

# ON FLUCTUATIONS OF COMPLEXITY MEASURES FOR THE FIND ALGORITHM

JASPER ISCHEBECK AND RALPH NEININGER

**ABSTRACT.** The FIND algorithm (also called Quickselect) is a fundamental algorithm to select ranks or quantiles within a set of data. It was shown by Grübel and Rösler that the number of key comparisons required by Find as a process of the quantiles  $\alpha \in [0, 1]$  in a natural probabilistic model converges after normalization in distribution within the càdlàg space  $D[0, 1]$  endowed with the Skorokhod metric. We show that the process of the residuals in the latter convergence after normalization converges in distribution to a mixture of Gaussian processes in  $D[0, 1]$  and identify the limit's conditional covariance functions. A similar result holds for the related algorithm QuickVal. Our method extends to other cost measures such as the number of swaps (key exchanges) required by Find or cost measures which are based on key comparisons but take into account that the cost of a comparison between two keys may depend on their values, an example being the number of bit comparisons needed to compare keys given by their bit expansions.

## 1. INTRODUCTION

In 1961, Hoare [11] introduced the algorithm FIND, also called Quickselect, to select a key (an element) of a given rank from a linearly ordered finite set of data. We assume that the data are distinct real numbers. To be definite a simple version of the FIND algorithm is given as follows:  $\text{FIND}(S, k)$  has as input a set  $S = \{s_1, \dots, s_n\}$  of distinct real numbers of size  $n$  and an integer  $1 \leq k \leq n$ . The algorithm FIND operates recursively as follows: If  $n = 1$  we have  $k = 1$  and FIND returns the single element of  $S$ . If  $n \geq 2$  and  $S = \{s_1, \dots, s_n\}$  the algorithm first chooses an elements from  $S$ , say  $s_j$ , called pivot, and generates the sets

$$S_{<} := \{s_i \mid s_i < s_j, i \in \{1, \dots, n\} \setminus \{j\}\}, \quad S_{\geq} := \{s_i \mid s_i \geq s_j, i \in \{1, \dots, n\} \setminus \{j\}\}.$$

If  $k = |S_{<}| + 1$ , the algorithm returns  $s_j$ . If  $k \leq |S_{<}|$ , recursively  $\text{FIND}(S_{<}, k)$  is applied. If  $k \geq |S_{<}| + 2$ , recursively  $\text{FIND}(S_{\geq}, k - |S_{<}| - 1)$  is applied. Note that  $\text{FIND}(S, k)$  returns the element of rank  $k$  from  $S$ . There are various variants of the algorithm, in particular regarding how the pivot element is chosen and how  $S$  is partitioned into the subsets  $S_{<}$  and  $S_{\geq}$ .

In a standard probabilistic model one assumes that the data are ordered, i.e. given as a vector  $(s_1, \dots, s_n)$ , and are randomly permuted, all permutations being equally likely. This can be achieved assuming that the data are given as  $(U_1, \dots, U_n)$  where  $(U_j)_{j \in \mathbb{N}}$  is a sequence of i.i.d. random variables with distribution  $\text{unif}[0, 1]$ , the uniform distribution over the unit interval  $[0, 1]$ . This is the probabilistic model considered below. Note that the randomness is within the data, while the algorithm is deterministic.

Various cost measures have been considered for FIND, mainly the number of key comparisons required which we analyze in detail below. At the end of this extended abstract we state related results for the number of swaps (key exchanges) required and for cost measures which

---

*Date:* March 12, 2024.

*2020 Mathematics Subject Classification.* 60F17, 68Q25, 68P10, 60C05.

are based on key comparisons, where the cost of a comparison may depend on the values of the keys  $s_i, s_j$ , the number of bit comparisons required to decide whether  $s_i < s_j$  or not being a prominent example.

For analysis purposes a related process, called QuickVal, has been considered, see [18, 8]. Informally, QuickVal for an  $\alpha \in [0, 1]$  mimics FIND to select (or to try to select) the value  $\alpha$  from the set of data, which, in our probabilistic model for large  $n$ , comes close to FIND selecting rank  $\lfloor \alpha n \rfloor$ . To be definite, QuickVal $((U_1, \dots, U_n), \alpha)$  compares the  $U_i$  with  $U_1$  to generate sublists

$$S_{<} := (U_{j_1}, \dots, U_{j_{m-1}}), \quad S_{\geq} := (U_{j_{m+1}}, \dots, U_{j_n}),$$

with  $U_{j_i} < U_1$  for  $i = 1, \dots, m-1$  and  $2 \leq j_1 < \dots < j_{m-1}$  and  $U_{j_i} \geq U_1$  for  $i = m+1, \dots, n$  and  $2 \leq j_{m+1} < \dots < j_n$ . Hence,  $m-1 \in \{0, \dots, n-1\}$  is the number of the  $U_i$ ,  $2 \leq i \leq n$ , being smaller than  $U_1$ . The algorithm recursively calls QuickVal( $S_{<}, \alpha$ ) if  $\alpha < U_1$  and  $|S_{<}| > 0$ . If  $\alpha \geq U_1$  and  $|S_{\geq}| > 0$  recursively QuickVal( $S_{\geq}, \alpha - U_1$ ) is called. The number of key comparisons required by QuickVal $((U_1, \dots, U_n), \alpha)$  is denoted by  $S_{\alpha, n}$ .

To describe the processes  $(S_{\alpha, n})_{\alpha \in [0, 1]}$  and their limit (after scaling) conveniently we also consider the binary search tree constructed from the data  $(U_i)_{i \in \mathbb{N}}$ . Part of the following definitions are depicted in Figure 1. The data are inserted into the rooted infinite binary tree, where we denote its nodes by the elements of  $\{0, 1\}^* := \bigcup_{n=0}^{\infty} \{0, 1\}^n$  as follows. Its root is denoted by the empty word  $\epsilon$  and for each node  $\phi \in \{0, 1\}^*$  we denote by  $\phi 0$  and  $\phi 1$  (the word  $\phi$  appended with a 0 resp. 1) its left and right child respectively. Moreover  $|\phi|$  denotes the length of the word  $\phi$ , which is the depth of the corresponding node in the tree. To construct the binary search tree for  $(U_1, \dots, U_n)$  the first key  $U_1$  is inserted into the root and occupies the root. Then, successively the following keys are inserted, where each key traverses the already occupied nodes starting at the root as follows: Whenever the key traversing is less than the occupying key at a node it moves on to the left child of that node, otherwise to its right child. The first empty node found is occupied by the key.

To describe the costs of the algorithms we organize, using notation of Fill and Nakama [8], the sub-intervals  $([L_\phi, R_\phi])_{\phi \in \{0, 1\}^*}$  implicitly generated starting with  $[0, 1) =: [L_\epsilon, R_\epsilon)$  and recursively setting

$$(1) \quad \begin{aligned} \tau_\phi &:= \inf\{i \in \mathbb{N} \mid L_\phi < U_i < R_\phi\}, \\ L_{\phi 0} &:= L_\phi, \quad R_{\phi 1} := R_\phi, \quad L_{\phi 1} := R_{\phi 0} := U_{\tau_\phi}, \quad I_\phi := R_\phi - L_\phi. \end{aligned}$$

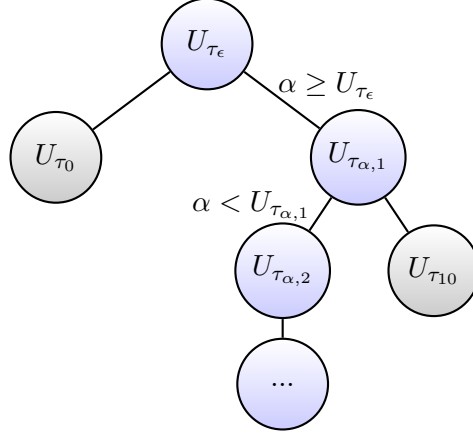
Now, if a sublist starting with pivot  $U_{\tau_\phi}$  has to be split by QuickVal, the keys which are inserted in the subtree rooted at  $U_{\tau_\phi}$  need to be compared with  $U_{\tau_\phi}$ . Hence, we get a contribution of key comparisons of

$$(2) \quad S_{\phi, n} = \sum_{\tau_\phi < k \leq n} \mathbf{1}_{[L_\phi, R_\phi)}(U_k).$$

Now, for  $\alpha \in [0, 1]$ , QuickVal $((U_1, \dots, U_n), \alpha)$  generates and splits sublists encoded by  $\phi(\alpha, k)$  for  $k = 0, 1, \dots$  for which we obtain by  $\phi(\alpha, 0) = \epsilon$  and

$$(3) \quad \phi(\alpha, k+1) = \begin{cases} \phi(\alpha, k)0, & \text{if } \alpha < U_{\tau_{\phi(\alpha, k)}}, \\ \phi(\alpha, k)1, & \text{if } \alpha \geq U_{\tau_{\phi(\alpha, k)}}. \end{cases}$$

When using the variables defined in (1) or (2), we abbreviate the notation  $\phi(\alpha, k)$  by  $\alpha, k$ , such as writing  $I_{\alpha, k} := I_{\phi(\alpha, k)}$  or  $S_{\alpha, k, n} := S_{\phi(\alpha, k), n}$ .



**Figure 1.** Part of the binary search tree. The pivots of sublists split by  $\text{QuickVal}((U_1, \dots, U_n), \alpha)$  for some  $\alpha \in [0, 1]$  are on the path indicated. Note that we have  $\tau_\epsilon = \tau_{\phi(\alpha, 0)} = \tau_{\alpha, 0} = 1$  and in this example  $\alpha \geq U_1$  and  $\alpha < U_{\tau_{\alpha, 1}}$  so that  $\phi(\alpha, 2) = 10 \in \{0, 1\}^2$ .

The number of key comparisons required by  $\text{QuickVal}((U_1, \dots, U_n), \alpha)$  is then given by the (finite) sum

$$S_{\alpha, n} = \sum_{k=1}^{\infty} S_{\alpha, k, n}.$$

Fill and Nakama [8, Theorem 3.2] showed (considering more general complexity measures) that for each  $\alpha \in [0, 1]$  almost surely

$$(4) \quad \frac{1}{n} S_{\alpha, n} \rightarrow S_\alpha := \sum_{k=0}^{\infty} I_{\alpha, k}, \quad (n \rightarrow \infty).$$

The latter convergence also holds in  $L_p$ , see Fill and Matterer [7, Proposition 6.1].

We take the point of view that such an almost sure asymptotic result may be considered a strong law of large numbers (SLLN). The subject of the present extended abstract is to study the fluctuations in such SLLN, sometimes called a central limit analogue. We study these fluctuations as processes in the metric space  $(D[0, 1], d_{\text{SK}})$  of càdlàg functions endowed with the Skorokhod metric; see Billingsley [2] for background on weak convergence of probability measures on metric spaces in general and on  $(D[0, 1], d_{\text{SK}})$  in particular. Note that, by definition,  $(S_{\alpha, n})_{\alpha \in [0, 1]}$  and  $(S_\alpha)_{\alpha \in [0, 1]}$  have càdlàg paths almost surely. As the normalized process of fluctuations we define

$$(5) \quad G_n := (G_{\alpha, n})_{\alpha \in [0, 1]} := \left( \frac{S_{\alpha, n} - nS_\alpha}{\sqrt{n}} \right)_{\alpha \in [0, 1]}.$$

Then we have the following result:

**Theorem 1.** *Let  $S_{\alpha, n}$  be the number of key comparisons required by  $\text{QuickVal}((U_1, \dots, U_n), \alpha)$  and  $(S_\alpha)_{\alpha \in [0, 1]}$  as in (4). Then for the fluctuation process  $G_n$  defined in (5) we have*

$$G_n \xrightarrow{d} G_\infty \text{ in } (D[0, 1], d_{\text{SK}}) \quad (n \rightarrow \infty),$$

where  $G_\infty$  is a mixture of centered Gaussian processes with random covariance function given by

$$(6) \quad \Sigma_{\infty, \alpha, \beta} := \sum_{k=0}^J \sum_{j=0}^{\infty} I_{\alpha, j \vee k} + \mathbf{1}_{\{\alpha \neq \beta\}} (J+1) \sum_{j=J+1}^{\infty} (I_{\beta, j}) - S_\alpha S_\beta, \quad \alpha, \beta \in [0, 1],$$

where  $J = J(\alpha, \beta) := \max\{k \in \mathbb{N}_0 \mid \tau_{\alpha, k} = \tau_{\beta, k}\} \in \mathbb{N}_0 \cup \{\infty\}$ .

*Remark 2.* In his PhD thesis, Matterer [14, Theorem 6.4] showed the convergence of the one-dimensional marginals for the functional limit law in Theorem 1.

*Remark 3.* An alternative representation of the random covariance function in (6) is as follows: With an independent random variable  $V$  uniformly distributed over  $[0, 1]$ , we have

$$(7) \quad \Sigma_{\infty, \alpha, \beta} = \text{Cov}(J(V, \alpha), J(V, \beta) \mid \mathcal{F}_\infty),$$

with the  $\sigma$ -algebra

$$(8) \quad \mathcal{F}_\infty := \sigma\{I_\phi \mid \phi \in \{0, 1\}^*\}.$$

*Remark 4.* A related functional limit law for the complexity of Radix Selection, an algorithm to select ranks based on the bit expansions of the data, with a limiting Gaussian process with a covariance function related to (7) can be found in [12, Theorem 1.2]. See [17, Theorem 1.1] for another related functional limit law.

The analysis of QuickVal is usually considered an intermediate step to analyze the original FIND algorithm. Grübel and Rösler [9] already pointed out that a version of FIND such as stated above with  $C_n^*(k)$  denoting the number of key comparisons for finding rank  $k$  within  $(U_1, \dots, U_n)$  does not lead to convergence within  $(D[0, 1], d_{\text{SK}})$  after the normalization  $\alpha \mapsto \frac{1}{n} C_n^*(\lfloor \alpha n \rfloor + 1)$ , where here and below the convention  $C_n^*(n+1) := C_n^*(n)$  is used. To overcome this problem they propose a version that does not stop in case a pivot turns out to be the rank to be selected by including the pivot in the list  $S_<$  and proceeding until a list of size 1 is generated. Moreover, their pivots are chosen uniformly at random. The number of key comparisons  $C'_n(k)$  for Grübel and Rösler's FIND-version has the property that

$$(9) \quad \left( \frac{1}{n} C'_n(\lfloor \alpha n \rfloor + 1) \right)_{\alpha \in [0, 1]} \xrightarrow{d} (S_\alpha)_{\alpha \in [0, 1]}, \quad \text{in } (D[0, 1], d_{\text{SK}}),$$

see [9, Theorem 4]. Without using random pivots we may also obtain right-continuous limits by just recursively calling  $\text{FIND}(S_{\geq}, 0)$  in case the pivot turns out to be the rank sought. We denote the number of key comparisons for this version by  $C_n(k)$ , which is close to Grübel and Rösler's FIND-version and also satisfies (9).

The convergence in (9) could only be stated weakly (not almost surely) since Grübel und Rösler's FIND-version due to randomization within the algorithm does not have a natural embedding on a probability space. Note that the formulation of the QuickVal complexity does have such an embedding which, e.g., makes the almost sure convergence in (4) possible. However, it is easy to see that we have the distributional equality

$$(10) \quad (C_n(|\{U_i \leq \alpha : 1 \leq i \leq n\}|))_{\alpha \in [0, 1]} \stackrel{d}{=} (S_{\alpha, n})_{\alpha \in [0, 1]}.$$

This allows to naturally couple the complexities on one probability space, which we call its *natural coupling*. See [7, page 807] for a related discussion of natural couplings.

To transfer Theorem 1 to FIND we need to align jumps to come up with a suitable fluctuation process. The conventions  $C_n(0) := C_n(1)$  and  $C_n(n+1) := C_n(n)$  are used.

**Corollary 5.** *Let  $C_n(k)$  be the number of key comparisons required to select rank  $1 \leq k \leq n$  within a set of  $n$  data by FIND with the natural coupling (10). Let  $\Lambda_n : [0, 1] \rightarrow [0, 1]$ ,  $n \in \mathbb{N}$ , be any (random) monotone increasing bijective function such that  $\Lambda_n(\frac{k}{n+1})$  is equal to the element of rank  $k$  within  $\{U_1, \dots, U_n\}$ . Then we have*

$$\left( \frac{C_n(\lfloor t(n+1) \rfloor) - nS_{\Lambda_n(t)}}{\sqrt{n}} \right)_{t \in [0,1]} \xrightarrow{d} G_\infty \text{ in } (D[0,1], d_{\text{SK}}),$$

where  $G_\infty$  is the process defined in Theorem 1.

The extended abstract is organized as follows: In Section 2 we introduce a novel perturbation argument which is the basis of our analysis. Section 3 contains a criterion for weak convergence of probability measures on  $(D[0,1], d_{\text{SK}})$ , which is applied in Section 4 to proof Theorem 1 and Corollary 5. In Section 5 further functional fluctuation results are stated for the number of swaps (key exchanges) required by QuickVal (depending on the specific algorithm used to partition  $S$  into the sublists  $S_<$  and  $S_>$ ) as well as functional fluctuation results for cost measures which are based on key comparisons, where the cost of a comparison may depend on the values of the keys.

## 2. PERTURBATION OF THE DATA

QuickVal splits an interval  $[L_\phi, R_\phi)$  by the first value falling into  $[L_\phi, R_\phi)$  denoted by  $U_{\tau_\phi}$ . Obviously, this implies dependencies between the data  $U_i$  and the lengths  $I_\phi$  of the intervals  $[L_\phi, R_\phi)$ . In the present section we construct a perturbed sequence  $(\tilde{U}_i)_{i \in \mathbb{N}}$  to the data  $(U_i)_{i \in \mathbb{N}}$  such that we gain independence of  $(\tilde{U}_i)_{i \in \mathbb{N}}$  from the  $\sigma$ -algebra  $\mathcal{F}_\infty$  of the interval lengths defined in (8). In particular, we aim that conditional on  $\mathcal{F}_\infty$  the number of data  $(\tilde{U}_1, \dots, \tilde{U}_n)$  falling into an interval  $[L_\phi, R_\phi)$  is binomial  $B(n, I_\phi)$  distributed, see Lemma 7 below.

Every value  $U_i$ ,  $i \in \mathbb{N}$ , falls successively into subintervals generated by QuickVal until becoming a pivot element. These subintervals correspond to the path between the root of the corresponding binary search tree and the node where  $U_i$  is inserted. Let  $\phi_i \in \{0, 1\}^*$  denote the node where  $U_i$  is inserted. Hence, we have  $\tau_{\phi_i} = i$  and  $U_i = L_{\phi_i} + I_{\phi_i}0$ .

Let  $(V_i)_{i \in \mathbb{N}}$  be a sequence of i.i.d.  $\text{unif}[0, 1]$  random variables being independent of  $(U_i)_{i \in \mathbb{N}}$ . We define

$$(11) \quad \tilde{U}_i := L_{\phi_i} + I_{\phi_i} V_i.$$

**Lemma 6.** *The sequence  $(\tilde{U}_i)_{i \in \mathbb{N}}$  defined in (11) consists of i.i.d.  $\text{unif}[0, 1]$  distributed random variables and is independent of  $\mathcal{F}_\infty$ .*

*Proof.* It suffices to show that  $\tilde{U}_i$  conditional on  $\mathcal{F}_\infty$  and  $\tilde{U}_1, \dots, \tilde{U}_{i-1}$  is uniformly distributed on  $[0, 1]$  for all  $i \in \mathbb{N}$ . We use infinitesimal notation to denote this claim by

$$\mathbb{P}(\tilde{U}_i \in du \mid \mathcal{F}_\infty, \tilde{U}_1, \dots, \tilde{U}_{i-1}) = \mathbf{1}_{[0,1]}(u)du, \quad i \in \mathbb{N}.$$

For each  $i \in \mathbb{N}$  the random variables  $\tilde{U}_i$  and  $U_i$  fall into the same interval  $[L_{\phi_i}, R_{\phi_i})$ , hence  $\phi_1, \dots, \phi_{i-1}$  are determined by  $\tilde{U}_1, \dots, \tilde{U}_{i-1}$ . Let us additionally condition on  $\phi_i$ , then, by definition,

$$\mathbb{P}(\tilde{U}_i \in du \mid \mathcal{F}_\infty, \tilde{U}_1, \dots, \tilde{U}_{i-1}, \phi_i) = \frac{1}{I_{\phi_i}} \mathbf{1}_{[L_{\phi_i}, R_{\phi_i})}(u)du.$$

Note that  $\phi_i$  denotes one of the  $i$  external nodes of the binary search tree with internal nodes denoted by  $\phi_1, \dots, \phi_{i-1}$ . We denote by  $\text{Ext}_{i-1}$  the set of the labels of these external nodes. Hence, conditional on  $\mathcal{F}_\infty, \phi_1, \dots, \phi_{i-1}$  the label  $\phi_i$  is chosen from  $\text{Ext}_{i-1}$  with probability given by the length of the corresponding interval, i.e.,  $\mathbb{P}(\phi_i = \phi \mid \mathcal{F}_\infty, \phi_1, \dots, \phi_{i-1}) = I_\phi$  for all  $\phi \in \text{Ext}_{i-1}$ . Thus, by the law of total probability we obtain

$$\mathbb{P}(\tilde{U}_i \in du \mid \mathcal{F}_\infty, \tilde{U}_1, \dots, \tilde{U}_{i-1}) = \sum_{\phi \in \text{Ext}_{i-1}} I_\phi \frac{1}{I_\phi} \mathbf{1}_{[L_\phi, R_\phi)}(u) du = \mathbf{1}_{[0,1]}(u) du.$$

This implies the assertion.  $\square$

The  $\tilde{U}_i$  are now coupled with the  $U_i$  but independent of the  $I_\phi$ . To compare with the number of key comparisons required by  $\text{QuickVal}((U_1, \dots, U_n), \alpha)$  we define

$$\tilde{S}_{\alpha,k,n} := \sum_{i=1}^n \mathbf{1}_{[L_{\alpha,k}, R_{\alpha,k})}(\tilde{U}_i).$$

**Lemma 7.** *Conditional on  $I_{\alpha,k}$  we have that  $\tilde{S}_{\alpha,k,n}$  has the binomial  $B(n, I_{\alpha,k})$  distribution. Moreover, for all  $\alpha \in [0, 1]$ ,  $n \in \mathbb{N}$  and  $0 \leq k \leq n$  we have*

$$(12) \quad S_{\alpha,k,n} \leq (\tilde{S}_{\alpha,k,n} - 1)^+ \leq S_{\alpha,k,n} + k - 1.$$

*Proof.* The conditional distribution of  $\tilde{S}_{\alpha,k,n}$  follows from Lemma 6. Recall that  $S_{\alpha,k,n}$  is defined as  $\sum_{i=\tau_{\alpha,k}}^n \mathbf{1}_{\{L_{\alpha,k-1} \leq U_i < R_{\alpha,k-1}\}}$ . By definition,  $U_i$  and  $\tilde{U}_i$  are in the interval  $[L_{\phi_i}, R_{\phi_i})$  for all  $i \in \mathbb{N}$ . If  $U_i \in (L_{\alpha,k}, R_{\alpha,k})$ , then  $U_i$  appears as a pivot after the  $k$ -th pivot. Hence, its interval  $[L_{\phi_i}, R_{\phi_i})$  and thus also  $\tilde{U}_i$  are contained in  $(L_{\alpha,k}, R_{\alpha,k})$ . The  $k$ -th pivot  $U_{\tau_{\alpha,k}}$  itself does not contribute to  $S_{\alpha,k,n}$ , which implies the left inequality stated in the present lemma.

For the right inequality, assume for some  $i \in \mathbb{N}$  that the perturbed value  $\tilde{U}_i$  is in  $(L_{\alpha,k}, R_{\alpha,k})$ , but  $U_i$  is not. Then the corresponding interval  $(L_{\phi_i}, R_{\phi_i})$  must contain  $(L_{\alpha,k}, R_{\alpha,k})$ , thus making  $U_i$  a pivot that appears before the  $k$ -th pivot. Since there are only  $k$  such pivots, the right inequality follows.  $\square$

### 3. ON WEAK CONVERGENCE IN $D[0, 1]$

To prove the convergence in distribution in Theorem 1 within the space  $(D[0, 1], d_{SK})$  we use the following Proposition 8. It can be proved by classical tools of weak convergence theory based on a study of the modulus of continuity and the Arzelà–Ascoli theorem in form of a general theorem of Billingsley [2, Theorem 13.2].

**Proposition 8.** *Let  $X_1, X_2, \dots$  be a sequence of random variables in  $(D[0, 1], d_{SK})$ . Suppose that for every  $K \in \mathbb{N}$ , there exist random càdlàg step functions  $X_1^K, X_2^K, \dots$  with all jumps contained in  $\{U_\phi \mid \phi \in \{0, 1\}^*, |\phi| < K\}$ . If*

- (i) *for all  $r \in \mathbb{N}$  and  $\alpha_1, \dots, \alpha_r \in [0, 1]$ , the marginals  $\mathcal{L}(X_n(\alpha_1), \dots, X_n(\alpha_r))$  converge weakly to some distribution  $\mu_{\alpha_1, \dots, \alpha_r}$ ,*
- (ii) *for all  $\varepsilon > 0$ ,*

$$(13) \quad \lim_{K \rightarrow \infty} \limsup_{n \rightarrow \infty} \mathbb{P}(\|X_n - X_n^K\|_\infty > \varepsilon) \rightarrow 0,$$

then  $(X_n)_{n \in \mathbb{N}}$  converges in distribution to a random variable  $X$  on  $(D[0, 1], d_{SK})$ , and for all  $r \in \mathbb{N}$  and  $\alpha_1, \dots, \alpha_r \in [0, 1]$  we have

$$(14) \quad \mathcal{L}(X(\alpha_1), \dots, X(\alpha_r)) = \mu_{\alpha_1, \dots, \alpha_r}.$$

#### 4. PROOF OF THEOREM 1

To split the contributions to the process  $G_n$  into costs resulting from above and below a level  $K \in \mathbb{N}$  we define

$$(15) \quad G_{\alpha, k, n} := \frac{S_{\alpha, k, n} - nI_{\alpha, k}}{\sqrt{n}}$$

as the normalized fluctuations of the contribution at level  $k$ , and set

$$(16) \quad G_{\alpha, n}^{\leq K} := \sum_{k=0}^K G_{\alpha, k, n}, \quad G_n^{\leq K} := (G_{\alpha, n}^{\leq K})_{\alpha \in [0, 1]}, \quad G_{\alpha, n}^{> K} := \sum_{k=K+1}^{\infty} G_{\alpha, k, n}.$$

Hence,  $G_{\alpha, n} = G_{\alpha, n}^{\leq K} + G_{\alpha, n}^{> K}$ . Analogously, for the perturbed values  $\tilde{S}_{k, n}$  we define

$$(17) \quad W_{\alpha, k, n} := \frac{\tilde{S}_{\alpha, k, n} - nI_{\alpha, k}}{\sqrt{n}}, \quad W_{\alpha, n}^{\leq K} := \sum_{k=0}^K W_{\alpha, k, n}, \quad W_n^{\leq K} := (W_{\alpha, n}^{\leq K})_{\alpha \in [0, 1]}.$$

**Lemma 9.** *For all  $K \in \mathbb{N}$  we have convergence in distribution of  $(G_n^{\leq K})_{n \in \mathbb{N}}$  towards a mixture  $G_{\infty}^{\leq K} = (G_{\alpha, \infty}^{\leq K})_{\alpha \in [0, 1]}$  of centered Gaussian processes within  $\|\cdot\|_{\infty}$ . Conditional on  $\mathcal{F}_{\infty}$ , the limit  $G_{\infty}^{\leq K}$  is a centered Gaussian process with covariance function given, for  $\alpha, \beta \in [0, 1]$  by*

$$(18) \quad \text{Cov}(G_{\alpha, \infty}^{\leq K}, G_{\beta, \infty}^{\leq K} \mid \mathcal{F}_{\infty}) = \sum_{k=0}^{K \wedge J} \sum_{j=0}^K I_{\alpha, j \vee k} + (1 + (K \wedge J)) \sum_{j=J+1}^K I_{\beta, j} - S_{\alpha}^{\leq K} S_{\beta}^{\leq K},$$

where  $J = J(\alpha, \beta)$  is as in Theorem 1 and  $S_{\alpha}^{\leq K} := \sum_{k=0}^K I_{\alpha, k}$ . The stated convergence in distribution also holds conditionally in  $\mathcal{F}_{\infty}$ , i.e., we have almost surely that  $\mathcal{L}(G_n^{\leq K} \mid \mathcal{F}_{\infty})$  converges weakly towards  $\mathcal{L}(G_{\infty}^{\leq K} \mid \mathcal{F}_{\infty})$ .

*Proof.* First note that by Lemma 7 we have  $\|G_n^{\leq K} - W_n^{\leq K}\|_{\infty} < K^2/\sqrt{n}$ , so it suffices to show the lemma for  $W_n^{\leq K}$ . Conditional on  $\mathcal{F}_{\infty}$ , the value of  $W_{\alpha, n}^{\leq K}$  is given by

$$(19) \quad W_{\alpha, n}^{\leq K} = \frac{1}{\sqrt{n}} \sum_{i=1}^n \sum_{k=0}^K \mathbf{1}\{L_{\alpha, k} \leq \tilde{U}_i < R_{\alpha, k}\} - (R_{\alpha, k} - L_{\alpha, k}),$$

thus the  $2^k$  different values of the process  $W_n^{\leq K}$  can be expressed as the sum of  $n$  centered, bounded i.i.d. random vectors, scaled by  $1/\sqrt{n}$ . By the multivariate central limit theorem, these converge towards a multivariate, centered normal variable. As the positions of the jumps, still conditional on  $\mathcal{F}_{\infty}$ , are fixed, we have convergence of  $W_n^{\leq K}$  and thus also of  $G_n^{\leq K}$  towards a Gaussian process. Define  $X_{\alpha, k} := \mathbf{1}\{L_{\alpha, k-1} \leq \tilde{U}_1 < R_{\alpha, k-1}\}$ . The covariance

function then is given by

$$\begin{aligned}
\text{Cov}(G_{\alpha,\infty}^{\leq K}, G_{\beta,\infty}^{\leq K} \mid \mathcal{F}_\infty) &= \text{Cov}\left(\sum_{k=0}^K X_{\alpha,k}, \sum_{j=0}^K X_{\beta,j} \mid \mathcal{F}_\infty\right) \\
&= \sum_{k=0}^K \sum_{j=0}^K \mathbb{E}\left[X_{\alpha,k} X_{\beta,j} \mid \mathcal{F}_\infty\right] - S_\alpha^{\leq K} S_\beta^{\leq K} \\
(20) \quad &= \sum_{k=0}^{K \wedge J} \sum_{j=0}^K I_{\alpha,k \vee j} + (1 + (K \wedge J)) \sum_{j=J+1}^K I_{\beta,j} - S_\alpha^{\leq K} S_\beta^{\leq K}.
\end{aligned}$$

The assertion follows.  $\square$

To see that the covariance functions in (20) converge towards the covariance function of  $G_\infty$  stated in Theorem 1 we restate a Lemma of Grübel and Rösler [9, Lemma 1] that the maximal length of the intervals at a level is decreasing geometrically with increasing levels. It is obtained observing that  $\mathbb{E}\left[\sum_{|\phi|=k} I_\phi^2\right] = (2/3)^k$  and states:

**Lemma 10.** *There exists an a.s. finite random variable  $K_1$  such that for all  $k \geq K_1$ :*

$$(21) \quad \max_{\alpha \in [0,1]} I_{\alpha,k} \leq k \left(\frac{2}{3}\right)^{k/2}.$$

Lemma 10 implies that the covariance functions of  $G_\infty^{\leq K}$  from (18) converge a.s. to the covariance function of  $G_\infty$  from (6).

For the costs from levels  $k > K$  we find:

**Proposition 11.** *For all  $\varepsilon, \eta > 0$  there are constants  $K, N \in \mathbb{N}$  such that for all  $n \geq N$*

$$(22) \quad \mathbb{P}\left(\|G_n^{>K}\|_\infty > \eta\right) < \varepsilon.$$

We postpone the proof of the latter proposition and first use Proposition 11 and Lemma 9 to show convergence of the finite-dimensional distributions, denoted fdd-convergence.

**Lemma 12.** *We have fdd-convergence of  $G_n$  towards  $G_\infty$ .*

*Proof.* For any  $K$ , we can split  $G_n = G_n^{\leq K} + G_n^{>K}$ . By Lemma 9, we have

$$G_n^{\leq K} \xrightarrow{\text{fdd}} G_\infty^{\leq K} \quad (n \rightarrow \infty).$$

Furthermore, because the covariance functions of the  $G_\infty^{\leq K}$  converge a.s., we obtain

$$G_\infty^{\leq K} \xrightarrow{\text{fdd}} G_\infty \quad (K \rightarrow \infty)$$

by Lévy's continuity theorem. Hence, for all  $\alpha_1, \dots, \alpha_\ell \in [0, 1]$  and all  $t_1, \dots, t_\ell \in \mathbb{R}$  we find a sequence  $(K_n)_{n \in \mathbb{N}}$  in  $\mathbb{N}$  such that

$$\mathbb{P}\left(G_{\alpha_1,n}^{\leq K_n} < t_1, \dots, G_{\alpha_\ell,n}^{\leq K_n} < t_\ell\right) \longrightarrow \mathbb{P}\left(G_{\alpha_1,\infty} < t_1, \dots, G_{\alpha_\ell,\infty} < t_\ell\right) \quad (n \rightarrow \infty).$$

Now, since  $\|G_n^{>K_n}\|_\infty \rightarrow 0$  in probability by Proposition 11 the claim of Lemma 12 follows by Slutsky's theorem.  $\square$

To prepare for the proof of Proposition 11, we show that the fluctuations on each level are also at least geometrically decreasing. Recall  $K_1$  from Lemma 10.



**Lemma 13.** *There exists a constant  $a > 1$  such that for all  $k, n \in \mathbb{N}$*

$$\mathbb{P} \left( \max_{\alpha \in [0,1]} |W_{\alpha,k,n}| > a^{-k}, K_1 \leq k \right) \leq b(k) + c(k, n)$$

with  $b(k), c(k, n) \geq 0$  such that

$$(23) \quad \sum_{k=1}^{\infty} b(k) < \infty \quad \text{and} \quad \sum_{k=1}^{\lfloor (9/2) \log n \rfloor} c(k, n) \rightarrow 0 \quad (n \rightarrow \infty).$$

For the proof of Lemma 13 we require the following Chernoff bound:

**Lemma 14.** *Let  $S_n$  be binomial  $B(n, p)$  distributed for some  $p \in [0, 1]$  and  $n \in \mathbb{N}$  and let  $\mu := \mathbb{E}[S_n]$ ,  $\varepsilon \geq 0$ . Then*

$$\mathbb{P} \left( S_n \notin ((1 - \varepsilon)\mu, (1 + \varepsilon)\mu) \right) \leq 2 \exp \left( -\frac{\varepsilon^2 \mu}{2 + \varepsilon} \right).$$

*Proof.* Combine upper and lower bound in McDiarmid [15, Theorem 2.3].  $\square$

*Proof of Corollary 13.* Fix some  $\alpha \in [0, 1]$ . Conditionally on  $I_{k,\alpha}$ , the costs  $\tilde{S}_{\alpha,k,n}$  are  $B(n, I_{k,\alpha})$ -distributed by Lemma 6. The Chernoff bound in Lemma 14 implies

$$(24) \quad \begin{aligned} \mathbb{P} \left( |W_{\alpha,k,n}| > a^{-k} \mid I_{\alpha,k} \right) &= \mathbb{P} \left( |\tilde{S}_{\alpha,k,n} - nI_{\alpha,k}| > \sqrt{n}a^{-k} \mid I_{\alpha,k} \right) \\ &\leq 2 \exp \left( -\frac{na^{-2k}}{nI_{\alpha,k}(2 + \sqrt{n}a^{-k}/(nI_{\alpha,k}))} \right) \\ &= 2 \exp \left( -\left( 2a^{2k}I_{\alpha,k} + a^k/\sqrt{n} \right)^{-1} \right). \end{aligned}$$

For the two summands in the exponent in (24) we have the following behavior: Summand  $2a^{2k}I_{\alpha,k}$  is falling geometrically with  $k$  for sufficiently small  $a > 1$ . Summand  $a^k/\sqrt{n}$  is falling with  $n$ , but growing with  $k$ . To separate these two contributions, note that  $\exp(x^{-1}) \geq \frac{1}{m!}x^{-m}$  and thus  $\exp(-x^{-1}) \leq m!x^m$  for all  $m \in \mathbb{N}$  and  $x \geq 0$ . Choosing  $m = 7$ , we obtain

$$\begin{aligned} \mathbb{P} \left( |W_{\alpha,k,n}| > a^{-k} \mid I_{\alpha,k} \right) &\leq 2 \cdot 7! \left( 2a^{2k}I_{\alpha,k} + a^k/\sqrt{n} \right)^7 \\ &\leq 2^{14}7!a^{14k}I_{\alpha,k}^7 + 2^77!a^{7k}n^{-7/2} \end{aligned}$$

by convexity of  $x \mapsto x^7$ . Note that the  $2^k$  intervals at level  $k$  have lengths  $I_{\alpha,k}$  summing to 1. Hence,

$$\mathbb{P} \left( \max_{\alpha \in [0,1]} |W_{\alpha,k,n}| > a^{-k} \mid I_{\alpha,k} \right) \leq 2^{14}7!a^{14k} \max_{\alpha \in [0,1]} I_{\alpha,k}^6 + 2^77!(2a^7)^k n^{-7/2}.$$

When furthermore  $K_1 \leq k$ , by Lemma 10 the length  $I_{\alpha,k}$  is bounded by  $k^{\frac{2}{3}}$ , hence

$$\mathbb{P} \left( \max_{\alpha \in [0,1]} |W_{\alpha,k,n}| > a^{-k}, K_1 \leq k \right) \leq 2^{14}7!ka^{14k} \left( \frac{2}{3} \right)^{3k} + 2^77!(2a^7)^k n^{-7/2}.$$

Define the first summand on the right hand side of the latter inequality by  $b(k)$  and the second summand by  $c(k, n)$ . For all  $1 < a < (3/2)^{3/14}$  the  $b(k)$  form a convergent series. To

also show the second statement in (23) note that

$$\sum_{k=1}^{\lfloor (9/2) \log n \rfloor} c(k, n) = O\left((2a^7)^{(9/2) \log n} n^{-7/2}\right) = O\left(n^{(9/2) \log 2 + (9/2) \cdot 7 \log a - 7/2}\right).$$

The latter O-term converges to 0 for  $(63/2) \log a < 7/2 - (9/2) \log 2 \approx 0.381 \dots$ , thus we may choose  $a$  as required.  $\square$

We are now prepared for the proof of Proposition 11.

*Proof of Proposition 11.* Let  $\varepsilon, \eta > 0$ . To  $K_1$  from Lemma 10 and  $a$  and  $b(k)$  from Lemma 13 we choose  $K$  sufficiently large such that

$$(25) \quad \mathbb{P}(K_1 > K) \leq \frac{\varepsilon}{4}, \quad \sum_{k=K}^{\infty} a^{-k} \leq \frac{\eta}{4} \quad \text{and} \quad \sum_{k=K}^{\infty} b(k) \leq \frac{\varepsilon}{4}.$$

Let  $H_n$  be the maximum amount of steps needed by  $\text{QuickVal}((U_1, \dots, U_n), \alpha)$  for any  $\alpha$ . Thus,  $H_n$  is also the height of the binary search tree built from  $U_1, \dots, U_n$ . Devroye [4] showed that the height has expectation  $\mathbb{E}[H_n] = \gamma \log n + o(\log n)$  with  $\gamma = 4.311 \dots$ . Reed [16] further showed that  $\text{Var}(H_n) = O(1)$ . Hence, we can choose  $N$  sufficiently large such that

$$(26) \quad \mathbb{P}(H_n > \lfloor (9/2) \log n \rfloor) < \frac{\varepsilon}{4}, \quad \sum_{k=1}^{\lfloor (9/2) \log n \rfloor} c(k, n) \leq \frac{\varepsilon}{4}, \quad \text{and} \quad \frac{((9/2) \log n)^2}{\sqrt{n}} \leq \frac{\eta}{4}$$

for all  $n \geq N$ . Subsequently we use the decomposition

$$G_{\alpha, n}^{>K} = \sum_{k=K+1}^{\lfloor (9/2) \log n \rfloor} \frac{S_{\alpha, k, n} - nI_{\alpha, k}}{\sqrt{n}} + \sum_{\lfloor (9/2) \log n \rfloor + 1}^{\infty} \frac{S_{\alpha, k, n} - nI_{\alpha, k}}{\sqrt{n}} =: \Gamma_n + G_{\alpha, n}^{>\lfloor (9/2) \log n \rfloor}$$

and consider the event

$$A_n := \{K_1 > K\} \cup \{H_n > \lfloor (9/2) \log n \rfloor\}.$$

We have  $\mathbb{P}(A_n) < \varepsilon/2$  for all  $n \geq N$ . Note that on  $A_n^c$  (the complement of  $A_n$ ) we have  $S_{\alpha, k, n} = 0$  for all  $k > \lfloor (9/2) \log n \rfloor$  and also the bound on  $I_{\alpha, k}$  from Lemma 10 applies, hence

$$\begin{aligned} \left| G_{\alpha, n}^{>\lfloor (9/2) \log n \rfloor} \right| &\leq \sum_{\lfloor (9/2) \log n \rfloor + 1}^{\infty} \sqrt{n} I_{\alpha, k} \leq \sum_{\lfloor (9/2) \log n \rfloor + 1}^{\infty} \sqrt{n} k \left( \frac{2}{3} \right)^{k/2} \\ &= O\left(n^{1/2 - (9/4) \log(3/2)} \log n\right) = o(1) \end{aligned}$$

since  $(9/4) \log(3/2) = 0.912 \dots$ . Hence, we can enlarge  $N$  so that on  $A_n^c$  we have  $\left| G_{\alpha, n}^{>\lfloor (9/2) \log n \rfloor} \right| < \eta/2$  for all  $n \geq N$ . This implies the bound

$$(27) \quad \begin{aligned} \mathbb{P}(|G_{\alpha, n}^{>K}| > \eta) &\leq \mathbb{P}(A_n) + \mathbb{P}(\{|G_{\alpha, n}^{>K}| > \eta\} \cap A_n^c) \\ &\leq \frac{\varepsilon}{2} + \mathbb{P}\left(\left\{|\Gamma_n| > \frac{\eta}{2}\right\} \cap A_n^c\right) + \mathbb{P}\left(\left\{|G_{\alpha, n}^{>\lfloor (9/2) \log n \rfloor}| > \frac{\eta}{2}\right\} \cap A_n^c\right). \end{aligned}$$

Note that the third summand in (27) is 0. Hence, it remains to bound the second summand in (27). To this end note that

$$\begin{aligned}
 |\Gamma_n| &\leq \sup_{\alpha \in [0,1]} \sum_{k=K+1}^{\lfloor (9/2) \log n \rfloor} \left( \left| \frac{\tilde{S}_{\alpha,k,n} - nI_{\alpha,k}}{\sqrt{n}} \right| + \left| \frac{S_{\alpha,k,n} - \tilde{S}_{\alpha,k,n}}{\sqrt{n}} \right| \right) \\
 (28) \quad &\leq \left( \sum_{k=K+1}^{\lfloor (9/2) \log n \rfloor} \max_{\alpha \in [0,1]} |W_{\alpha,k,n}| \right) + \frac{\lfloor (9/2) \log n \rfloor^2}{\sqrt{n}},
 \end{aligned}$$

where Lemma 7 is used. The third relation in (26) assures that the second term in (28) is smaller than  $\eta/4$ . In view of the second relation in (25) and (28), we have

$$\left\{ |\Gamma_n| > \frac{\eta}{2} \right\} \cap A_n^c \subset \bigcup_{k=K}^{\lfloor (9/2) \log n \rfloor} \left\{ \max_{\alpha \in [0,1]} |W_{\alpha,k,n}| > a^{-k}, K_1 \leq k \right\}.$$

Thus, Lemma 13 together with (25) and (26) imply that the second summand in (27) is bounded by  $\varepsilon/2$ . This implies the assertion.  $\square$

*Proof of Theorem 1.* We apply Lemma 8 to  $G_n$  and  $G_n^{\leq K}$ . The first condition, fdd convergence, is Lemma 12, the second condition is Proposition 11.  $\square$

We now transfer the fluctuation result for QuickVal in Theorem 1 to the original Find process in Corollary 5.

*Proof of Corollary 5.* Let  $\tilde{F}_n$  be the inverse of  $\Lambda_n$  in the statement of Corollary 5. By definition of  $\Lambda_n$ , the value of the element  $U_{(k)}$  of rank  $k$  within  $U_1, \dots, U_n$  is given by  $\frac{k}{n+1}$ , so

$$(29) \quad \lfloor (n+1)\tilde{F}_n(\alpha) \rfloor = |\{U_i \leq \alpha \mid 1 \leq i \leq n\}|$$

for all  $\alpha \in [0, 1)$ . Thus,  $C_n(\lfloor (n+1)\tilde{F}_n(\alpha) \rfloor) = S_{\alpha,n}$  a.s. for all  $\alpha \in [0, 1]$ , see (10). For  $\alpha = 1$  note that  $\tilde{F}_n(\alpha) = 1$  and  $C_n(n+1) = C_n(n)$  by definition. The Skorokhod distance  $d_{\text{SK}}$  is then bounded by

$$\begin{aligned}
 d_{\text{SK}} \left( G_n, \left( \frac{C_n(\lfloor t(n+1) \rfloor) - nS_{\Lambda_n(t)}}{\sqrt{n}} \right)_{t \in [0,1]} \right) &= d_{\text{SK}} \left( G_n, \left( \frac{S_{\Lambda_n(t),n} - nS_{\Lambda_n(t)}}{\sqrt{n}} \right)_{t \in [0,1]} \right) \\
 &= d_{\text{SK}}(G_n, (G_{\Lambda_n(t),n})_{t \in [0,1]}) \\
 &\leq \|\tilde{F}_n - \text{id}\|_{\infty}.
 \end{aligned}$$

By (29),  $\tilde{F}_n$  is close to the empirical distribution function and thus converges a.s. uniformly to the identity  $\text{id}$  by the Glivenko–Cantelli theorem. The statement of Corollary 5 then follows from Slutsky's theorem.  $\square$

## 5. FURTHER COST MEASURES

In this section we sketch results analogous to Theorem 1 for other cost measures than the number of key comparisons. We consider the number of swaps required by QuickVal which however depends on the implementation of the procedure to partition the input  $(U_1, \dots, U_n)$  into the sublists  $S_{<}$  and  $S_{\geq}$ . We consider two such procedures, the one originally proposed

by Hoare [10] and one that is attributed to Lomuto, see [1, 3, 13]. Our results are stated in Subsection 5.1.

As a further cost measure we consider the model where the costs to compare two keys may depend on their values, e.g., the number of bit comparisons required to compare them when they are given by their binary expansions. The total cost for all key comparisons required by  $\text{QuickVal}((U_1, \dots, U_n), \alpha)$  or  $\text{FIND}((U_1, \dots, U_n), k)$  is no longer determined by the fact that the ranks of  $(U_1, \dots, U_n)$  form an uniformly random permutation. Here, the distribution of the  $U_i$  matters. We only consider the uniform distribution as in the previous sections and hope to report on other distributions in the full paper version of this extended abstract. Our results are stated in Subsection 5.2. A probabilistic analysis for the number of bit comparisons of the related Quicksort algorithms was given in [6, 5].

**5.1. Number of swaps.** Usually, QuickSelect is implemented in-place, meaning that it only requires the memory for the list  $S$  of values and a bounded amount of additional memory. This is achieved by swapping values within  $S$  so that the elements of  $S_{<}$  and  $S_{\geq}$  are contained in contiguous parts of the list. Such a procedure is called *partition*. There are various procedures of *partition*.

The original procedure by Hoare [10] searches the list  $S$  from both ends at once: It repeatedly finds the index  $i = \min\{2 \leq i \leq n \mid U_i > U_1\}$  of the leftmost element bigger than the pivot and the index  $j := \max\{2 \leq j \leq n \mid U_j < U_1\}$  of the rightmost element smaller than the pivot. If  $i < j$ , it swaps  $U_i$  and  $U_j$ . Else the algorithm terminates.

A simpler, but less efficient implementation is the so-called Lomuto partition scheme [1, 3, 13] that only searches from one end of  $S$ . It keeps track of the amount  $i$  of elements at the start of the list it has already swapped. In every step, it finds the index  $j := \max\{2 \leq j \leq n \mid U_j < U_1\}$  of the rightmost element smaller than the pivot. If  $i + 1 < j$ , it swaps  $U_{i+1}$  and  $U_j$  and increases  $i$  by one. Otherwise the algorithm terminates.

Both partition schemes only compare elements to the pivot, so the model of randomness is preserved within the sublists  $S_{<}$  and  $S_{\geq}$ . However, their original order is not preserved, so QuickSelect run on  $U_1, \dots, U_n$  will usually not select the same pivots as QuickSelect on  $U_1, \dots, U_{n+1}$ . For convenience, we assume that the pivot to split a sublist  $S'$  of  $S$  is the element of  $S'$  that came first in the original list  $S$ . We call this choice of the pivots a *suitable embedding*.

**5.1.1. Hoare's partition.** For Hoare's partition, via a hypergeometric distribution the expected number of swaps in step  $k$  given  $\mathcal{F}_\infty$  is approximately  $nI_{\alpha,k+1}(I_\alpha - I_{\alpha,k})/I_{\alpha,k}$ , which leads to the limit process  $L = (L_\alpha)_{\alpha \in [0,1]}$  given by

$$L_\alpha := \sum_{k=0}^{\infty} \frac{I_{\alpha,k+1}(I_\alpha - I_{\alpha,k})}{I_{\alpha,k}}, \quad \alpha \in [0, 1].$$

It is now possible to study the fluctuations by their contributions from the individual levels and combine them as for the number of key comparisons above. Since we are still in the range of the central limit theorem we again obtain a mixture of centered Gaussian processes. To be explicit, first denote by  $Z_\phi$  the limit of the  $G_{\alpha,k,n}$  as  $n \rightarrow \infty$  where  $\phi = \phi(\alpha, k) \in \{0, 1\}^k$ . Further, denote by  $\{Y_\phi \mid \phi \in \{0, 1\}^*\}$  a set of i.i.d.  $\mathcal{N}(0, 1)$  random variables being independent

of  $\{Z_\phi \mid \phi \in \{0,1\}^*\}$  and of  $\mathcal{F}_\infty$ . Then the limiting process  $G^{\text{swap}} = (G_\alpha^{\text{swap}})_{\alpha \in [0,1]}$  is given by

$$(30) \quad G_\alpha^{\text{swap}} := \sum_{\phi \in \{\phi(k,\alpha) \mid k \in \mathbb{N}_0\}} Y_\phi \frac{I_{\phi 0} I_{\phi 1}}{I_\phi^{3/2}} + Z_{\phi 0} \frac{I_{\phi 1}}{I_\phi} + Z_{\phi 1} \frac{I_{\phi 0}}{I_\phi} - Z_\phi \frac{I_{\phi 0} I_{\phi 1}}{I_\phi}, \quad \alpha \in [0,1].$$

Then we have the following result for key exchanges corresponding to Theorem 1.

**Theorem 15.** *Let  $K_{\alpha,n}$  be the number of key exchanges required by QuickVal $((U_1, \dots, U_n), \alpha)$  with Hoare's partition algorithm in a suitable embedding. Then, as  $n \rightarrow \infty$ , we have*

$$\left( \frac{K_{\alpha,n} - nL_\alpha}{\sqrt{n}} \right)_{\alpha \in [0,1]} \xrightarrow{d} G^{\text{swap}} \quad \text{in } (D[0,1], d_{\text{SK}}).$$

**5.1.2. Lomuto's partition.** The Lomuto partition is simpler to implement and much easier to analyze. The Lomuto partition swaps every element smaller than the pivot, so the amount of swaps at some path  $\phi \in \{0,1\}^*$  is given by  $S_{\phi 0} + 1$ . With the  $Z_\phi$  introduced in Subsection 5.1.1 we find that  $(Z_\phi)_{\phi \in \{0,1\}^*}$  is a mixture of centered Gaussian processes with conditional covariance function given by

$$\text{Cov}(Z_\phi, Z_\psi \mid \mathcal{F}_\infty) = I_{\phi \vee \psi} - I_\phi I_\psi, \quad \phi, \psi \in \{0,1\}^*,$$

where  $I_{\phi \vee \psi}$  is the length of the interval  $[L_\phi, R_\phi) \cap [L_\psi, R_\psi)$ , thus  $I_{\phi \vee \psi}$  is only nonzero if one of  $\psi$  and  $\phi$  is a prefix of the other. Then the limiting process  $G^{\text{Lo}} = (G_\alpha^{\text{Lo}})_{\alpha \in [0,1]}$  is given by

$$G_\alpha^{\text{Lo}} = \sum_{k=0}^{\infty} Z_{\phi(\alpha,k)0}, \quad \alpha \in [0,1].$$

We can directly apply Lemma 13 and our proof for the number of key comparisons can be straightforwardly transferred.

**Theorem 16.** *Let  $K_{\alpha,n}^{\text{Lo}}$  be the number of key exchanges required by QuickVal $((U_1, \dots, U_n), \alpha)$  with Lomuto's partition procedure in a suitable embedding. Then, as  $n \rightarrow \infty$ , we have*

$$\left( \frac{K_{\alpha,n}^{\text{Lo}} - n \sum_{k=0}^{\infty} I_{\phi(\alpha,k)0}}{\sqrt{n}} \right)_{\alpha \in [0,1]} \xrightarrow{d} G^{\text{Lo}} \quad \text{in } (D[0,1], d_{\text{SK}}).$$

**5.2. Number of bit comparisons.** We now consider the model where the cost to compare two keys depends on their values. These costs are described by a measurable *cost function*  $\beta : [0,1]^2 \rightarrow [0,\infty)$ , and we require that they have a polynomial tail, that is: There are constants  $c, \varepsilon > 0$  such that for all  $u \in [0,1], x \in \mathbb{N}$  and for  $V \sim \text{unif}[0,1]$

$$\mathbb{P}(\beta(u, V) \geq x) \leq cx^{-1/\varepsilon}.$$

This condition is called  $(c, \varepsilon)$ -*tameness*, see Matterer [14], and  $\beta$  is called to be  $\varepsilon$ -*tame* if it is  $(c, \varepsilon)$ -tame for some  $c > 0$ . Note that, e.g.,  $\beta$  counting the number of bit comparisons is  $\varepsilon$ -tame for all  $\varepsilon > 0$ . The costs of QuickVal $((U_1, \dots, U_n), \alpha)$  in this model are given by

$$S_{\alpha,n}^\beta := \sum_{k=0}^{\infty} \sum_{\tau_{\alpha,k} < i \leq n} \mathbf{1}_{[L_{\alpha,k}, R_{\alpha,k})}(U_i) \beta(U_{\tau_{\alpha,k}}, U_i)$$

and the limit is, with  $V \sim \text{unif}[0,1]$  being independent of the  $U_1, \dots, U_n$ , given as

$$S_{\alpha,\infty}^\beta := \sum_{k=0}^{\infty} \mathbb{E} \left[ \mathbf{1}_{[L_{\alpha,k}, R_{\alpha,k})}(V) \beta(U_{\tau_{\alpha,k}}, V) \mid \mathcal{F}_\infty \right]$$

Matterer [14, Theorem 6.4 and Theorem 6.14] shows for  $\varepsilon < \frac{1}{2}$  that for fixed  $\alpha \in [0, 1]$  the resulting residual

$$G_{\alpha,n}^\beta := \frac{S_{\alpha,n}^\beta - nS_{\alpha,\infty}^\beta}{\sqrt{n}}$$

converges to a mixed centered Gaussian random variable  $G_{\alpha,\infty}^\beta$  in distribution and with all moments. It is possible to combine them to a mixture of centered Gaussian processes

$$(31) \quad G_\infty^\beta = (G_{\alpha,\infty}^\beta)_{\alpha \in [0,1]},$$

defined by the conditional covariance functions given, with  $X_{\alpha,k}^\beta := \mathbf{1}_{[L_{\alpha,k}, R_{\alpha,k})}(V) \cdot \beta(U_{\tau_{\alpha,k,n}}, V)$ , by

$$(32) \quad \text{Cov}(G_{\alpha,\infty}^\beta, G_{\gamma,\infty}^\beta \mid \mathcal{F}_\infty) = \text{Cov}\left(\sum_{k=0}^{\infty} X_{\alpha,k}^\beta, \sum_{k=0}^{\infty} X_{\gamma,k}^\beta \mid \mathcal{F}_\infty\right), \quad \alpha, \gamma \in [0, 1].$$

It can be shown that the latter expression is a.s. well-defined. We have the following result corresponding to Theorem 1.

**Theorem 17.** *Let  $\beta$  be an  $\varepsilon$ -tame cost function with  $\varepsilon < \frac{1}{4}$ . Then we have*

$$\left( \frac{S_{\alpha,n}^\beta - nS_{\alpha,\infty}^\beta}{\sqrt{n}} \right)_{\alpha \in [0,1]} \xrightarrow{d} G_\infty^\beta \quad \text{in } (D[0, 1], d_{\text{SK}}),$$

where  $G_\infty^\beta$  is the mixture of centered Gaussian processes defined in (31).

## REFERENCES

- [1] Jon Bentley, *Programming pearls*, Addison-Wesley, Reading, Mass., 2000.
- [2] Patrick Billingsley, *Convergence of probability measures*, Wiley Series in Probability and Statistics: Probability and Statistics, John Wiley & Sons Inc., New York, 1999, A Wiley-Interscience Publication. MR MR1700749 (2000e:60008)
- [3] Thomas H. Cormen, Charles E. Leiserson, Ronald L. Rivest, and Clifford Stein, *Introduction to algorithms*, third ed., MIT Press, Cambridge, MA, 2009. MR 2572804
- [4] Luc Devroye, *A note on the height of binary search trees*, J. ACM **33** (1986), no. 3, 489–498.
- [5] James Allen Fill, *Distributional convergence for the number of symbol comparisons used by QuickSort*, Ann. Appl. Probab. **23** (2013), no. 3, 1129–1147. MR 3076680
- [6] James Allen Fill and Svante Janson, *The number of bit comparisons used by Quicksort: an average-case analysis*, Electron. J. Probab. **17** (2012), no. 43, 22. MR 2928726
- [7] James Allen Fill and Jason Matterer, *QuickSelect tree process convergence, with an application to distributional convergence for the number of symbol comparisons used by worst-case find*, Combin. Probab. Comput. **23** (2014), no. 5, 805–828. MR 3249225
- [8] James Allen Fill and Takehiko Nakama, *Distributional convergence for the number of symbol comparisons used by quickselect*, Advances in Applied Probability **45** (2013), no. 2, 425–450.
- [9] Rudolf Grübel and Uwe Rösler, *Asymptotic distribution theory for Hoare's selection algorithm*, Adv. in Appl. Probab. **28** (1996), no. 1, 252–269. MR 1372338
- [10] C. A. R. Hoare, *Algorithm 63: Partition*, Commun. ACM **4** (1961), no. 7, 321.
- [11] ———, *Algorithm 65: Find*, Commun. ACM **4** (1961), no. 7, 321–322.
- [12] Kevin Leckey, Ralph Neininger, and Henning Sulzbach, *Process convergence for the complexity of radix selection on markov sources*, Stochastic Processes and their Applications **129** (2019), no. 2, 507–538.
- [13] Hosam M. Mahmoud, *Distributional analysis of swaps in quick select*, Theoretical Computer Science **411** (2010), no. 16, 1763–1769.
- [14] Jason Matterer, *Quickselect process and quickval residual convergence*, Ph.D. thesis, The Johns Hopkins University, Baltimore, Maryland, 2015.

- [15] Colin McDiarmid, *Concentration*, Probabilistic Methods for Algorithmic Discrete Mathematics (Michel Habib, Colin McDiarmid, Jorge Ramirez-Alfonsin, and Bruce Reed, eds.), Springer Berlin Heidelberg, Berlin, Heidelberg, 1998, pp. 195–248.
- [16] Bruce Reed, *The height of a random binary search tree*, J. ACM **50** (2003), no. 3, 306–332.
- [17] Henning Sulzbach, Ralph Neininger, and Michael Drmota, *A Gaussian limit process for optimal FIND algorithms*, Electronic Journal of Probability **19** (2014), no. none, 1 – 28.
- [18] Brigitte Vallée, Julien Clément, James Allen Fill, and Philippe Flajolet, *The number of symbol comparisons in QuickSort and QuickSelect*, Automata, languages and programming. Part I, Lecture Notes in Comput. Sci., vol. 5555, Springer, Berlin, 2009, pp. 750–763. MR 2544890

INSTITUTE OF MATHEMATICS, GOETHE UNIVERSITY FRANKFURT, FRANKFURT A.M., GERMANY  
*Email address:* `ischebec@math.uni-frankfurt.de`

INSTITUTE OF MATHEMATICS, GOETHE UNIVERSITY FRANKFURT, FRANKFURT A.M., GERMANY  
*Email address:* `neininger@math.uni-frankfurt.de`