	1.51	1.03	1.52	0.86	0.92	8.50	0.44
	16M	1.81	1.93	1.71	1.52	1.53	1.51	1.08	1.52	0.95	1.00	8.43	0.35

7.3 Computational results

We varied the input size n from 50,000 to 16,000,000. For the random, onezero and twofaced sequences, for each input size, 20 instances were randomly generated; for the deterministic sequences, five runs were made to measure the solution time.

The performance of SELECT on randomly generated inputs is summarized in Table 7.1, where the average, maximum and minimum solution times are in seconds, and the comparison counts are in multiples of n ; e.g., column six gives C_{avg}/n , where C_{avg} is the average number of comparisons made over all instances. Thus $\gamma_{\text{avg}} := (C_{\text{avg}} - 1.5n)_+/f(n)$ estimates the constant γ in the bound (4.2); moreover, we have $C_{\text{avg}} \approx L_{\text{avg}}$, where L_{avg} is the average sum of sizes of partitioned arrays. Further, P_{avg} is the average number of SELECT partitions, whereas N_{avg} is the average number of calls to sSelect and p_{avg} is the average number of sSelect partitions per call; both P_{avg} and N_{avg} grow slowly with $\ln n$ (linearly on the onezero inputs). Finally, s_{avg} is the average sum of sample sizes; $s_{\text{avg}}/n^{2/3}$ drops from 0.97 for $n = 50K$ to 0.88 for $n = 16M$ on the random and twofaced inputs, and oscillates about 0.7 on the onezero inputs, whereas the initial $s/n^{2/3} \approx \alpha = 0.5$. The results for the random and twofaced sequences are very similar: the average solution times grow linearly with n (except for small inputs whose solution times couldn't be

Table 7.2: Performance of SELECT on deterministic inputs.

Sequence	Size n	Time [sec]			Comparisons $[n]$	γ_{avg}	L_{avg} $[n]$	P_{avg} $[\ln n]$	N_{avg} $[\ln n]$	p_{avg}	s_{avg} $[\%n]$
		avg	max	min							
sorted	50K	0.02	0.06	0.01	1.76	2.92	1.76	0.55	0.65	6.43	2.88
	100K	0.02	0.06	0.01	1.73	3.11	1.73	0.69	0.78	6.44	2.33
	500K	0.06	0.06	0.06	1.56	1.35	1.56	0.69	0.76	8.50	1.19
	1M	0.09	0.11	0.05	1.56	1.51	1.56	0.65	0.72	7.00	0.94
	2M	0.17	0.17	0.16	1.56	1.90	1.56	0.90	0.96	8.21	0.75
	4M	0.32	0.33	0.27	1.54	1.73	1.54	0.79	0.86	9.85	0.58
	8M	0.62	0.66	0.60	1.53	1.64	1.53	1.07	1.13	7.56	0.45
	16M	1.21	1.21	1.20	1.53	1.71	1.53	1.02	1.09	6.61	0.36
rotated	50K	0.01	0.05	0.01	1.78	3.16	1.78	0.55	0.65	7.29	2.87
	100K	0.02	0.06	0.01	1.73	3.13	1.73	0.61	0.69	7.88	2.30
	500K	0.04	0.05	0.01	1.56	1.33	1.56	0.61	0.69	9.78	1.18
	1M	0.09	0.11	0.06	1.56	1.52	1.56	0.72	0.80	7.36	0.94
	2M	0.16	0.17	0.16	1.56	1.92	1.56	0.90	0.96	8.57	0.75
	4M	0.32	0.33	0.27	1.54	1.71	1.54	0.86	0.92	8.93	0.58
	8M	0.61	0.66	0.60	1.53	1.63	1.53	1.01	1.07	8.00	0.45
	16M	1.21	1.21	1.21	1.53	1.72	1.53	1.02	1.09	8.22	0.36
organpipe	50K	0.01	0.05	0.01	1.66	1.74	1.66	0.46	0.55	8.33	2.59
	100K	0.01	0.06	0.01	1.59	1.24	1.59	0.52	0.61	7.57	2.02
	500K	0.05	0.06	0.01	1.54	0.89	1.54	0.61	0.69	8.78	1.15
	1M	0.11	0.11	0.11	1.53	0.68	1.53	0.51	0.58	9.13	0.89
	2M	0.20	0.22	0.16	1.54	1.33	1.54	0.76	0.83	8.00	0.73
	4M	0.40	0.44	0.38	1.53	1.12	1.53	0.92	0.99	7.40	0.57
	8M	0.77	0.77	0.76	1.53	1.41	1.53	1.01	1.07	7.53	0.45
	16M	1.53	1.54	1.48	1.52	1.38	1.52	1.09	1.15	7.42	0.35
m3killer	50K	0.01	0.06	0.01	1.65	1.67	1.65	0.46	0.55	8.83	2.55
	100K	0.02	0.05	0.01	1.59	1.17	1.59	0.61	0.69	7.75	2.05
	500K	0.05	0.06	0.05	1.54	0.92	1.54	0.61	0.69	6.33	1.16
	1M	0.11	0.11	0.11	1.53	0.83	1.53	0.58	0.65	9.89	0.90
	2M	0.22	0.22	0.22	1.54	1.33	1.54	0.76	0.83	8.58	0.73
	4M	0.43	0.44	0.38	1.53	1.19	1.53	0.92	0.99	8.87	0.57
	8M	0.77	0.77	0.77	1.54	1.85	1.54	1.13	1.20	8.37	0.46
	16M	1.48	1.49	1.48	1.52	1.08	1.52	1.09	1.15	8.16	0.35

measured accurately), and the differences between maximum and minimum times are quite small (and also partly due to the operating system). Except for the smallest inputs, the maximum and minimum numbers of comparisons are quite close, and C_{avg} nicely approaches the theoretical lower bound of $1.5n$; this is reflected in the values of γ_{avg} . The results for the onezero inputs essentially average two cases: the first pass eliminates either almost all or about half of the elements.

Table 7.2 exhibits similar features of SELECT on the deterministic inputs. The results for the sorted and rotated sequences are very similar, whereas the solution times on the organpipe and m3killer sequences are between those for the sorted and random sequences.

The results of Tabs. 7.1–7.2 were obtained with scheme A of §6.1; to save space, Table 7.3 gives only selected results for scheme B. Scheme B was slower than scheme A by about

Table 7.3: Performance of SELECT with ternary scheme B.

Sequence	Size n	Time [sec]			Comparisons [n]			γ_{avg}	L_{avg} [n]	P_{avg} [$\ln n$]	N_{avg} [$\ln n$]	p_{avg}	s_{avg} [% n]
		avg	max	min	avg	max	min						
random	2M	0.28	0.33	0.27	1.54	1.57	1.52	1.18	1.54	0.77	0.83	8.18	0.72
	4M	0.55	0.61	0.49	1.53	1.55	1.51	1.14	1.53	0.87	0.93	8.12	0.57
	8M	1.09	1.16	1.04	1.52	1.54	1.51	1.05	1.52	0.89	0.95	8.14	0.44
	16M	2.16	2.26	2.08	1.52	1.53	1.51	1.09	1.52	0.94	1.00	8.29	0.35
onezero	2M	0.49	0.55	0.38	1.30	1.50	1.00	0.00	1.30	0.26	0.14	1.19	0.56
	4M	0.93	1.10	0.77	1.25	1.50	1.00	0.00	1.25	0.22	0.12	1.24	0.43
	8M	1.87	2.20	1.53	1.25	1.50	1.00	0.00	1.25	0.22	0.12	1.28	0.34
	16M	3.92	4.40	3.02	1.33	1.50	1.00	0.00	1.33	0.22	0.12	1.12	0.29
sorted	2M	0.22	0.22	0.22	1.56	1.56	1.56	1.90	1.56	0.90	0.96	8.43	0.75
	4M	0.46	0.50	0.44	1.54	1.54	1.54	1.71	1.54	0.79	0.79	9.42	0.58
	8M	0.88	0.88	0.87	1.53	1.53	1.53	1.62	1.53	1.13	1.20	6.79	0.45
	16M	1.73	1.76	1.70	1.53	1.53	1.53	1.71	1.53	1.02	1.09	6.61	0.36

20% on the random, twofaced, organpipe and m3killer inputs, 40% on the sorted and rotated ones, and 50% on the onezero inputs.

The performance of BSELECT with the binary scheme C on the same inputs is given in Table 7.4. SELECT with scheme A was slower than BSELECT by about 5% on the random, twofaced, organpipe and m3killer inputs. Both were equivalent on the sorted and rotated inputs. However, on the onezero inputs, BSELECT was slower by about 40% and made excessively many comparisons and partitions.

The preceding results were obtained with the modified choice (5.1) of i_v . For brevity, Table 7.5 gives results for SELECT with scheme A and the standard choice (2.1) of i_v on the random inputs only, since these inputs are most frequently used in theory and practice for evaluating sorting and selection methods. The modified choice typically requires fewer comparisons for small inputs, but its advantages are less pronounced for larger inputs. A similar behavior was observed for SELECT with scheme B and for BSELECT.

For comparison, Table 7.6 extracts from [Kiw02] some results of QSELECT for the samples (4.1). As noted in §1, QSELECT is slightly faster than SELECT on larger inputs because most of its work occurs on the first partition (cf. L_{avg} in Tabs. 7.1 and 7.6). In Table 7.7 we give corresponding results for rISELECT, a Fortran version of the algorithm of [Val00]. For these inputs, rISELECT behaves like FIND with median-of-three partitioning (because the average numbers of randomization steps, N_{rnd} , are negligible); hence the expected value of C_{avg} is of order $2.75n$ [KMP97].

Our final Table 7.8 shows that SELECT beats its competitors with respect to the numbers of comparisons made (expressed as multiples of n) on small random inputs (100 instances for each input size n).

Our computational results, combined with those in [Kiw02], suggest that both SELECT and QSELECT may compete with FIND in practice.

Acknowledgment. I would like to thank Olgierd Hryniewicz, Roger Koenker, Ronald L. Rivest and John D. Valois for useful discussions.

Table 7.4: Performance of BSELECT with binary scheme C.

Sequence	Size n	Time [sec]			Comparisons [n]			γ_{avg}	L_{avg} [n]	P_{avg} [$\ln n$]	N_{avg} [$\ln n$]	P_{avg}	s_{avg} [% n]
		avg	max	min	avg	max	min						
random	2M	0.23	0.27	0.22	1.53	1.57	1.52	1.15	1.53	0.75	0.82	8.49	0.72
	4M	0.45	0.50	0.38	1.53	1.55	1.52	1.15	1.53	0.85	0.92	8.32	0.57
	8M	0.87	0.94	0.82	1.52	1.54	1.51	1.05	1.52	0.88	0.94	8.23	0.44
	16M	1.75	1.86	1.64	1.52	1.53	1.51	1.09	1.52	0.94	1.00	8.39	0.35
onezero	2M	0.44	0.49	0.43	2.69	2.69	2.69	39.42	2.69	3.08	3.15	6.05	1.56
	4M	0.90	0.94	0.87	2.69	2.69	2.69	48.40	2.69	3.66	3.73	6.05	1.24
	8M	1.77	1.82	1.75	2.69	2.69	2.69	59.52	2.69	4.20	4.26	6.13	0.98
	16M	3.53	3.58	3.46	2.68	2.69	2.68	73.29	2.68	4.76	4.82	6.18	0.78
sorted	2M	0.17	0.17	0.16	1.56	1.56	1.56	1.88	1.56	0.90	0.96	7.43	0.75
	4M	0.33	0.33	0.33	1.54	1.54	1.54	1.72	1.54	0.79	0.86	9.31	0.58
	8M	0.62	0.66	0.61	1.53	1.53	1.53	1.63	1.53	1.01	1.07	7.71	0.45
	16M	1.21	1.21	1.20	1.53	1.53	1.53	1.70	1.53	1.02	1.09	7.50	0.36

Table 7.5: Performance of SELECT with the standard choice of i_v .

Sequence	Size n	Time [sec]			Comparisons [n]			γ_{avg}	L_{avg} [n]	P_{avg} [$\ln n$]	N_{avg} [$\ln n$]	P_{avg}	s_{avg} [% n]
		avg	max	min	avg	max	min						
random	50K	0.01	0.06	0.01	1.82	1.94	1.67	3.62	1.82	0.59	0.68	8.36	2.98
	100K	0.02	0.06	0.01	1.77	1.86	1.71	3.71	1.77	0.76	0.85	7.57	2.40
	500K	0.07	0.11	0.05	1.65	1.69	1.61	3.18	1.65	0.80	0.87	8.62	1.29
	1M	0.14	0.17	0.11	1.60	1.63	1.55	2.80	1.60	0.89	0.96	8.32	0.99
	2M	0.26	0.28	0.22	1.58	1.61	1.54	2.76	1.58	0.96	1.03	8.46	0.77
	4M	0.52	0.55	0.49	1.57	1.59	1.54	2.77	1.57	1.16	1.23	8.13	0.60
	8M	0.99	1.04	0.98	1.55	1.57	1.53	2.60	1.55	1.18	1.24	8.29	0.47
	16M	1.94	1.98	1.92	1.54	1.55	1.53	2.54	1.54	1.19	1.25	8.72	0.36

Table 7.6: Performance of quintary QSELECT on random inputs.

Sequence	Size n	Time [sec]			Comparisons [n]			γ_{avg}	L_{avg} [n]	P_{avg} [$\ln n$]	N_{avg} [$\ln n$]	P_{avg}	s_{avg} [% n]
		avg	max	min	avg	max	min						
random	50K	0.01	0.06	0.01	1.79	1.84	1.74	4.91	1.21	0.46	1.01	7.40	4.10
	100K	0.02	0.06	0.01	1.73	1.77	1.70	4.77	1.15	0.43	0.96	8.03	3.20
	500K	0.06	0.11	0.05	1.62	1.63	1.61	4.06	1.08	0.56	1.20	8.00	1.86
	1M	0.12	0.17	0.11	1.59	1.60	1.58	3.95	1.06	0.67	1.40	7.95	1.47
	2M	0.22	0.22	0.21	1.57	1.58	1.56	3.76	1.04	0.76	1.59	7.90	1.16
	4M	0.43	0.44	0.38	1.56	1.56	1.55	3.63	1.03	0.95	1.95	7.29	0.92
	8M	0.83	0.88	0.82	1.54	1.55	1.54	3.54	1.03	0.98	2.00	7.41	0.72
	16M	1.62	1.65	1.59	1.53	1.54	1.53	3.39	1.02	1.00	2.05	7.77	0.57

Table 7.7: Performance of riSELECT on random inputs.

Sequence	Size n	Time [sec]			Comparisons [n]			L_{avg} [n]	N_{rnd}
		avg	max	min	avg	max	min		
random	50K	0.01	0.06	0.01	3.10	4.32	1.88	3.10	0.40
	100K	0.03	0.06	0.01	2.61	4.20	1.77	2.61	0.25
	500K	0.10	0.11	0.05	2.90	4.23	1.69	2.90	0.20
	1M	0.18	0.22	0.11	2.81	3.64	1.84	2.81	0.35
	2M	0.34	0.44	0.22	2.60	3.57	1.83	2.60	0.30
	4M	0.77	1.38	0.44	2.88	4.81	1.83	2.88	0.55
	8M	1.38	1.70	1.05	2.60	3.48	1.80	2.60	0.45
16M	3.00	4.01	1.75	2.99	4.49	1.73	2.99	0.45	

Table 7.8: Numbers of comparisons made on small random inputs.

Size		1000	2500	5000	7500	10000	12500	15000	17500	20000	25000
SELECT	avg	2.47	2.07	1.94	1.86	1.84	1.81	1.76	1.75	1.74	1.73
	max	4.79	2.73	2.53	2.12	2.02	2.02	2.01	2.16	1.95	1.91
	min	1.51	1.72	1.64	1.63	1.64	1.62	1.07	1.59	1.04	1.61
QSELECT	avg	2.82	2.54	2.26	2.14	2.08	2.04	1.99	1.96	1.95	1.91
	max	3.90	3.37	2.61	2.45	2.31	2.20	2.12	2.13	2.12	2.07
	min	2.05	2.09	1.99	1.92	1.85	1.87	1.86	1.82	1.84	1.82
riSELECT	avg	2.72	2.84	2.66	2.71	2.72	2.82	2.78	2.75	2.75	2.84
	max	4.41	4.51	4.74	4.38	4.57	4.65	4.66	4.56	4.61	4.64
	min	1.68	1.83	1.75	1.59	1.70	1.77	1.78	1.67	1.90	1.71

References

- [BoM93] J. L. Bentley and M. D. McIlroy, *Engineering a sort function*, Software–Practice and Experience **23** (1993) 1249–1265.
- [BoS97] J. L. Bentley and R. Sedgwick, *Fast algorithms for sorting and searching strings*, in Proceedings of the 8th Annual ACM-SIAM Symposium on Discrete Algorithms (SODA'97), SIAM, Philadelphia, 1997, pp. 360–369.
- [BFP⁺72] M. R. Blum, R. W. Floyd, V. R. Pratt, R. L. Rivest and R. E. Tarjan, *Time bounds for selection*, J. Comput. System Sci. **7** (1972) 448–461.
- [Bro76] T. Brown, *Remark on algorithm 489*, ACM Trans. Math. Software **3** (1976) 301–304.
- [Chv79] V. Chvátal, *The tail of the hypergeometric distribution*, Discrete Math. **25** (1979) 285–287.
- [CuM89] W. Cunto and J. I. Munro, *Average case selection*, J. of the ACM **36** (1989) 270–279.
- [DHU201] D. Dor, J. Håstad, S. Ulfberg and U. Zwick, *On lower bounds for selecting the median*, SIAM J. Discrete Math. **14** (2001) 299–311.
- [DoZ99] D. Dor and U. Zwick, *Selecting the median*, SIAM J. Comput. **28** (1999) 1722–1758.
- [DoZ01] ———, *Median selection requires $(2 + \epsilon)N$ comparisons*, SIAM J. Discrete Math. **14** (2001) 312–325.
- [FIR75a] R. W. Floyd and R. L. Rivest, *The algorithm SELECT—for finding the i th smallest of n elements (Algorithm 489)*, Comm. ACM **18** (1975) 173.
- [FIR75b] ———, *Expected time bounds for selection*, Comm. ACM **18** (1975) 165–172.

- [Grü99] R. Grübel, *On the median-of-k version of Hoare's selection algorithm*, Theor. Inform. Appl. **33** (1999) 177–192.
- [Hoa61] C. A. R. Hoare, *FIND (Algorithm 65)*, Comm. ACM **4** (1961) 321–322.
- [Hoe63] W. Hoeffding, *Probability inequalities for sums of bounded random variables*, J. Amer. Statist. Assoc. **58** (1963) 13–30.
- [Kiw02] K. C. Kiwiel, *Randomized selection revisited*, Tech. report, Systems Research Institute, Warsaw, 2002.
- [KMP97] P. Kirschenhofer, C. Martínez and H. Prodinger, *Analysis of Hoare's Find algorithm with median-of-three partition*, Random Structures and Algorithms **10** (1997) 143–156.
- [Knu97] D. E. Knuth, *The Art of Computer Programming. Volume I: Fundamental Algorithms*, third ed., Addison-Wesley, Reading, MA, 1997.
- [Knu98] ———, *The Art of Computer Programming. Volume III: Sorting and Searching*, second ed., Addison-Wesley, Reading, MA, 1998.
- [MaR01] C. Martínez and S. Roura, *Optimal sampling strategies in quicksort and quickselect*, SIAM J. Comput. **31** (2001) 683–705.
- [Mus97] D. R. Musser, *Introspective sorting and selection algorithms*, Software–Practice and Experience **27** (1997) 983–993.
- [PRKT83] J. T. Postmus, A. H. G. Rinnooy Kan and G. T. Timmer, *An efficient dynamic selection method*, Comm. ACM **26** (1983) 878–881.
- [Sed77] R. Sedgewick, *Quicksort with equal keys*, SIAM J. Comput. **6** (1977) 240–287.
- [SPP76] A. Schönhage, M. Paterson and N. Pippenger, *Finding the median*, J. Comput. System Sci. **13** (1976) 184–199.
- [Val00] J. D. Valois, *Introspective sorting and selection revisited*, Software–Practice and Experience **30** (2000) 617–638.

