

DISTRIBUTIONAL CONVERGENCE FOR THE NUMBER OF SYMBOL COMPARISONS USED BY QUICKSORT

JAMES ALLEN FILL

ABSTRACT

Most previous studies of the sorting algorithm `QuickSort` have used the number of key comparisons as a measure of the cost of executing the algorithm. Here we suppose that the n independent and identically distributed (iid) keys are each represented as a sequence of symbols from a probabilistic source and that `QuickSort` operates on individual symbols, and we measure the execution cost as the number of symbol comparisons. Assuming only a mild “tameness” condition on the source, we show that there is a limiting distribution for the number of symbol comparisons after normalization: first centering by the mean and then dividing by n . Additionally, under a condition that grows more restrictive as p increases, we have convergence of moments of orders p and smaller. In particular, we have convergence in distribution and convergence of moments of every order whenever the source is memoryless, i.e., whenever each key is generated as an infinite string of iid symbols. This is somewhat surprising: Even for the classical model that each key is an iid string of unbiased (“fair”) bits, the mean exhibits periodic fluctuations of order n .

1. INTRODUCTION, REVIEW OF RELATED LITERATURE, AND SUMMARY

1.1. Introduction. We consider Hoare’s [11] `QuickSort` algorithm applied to n distinct random keys X_1, \dots, X_n , each represented as a word (i.e., infinite string of symbols such as bits) from some specified finite or countably infinite alphabet. We will consider various probabilistic mechanisms [called (*probabilistic*) *sources*] for generating the symbols within a key, but we will always assume that the keys themselves are iid (independent and identically distributed), and we will later place conditions on the source that rule out the generation of equal keys.

`QuickSort`(X_1, \dots, X_n) chooses one of the n keys X_1, \dots, X_n (called the “pivot”) uniformly at random, compares each of the other keys to it, and then proceeds recursively to sort both the keys smaller than the pivot and those larger than it.

Key observation (coupling): Because of the assumption that the keys are iid, we may take the pivot to be the *first* key in the sequence, X_1 . Thus if X_1, X_2, \dots is an infinite sequence of keys and C_n is any measure of the cost of sorting n random keys using any cost function c (for example, the number of key comparisons or the number of symbol comparisons), then we can place all the random variables C_n on a common probability space by using $C_n = c(X_1, \dots, X_n)$. Notice that then C_n is nondecreasing in n . We will assume throughout that this natural coupling of the random variables C_n has been used. The coupling opens up the possibility of

Date: Revised March 28, 2010.

Research supported by the Acheson J. Duncan Fund for the Advancement of Research in Statistics.

establishing stronger forms of convergence than convergence in distribution, such as almost sure convergence and convergence in L^p , for suitably normalized C_n .

Many authors (Knuth [15], Régnier [18], Rösler [20], Knessl and Szpankowski [14], Fill and Janson [5] [6], Neininger and Ruschendorf [17], and others) have studied K_n , the (random) number of key comparisons performed by the algorithm. This is an appropriate measure of the cost of the algorithm if each comparison has the same cost. On the other hand, if keys are represented as words and comparisons are done by scanning the words from left to right, comparing the symbols of matching index one by one, then the cost of comparing two keys is determined by the number of symbols compared until a difference is found. We call this number the number of *symbol comparisons* for the key comparison, and let S_n denote the total number of symbol comparisons when n keys are sorted by `QuickSort`. Symbol-complexity analysis allows us to compare key-based algorithms such as `QuickSort` with digital algorithms such as those utilizing digital search trees.

The goal of the present work is to establish a limiting distribution for the normalized sequence of random variables $(S_n - \mathbf{E} S_n)/n$. Both exact and limiting distributions of S_n will depend on the source, unlike for K_n .

1.2. Review of closely related literature (QuickSort and QuickSelect). Until now, study of asymptotics for `QuickSort`'s S_n has been limited mainly to the expected value $\mathbf{E} S_n$. Fill and Janson [7] were the pioneers in that regard, obtaining, *inter alia*, exact and asymptotic expressions for $\mathbf{E} S_n$ [consult their Theorem 1.1, and note that the asymptotic expansion extends through terms of order n with a $O(\log n)$ remainder] when the keys are infinite binary strings and the bits within a key result from iid fair coin tosses. (We will refer to this model for key-generation as “the standard binary source”. Equivalently, a key is generated by sampling uniformly from the unit interval, representing the result in binary notation, and dropping the leading “binary point”.) They found that the expected number of bit comparisons required by `QuickSort` to sort n keys is asymptotically equivalent to $\frac{1}{\ln 2} n \ln^2 n$, whereas the lead-order term of the expected number of *key* comparisons is $2n \ln n$, smaller by a factor of order $\log n$. Now suppose that $N = (N(t) : 0 \leq t < \infty)$ is a Poisson process with rate 1 and is independent of the generation of the keys, and let $S(t) := S_{N(t)}$. The authors also found for each fixed $1 \leq p < \infty$ an upper bound independent of $t \geq 1$ on the L^p -norm of

$$(1.1) \quad Y(t) := \frac{S(t) - \mathbf{E} S(t)}{t}$$

[see their Remark 5.1(a)], leading them to speculate that $Y(t)$ might have a limiting distribution as $t \rightarrow \infty$. We will see that a limiting distribution does indeed exist, not only for the standard binary source but for a wide range of sources, as well.

Vallée, Clément, Fill, and Flajolet [22] greatly extended the scope of [7] by establishing for much more general sources both an exact expression for $\mathbf{E} S_n$ [consult their Proposition 3 and display (8)] and an asymptotic expansion (see their Theorem 1) through terms of order n with a $o(n)$ remainder. For the broad class of sources \mathcal{S} considered, the expected number of symbol comparisons is of lead order $\frac{1}{h(\mathcal{S})} n \ln^2 n$, where $h(\mathcal{S})$ is the entropy of the source (see their Figure 1 for a definition).

Building on work of Fill and Nakama [3], who had in turn followed closely along the lines of [7], Vallée et al. [22] also studied the expected number of symbol

comparisons required by the algorithm `QuickSelect`(n, m). This algorithm [aka `Find`(n, m)], a close cousin of `QuickSort` also devised by Hoare [12], finds a key of specified rank m from a list of n keys. The authors of [22] considered the case where $m = \alpha n + o(n)$ for general $\alpha \in [0, 1]$ [note: we will sometimes refer to `QuickQuant`(n, α), rather than `QuickSelect`(n, m), in this case] and a broad class of sources \mathcal{S} . They found that the expected number of symbol comparisons asymptotically has lead term $\rho_{\mathcal{S}}(\alpha)n$, where $\rho_{\mathcal{S}}(\alpha)$ is described in their Figure 1. Unlike in the case of `QuickSort`, this is only a constant times larger than the expected number of key comparisons, which is well known to be asymptotically $\kappa(\alpha)n$ with

$$\kappa(\alpha) := 2[1 - \alpha \ln \alpha - (1 - \alpha) \ln(1 - \alpha)].$$

For either `QuickSelect` or `QuickSort`, a deeper probabilistic analysis of the numbers of key comparisons and symbol comparisons is obtained by treating entire distributions and not just expectations—in particular, by finding limiting distributions for suitable normalizations of these counts and, if possible, establishing corresponding convergence of moments. Consider `QuickQuant`(n, α) first. For both key comparisons and symbol comparisons a suitable normalization is to divide by n , with no need to center first. For a literature review on the number of key comparisons, we refer the reader to [9, Section 2.2]; the number of symbol comparisons is discussed next.

Fill and Nakama [9] (see also [16]) were the first to establish a limiting distribution for the number of symbol comparisons for any sorting or searching algorithm. They considered `QuickQuant`(n, α) for a broad class of sources and found a limiting distribution (depending on α , and of course also on the source) for the number $S_n(\alpha)$ of symbol comparisons (after division by n). It would take us a bit too far afield to describe the limiting random variable $S(\alpha)$, so we refer the reader to [9, Section 3.1, see (3.7)] for an explicit description. In their paper they use the natural coupling discussed in Section 1.1 and prove, for each α , that $S_n(\alpha)/n$ converges to $S(\alpha)$ both (i) almost surely and, under ever stronger conditions on the source as p increases, (ii) in L^p . Either conclusion implies convergence in distribution, and (ii) implies convergence of moments of order $\leq p$. The approach taken in [9] is sufficiently general that the authors were able to unify treatment of key comparisons and symbol comparisons and to consider various other cost functions: see their Example 2.1.

Now we turn our attention back to `QuickSort`, the focus of this paper. Let K_n (respectively, S_n) denote the random number of key (resp., symbol) comparisons required by `QuickSort` to sort a list of n keys. We first consider K_n , for which we know the following convergence in law, for some random variable T (where the immaterial choice of scaling by $n + 1$, rather than n , matches with [18]):

$$(1.2) \quad \frac{K_n - \mathbf{E} K_n}{n + 1} \xrightarrow{\mathcal{L}} T.$$

This was proved (i) by Régnier [18], who used the natural coupling and martingale techniques to establish convergence both almost surely and in L^p for every finite p ; and (ii) by Rösler [20], who used the *contraction method* (see Rösler and Rüschendorf [21] for a general discussion) to prove convergence in the so-called *minimal L^p metric* for every finite p [from which (1.2), with convergence of all moments, again follows]. An advantage of Rösler’s approach was identification of the distribution of the limiting T as the unique distribution of a zero-mean random

variable with finite variance satisfying the distributional fixed-point equation

$$(1.3) \quad T \stackrel{\mathcal{L}}{=} UT + (1 - U)T^* + g(U),$$

with $g(u) := 1 + 2u \ln u + 2(1 - u) \ln(1 - u)$ and where, on the right, T , T^* , and U are independent random variables; T^* has the same distribution as T ; and U is distributed uniformly over $(0, 1)$. Later, Fill and Janson [4] showed that uniqueness of the zero-mean solution $\mathcal{L}(T)$ to (1.3) continues to hold without the assumption of finite variance, or indeed any other assumption.

1.3. Summary. This paper establishes, for a broad class of sources, a limiting distribution for the number S_n of symbol comparisons for `QuickSort`. We tried without success to mimic the approach used in [9] for `QuickQuant`. The approach used in this paper, very broadly put, is to relate the count S_n of symbol comparisons to various counts of key comparisons and then rely (heavily) on the result of Régnier [18]. Like Fill and Janson [7], we will find it much more convenient to work in continuous time than in discrete time. (We hope to “de-Poissonize” our result in the full-length paper.) In the continuous-time setting and notation established at (1.1) (but without limiting attention to the standard binary source), we will prove in this extended abstract, assuming that the source is suitably “tame” (in a sense to be made precise), that

$$(1.4) \quad Y(t) = \frac{S(t) - \mathbf{E}S(t)}{t} \xrightarrow{\mathcal{L}} Y$$

for some random variable Y . Following the lead of [18] and [9], we will use the natural coupling discussed in Section 1.1. Under a mild tameness condition that becomes more stringent as $p \in [2, \infty)$ increases we will in fact establish convergence in L^p . In particular, for any g -tamed source as defined in Remark 2.3(a)—for example, for any (nondegenerate) memoryless source—we have convergence in L^p for every finite p .

Outline of the paper. After carefully describing in Section 2.1 the probabilistic models used to govern the generation of keys, reviewing in Section 2.2 four known results about the number of key comparisons we will need in our analysis of symbol comparisons, and listing in Section 2.3 the other basic probability tools we will need, in Section 3 we state and prove our main result about convergence in distribution for the number of symbol comparisons.

2. BACKGROUND AND PRELIMINARIES

2.1. Probabilistic source models for the keys. In this subsection, extracted with only small modifications from [9], we describe what is meant by a probabilistic source—our model for how the iid keys are generated—using the terminology and notation of Vallée et al. [22].

Let Σ denote a totally ordered alphabet (set of symbols), assumed to be isomorphic either to $\{0, \dots, r - 1\}$ for some finite r or to the full set of nonnegative integers, in either case with the natural order; a *word* is then an element of Σ^∞ , i.e., an infinite sequence (or “string”) of symbols. We will follow the customary practice of denoting a word $w = (w_1, w_2, \dots)$ more simply by $w_1 w_2 \dots$.

We will use the word “prefix” in two closely related ways. First, the symbol strings belonging to Σ^k are called *prefixes* of length k , and so $\Sigma^* := \cup_{0 \leq k < \infty} \Sigma^k$

denotes the set of all prefixes of any nonnegative finite length. Second, if $w = w_1w_2\cdots$ is a word, then we will call

$$(2.1) \quad w(k) := w_1w_2\cdots w_k \in \Sigma^k$$

its *prefix of length k* .

Lexicographic order is the linear order (to be denoted in the strict sense by \prec) on the set of words specified by declaring that $w \prec w'$ if (and only if) for some $0 \leq k < \infty$ the prefixes of w and w' of length k are equal but $w_{k+1} < w'_{k+1}$. Then the symbol-comparisons cost of determining $w \prec w'$ for such words is just $k + 1$, the number of symbol comparisons.

A *probabilistic source* is simply a stochastic process $W = W_1W_2\cdots$ with state space Σ (endowed with its total σ -field) or, equivalently, a random variable W taking values in Σ^∞ (with the product σ -field). According to Kolmogorov's consistency criterion (e.g., [1, Theorem 3.3.6]), the distributions μ of such processes are in one-to-one correspondence with consistent specifications of finite-dimensional marginals, that is, of the probabilities

$$p_w := \mu(\{w_1\cdots w_k\} \times \Sigma^\infty), \quad w = w_1w_2\cdots w_k \in \Sigma^*.$$

Here the *fundamental probability* p_w is the probability that a word drawn from μ has $w_1\cdots w_k$ as its length- k prefix.

Because the analysis of **QuickSort** is significantly more complicated when its input keys are not all distinct, we will restrict attention to probabilistic sources with continuous distributions μ . Expressed equivalently in terms of fundamental probabilities, our continuity assumption is that for any $w = w_1w_2\cdots \in \Sigma^\infty$ we have $p_{w(k)} \rightarrow 0$ as $k \rightarrow \infty$, recalling the prefix notation (2.1).

Example 2.1. We present a few classical examples of sources. For more examples, and for further discussion, see Section 3 of [22].

(a) In computer science jargon, a *memoryless source* is one with W_1, W_2, \dots iid. Then the fundamental probabilities p_w have the product form

$$p_w = p_{w_1}p_{w_2}\cdots p_{w_k}, \quad w = w_1w_2\cdots w_k \in \Sigma^*.$$

(b) A *Markov source* is one for which $W_1W_2\cdots$ is a Markov chain.

(c) An intermittent source (a model for long-range dependence) over the finite alphabet $\Sigma = \{0, \dots, r-1\}$ is defined by specifying the conditional distributions $\mathcal{L}(W_j | W_1, \dots, W_{j-1})$ ($j \geq 2$) in a way that pays special attention to a particular symbol $\underline{\sigma}$. The source is said to be *intermittent of exponent $\gamma > 0$ with respect to $\underline{\sigma}$* if $\mathcal{L}(W_j | W_1, \dots, W_{j-1})$ depends only on the maximum value k such that the last k symbols in the prefix $W_1\cdots W_{j-1}$ are all $\underline{\sigma}$ and (i) is the uniform distribution on Σ , if $k = 0$; and (ii) if $1 \leq k \leq j-1$, assigns mass $[k/(k+1)]^\gamma$ to $\underline{\sigma}$ and distributes the remaining mass uniformly over the remaining elements of Σ .

For our results, the quantity

$$(2.2) \quad \pi_k := \max\{p_w : w \in \Sigma^k\}$$

will play an important role, as it did in [22, eqn. (7)] in connection with the generalized Dirichlet series $\Pi(s) := \sum_{k \geq 0} \pi_k^{-s}$. In particular, it will be sufficient to obtain L^p convergence in our main result (Theorem 3.1) that

$$(2.3) \quad \Pi\left(-\frac{1}{p}\right) = \sum_{k \geq 0} \pi_k^{1/p} < \infty;$$

a sufficient condition for this, in turn, is of course that the source is Π -tamed with $\gamma > p$ in the sense of the following definition:

Definition 2.2. Let $0 < \gamma < \infty$ and $0 < A < \infty$. We say that the source is Π -tamed (with parameters γ and A) if the sequence (π_k) at (2.2) satisfies

$$\pi_k \leq A(k+1)^{-\gamma} \text{ for every } k \geq 0.$$

Observe that a Π -tamed source is always continuous.

Remark 2.3. (a) Many common sources have geometric decrease in π_k (call these “g-tamed”) and so for *any* γ are Π -tamed with parameters γ and A for suitably chosen $A \equiv A_\gamma$.

For example, a memoryless source satisfies $\pi_k = p_{\max}^k$, where

$$p_{\max} := \sup_{w \in \Sigma^1} p_w$$

satisfies $p_{\max} < 1$ except in the highly degenerate case of an essentially single-symbol alphabet. We also have $\pi_k \leq p_{\max}^k$ for any Markov source, where now p_{\max} is the supremum of all one-step transition probabilities, and so such a source is g-tamed provided $p_{\max} < 1$. Expanding dynamical sources (cf. [2]) are also g-tamed.

(b) For an intermittent source as in Example 2.1, for all large k the maximum probability π_k is attained by the prefix $\underline{\sigma}^k$ and equals

$$\pi_k = r^{-1}k^{-\gamma}.$$

Intermittent sources are therefore examples of Π -tamed sources for which π_k decays at a truly inverse-polynomial rate, not an exponential rate as in the case of g-tamed sources.

2.2. Known results for the numbers of key comparisons for QuickSort. In this subsection we review four known QuickSort key-comparisons results—the first two formulated in discrete time and the next two in continuous time—that will be useful in proving our main Theorem 3.1. The first gives exact and asymptotic formulas for the expected number of key comparisons in discrete time and is extremely basic and well known. (See, e.g., (2.1)–(2.2) in [7].)

Lemma 2.4. Let K_n denote the number of key comparisons required to sort a list of n distinct keys. Then

$$\begin{aligned} \mathbf{E} K_n &= 2(n+1)H_n - 4n \\ (2.4) \quad &= 2n \ln n - (4 - 2\gamma)n + 2 \ln n + (2\gamma + 1) + O(1/n). \end{aligned}$$

The second result—mentioned previously at (1.2)—is due to Régnier [18], who also proved convergence in L^p for every finite p . Recall the *natural coupling* discussed in Section 1.1.

Lemma 2.5 ([18]). Under the natural coupling, there exists a random variable T satisfying

$$(2.5) \quad \frac{K_n - \mathbf{E} K_n}{n+1} \rightarrow T \text{ almost surely.}$$

We now shift to continuous time by assuming that the successive keys are generated at the arrival times of a Poisson process with unit rate. The number of key comparisons through epoch t is then $K_{N(t)}$, which we will abbreviate as $K(t)$; while

the sequence (K_n) is thereby naturally embedded in the continuous-time process, the random variables $K(n)$ and K_n are not to be confused. We will use such abbreviations throughout this extended abstract; for example, we will also write $S_{N(t)}$ as $S(t)$.

The third result we review is the continuous-time analogue of Lemma 2.4. Note the difference in constant terms and the much smaller error term in continuous time.

Lemma 2.6 ([7, Lemma 5.1]; proved in [8]). *In the continuous-time setting, the expected number of key comparisons is given by*

$$\mathbf{E} K(t) = 2 \int_0^t (t-y)(e^{-y} - 1 + y)y^{-2} dy.$$

Asymptotically, as $t \rightarrow \infty$ we have

$$(2.6) \quad \mathbf{E} K(t) = 2t \ln t - (4 - 2\gamma)t + 2 \ln t + (2\gamma + 2) + O(e^{-t} t^{-2}).$$

The fourth result gives bounds on the moments of $K(t)$. For real $p \in [1, \infty)$, we let $\|W\|_p := (\mathbf{E} |W|^p)^{1/p}$ denote L^p -norm.

Lemma 2.7 ([7, Lemma 5.2]; proved in [8]). *For every real $p \in [1, \infty)$, there exists a constant $c_p < \infty$ such that*

$$\begin{aligned} \|K(t) - \mathbf{E} K(t)\|_p &\leq c_p t && \text{for } t \geq 1, \\ \|K(t)\|_p &\leq c_p t^{2/p} && \text{for } t \leq 1. \end{aligned}$$

In the special case $p = 2$, it follows immediately from Lemma 2.7 that

$$(2.7) \quad \mathbf{Var} K(t) \leq c_2^2 t^2 \quad \text{for } 0 \leq t < \infty.$$

2.3. Basic probability tools. The following elementary lemma is the basic tool we will use for L^p -convergence. For completeness and the reader's convenience, we supply a proof.

Lemma 2.8. *Let $Y_k(t)$ be random variables, all defined on a common probability space, for $k = 0, 1, 2, \dots$ and $0 \leq t \leq \infty$. Fix $t_0 \in [0, \infty)$ and $1 \leq p < p' < \infty$ and suppose for some sequences (b_k) and (b'_k) that*

- (i) *for each k we have $Y_k(t) \rightarrow Y_k(\infty)$ almost surely as $t \rightarrow \infty$,*
- (ii) *for each k we have $\|Y_k(t)\|_p \leq b_k$ for all $t_0 \leq t < \infty$,*
- (ii') *for each k we have $\|Y_k(t)\|_{p'} \leq b'_k < \infty$ for all $t_0 \leq t < \infty$, and*
- (iii) $\sum_{k=0}^{\infty} b_k < \infty$.

Then

- (a) *for each $t_0 \leq t \leq \infty$ the series $\sum_{k=0}^{\infty} Y_k(t)$ converges in L^p to some random variable $Y(t)$, and moreover*
- (b) $Y(t) \rightarrow Y(\infty)$ in L^p as $t \rightarrow \infty$.

Proof. We assume without loss of generality that $t_0 = 0$. Note that hypotheses (ii) and (ii') extend to $t = \infty$ by Fatou's lemma.

(a) From (ii) and (iii) it follows for each $0 \leq t \leq \infty$ that the sequence of partial sums $\sum_{k=0}^K Y_k(t)$, $K = 0, 1, \dots$, is a Cauchy sequence in the Banach space L^p and so converges to some random variable $Y(t)$.

(b) We first claim for each k that $Y_k(t) \rightarrow Y_k(\infty)$ in L^p , i.e., $|Y_k(t) - Y_k(\infty)|^p \rightarrow 0$ in L^1 as $t \rightarrow \infty$. To see this, from (ii') it follows using [1, Exercise 4.5.8] that $|Y_k(t)|^p$ is uniformly integrable in t , as therefore is $|Y_k(t) - Y_k(\infty)|^p$. Our claim then follows from (i), since almost-sure convergence to 0 implies convergence in probability to 0, and that together with uniform integrability implies convergence in L^1 (e.g., [1, Theorem 4.5.4]).

Using the triangle inequality for L^p -norm, the claim proved in the preceding paragraph, and the extended condition (ii), we find for any K that

$$\limsup_{t \rightarrow \infty} \|Y(t) - Y(\infty)\|_p \leq \limsup_{t \rightarrow \infty} \sum_{k=K+1}^{\infty} \|Y_k(t) - Y_k(\infty)\|_p \leq 2 \sum_{k=K+1}^{\infty} b_k.$$

Now let $K \rightarrow \infty$, using (iii), to complete the proof. \square

Later (Lemma 3.3) we will transfer Lemma 2.5 to continuous time. When we do so, the following result will prove useful. This *law of the iterated logarithm* (LIL) is well known, and for example can be found for general renewal processes in [13, Theorem 12.13].

Lemma 2.9 (LIL for a Poisson process). *Abbreviate the natural logarithm function as L . For a Poisson process N with unit rate,*

$$(2.8) \quad \mathbf{P} \left(\limsup_{t \rightarrow \infty} \frac{N(t) - t}{\sqrt{2tLLt}} = 1, \quad \liminf_{t \rightarrow \infty} \frac{N(t) - t}{\sqrt{2tLLt}} = -1 \right) = 1.$$

3. MAIN RESULTS

3.1. Convergence in L^p (and therefore in distribution). The following theorem, which adopts the natural coupling discussed in Section 1.1 and utilizes the terminology and notation of Section 2.1 for probabilistic sources, is our main result.

Theorem 3.1. *Consider the continuous-time setting in which keys are generated from a probabilistic source at the arrival times of a Poisson process N with unit rate. Let $S(t) = S_{N(t)}$ denote the number of symbol comparisons required by QuickSort to sort the keys generated through epoch t , and let*

$$(3.1) \quad Y(t) := \frac{S(t) - \mathbf{E}S(t)}{t}, \quad 0 < t < \infty.$$

Let $p \in [2, \infty)$ and assume that

$$(3.2) \quad \sum_{k=0}^{\infty} \left(\sum_{w \in \Sigma^k} p_w^2 \right)^{1/p} < \infty.$$

Then there exists a random variable Y such that $Y(t) \rightarrow Y$ in L^p . Thus $Y(t) \xrightarrow{\mathcal{L}} Y$, with convergence of moments of orders $\leq p$; in particular (because $\mathbf{E}Y(t) \rightarrow \mathbf{E}Y$), we have $\mathbf{E}Y = 0$.

Remark 3.2. (a) Observe that $\sum_{w \in \Sigma^k} p_w = 1$ for each k . Thus $\sum_{w \in \Sigma^k} p_w^2 \leq 1$, and condition (3.2) grows increasingly stronger as p increases.

(b) Under the weakest instance $p = 2$ of the assumption (3.2) we have $Y(t) \rightarrow Y$ in L^2 , and so $Y(t) \rightarrow Y$ in law with convergence of means and variances. The random variable Y in Theorem 3.1 of course does not (more precisely, can be taken not to) depend on the value of p considered (because a limit in L^p for any p is also a limit in probability, and limits in probability are almost surely unique).

(c) The expected number of symbol comparisons in comparing two independent keys generated by the given source is $\sum_{w \in \Sigma^*} p_w^2 = \sum_{k=0}^{\infty} \sum_{w \in \Sigma^k} p_w^2$. So (3.2) is certainly sufficient to imply that $\mathbf{E}S(t) < \infty$ for every t [in fact, it follows from calculations to be performed in the proof of Theorem 3.1 for $p = 2$ that $\mathbf{E}S^2(t) < \infty$ for every t] and that with probability one $S(t) < \infty$ for all t .

(d) The sum on w in (3.2) is bounded above by the max-prefix probability π_k defined at (2.2), and so (2.3) (namely, $\sum_k \pi_k^{1/p} < \infty$) is sufficient for (3.2). Thus from the discussion in Section 2.1 we see that Theorem 3.1 gives L^p -convergence for $Y(t)$ for all Π -tamed sources with parameter $\gamma > p$. In particular, for any g -tamed source, such as any (nondegenerate) memoryless source, we have $Y(t) \rightarrow Y$ in L^p for every $p < \infty$.

(e) The standard binary source is a classical example of a periodic memoryless source (cf. [22]—specifically, Definition 3(d), Theorem 1(ii), and the discussion (ii) in Section 3). Fill and Janson [7, eqn. (1.3)] (proved as Proposition 5.4 in [8]) show explicitly for the standard binary source that

$$\mathbf{E}S(t) = \frac{1}{\ln 2} t \ln^2 t - c_1 t \ln t + c_2 t + \pi_t t + O(\log t) \text{ as } t \rightarrow \infty,$$

where c_1, c_2 are explicitly given constants and π_t is a certain periodic function of $\log t$. Given the periodic term of order t in the mean for this periodic source, we find it surprising that Theorem 3.1 nevertheless applies.

(f) We wonder (but have not yet considered): Under what conditions do we have $Y(t) \rightarrow Y$ almost surely?

To prepare for the proof of Theorem 3.1, we “Poissonize” Lemma 2.5.

Lemma 3.3. *In the continuous-time setting of Theorem 3.1, let $K(t) = K_{N(t)}$ denote the number of key comparisons required by QuickSort. Then for the same random variable T as in the discrete-time Lemma 2.5 we have*

$$\frac{K(t) - \mathbf{E}K(t)}{t} \rightarrow T \text{ almost surely as } t \rightarrow \infty.$$

Proof. This is routine. According to Lemmas 2.5 and 2.4,

$$\frac{K_n - [2n \ln n - (4 - 2\gamma)n]}{n + 1} \rightarrow T \text{ almost surely as } n \rightarrow \infty.$$

Since $N(t) \rightarrow \infty$ almost surely as $t \rightarrow \infty$, it follows that

$$\frac{K(t) - [2N(t) \ln N(t) - (4 - 2\gamma)N(t)]}{N(t) + 1} \rightarrow T \text{ almost surely as } t \rightarrow \infty.$$

Using the strong law of large numbers (SLLN) for N [namely, $N(t)/t \rightarrow 1$ almost surely, for which Lemma 2.9 is plenty sufficient], we deduce

$$\frac{K(t) - [2N(t) \ln N(t) - (4 - 2\gamma)t]}{t} \rightarrow T \text{ almost surely as } t \rightarrow \infty.$$

From the mean value theorem it follows that $|y \ln y - x \ln x| \leq |y - x|(1 + \ln x + \ln y)$ for $x, y \geq 1$. Applying this with $x = t$ and $y = N(t)$ and invoking the SLLN and the LIL (Lemma 2.9), we find almost surely that for large t we have

$$\begin{aligned} |N(t) \ln N(t) - t \ln t| &\leq |N(t) - t| [1 + \ln N(t) + \ln t] \leq \sqrt{3t \ln \ln t} [2 \ln t + 1 + o(1)] \\ &= O(\sqrt{t \ln \ln t} \times \ln t) = o(t), \end{aligned}$$

and so

$$\frac{K(t) - [2t \ln t - (4 - 2\gamma)t]}{t} \rightarrow T \text{ almost surely as } t \rightarrow \infty.$$

The desired result now follows from (2.6) in Lemma 2.6. \square

We are now ready for the

Proof of Theorem 3.1. We use an idea of Fill and Janson [7, Sec. 5] and decompose $S(t)$ as $\sum_{k=0}^{\infty} S_k(t)$, and each $S_k(t)$ further as $\sum_{w \in \Sigma^k} S_w(t)$, where for an integer k and a prefix $w \in \Sigma^k$ we define (with little possibility of notational confusion)

$$\begin{aligned} S_k(t) &:= \text{number of comparisons of } (k+1)\text{st symbols,} \\ S_w(t) &:= \text{number of comparisons of } (k+1)\text{st symbols between keys with prefix } w. \end{aligned}$$

A major advantage of working in continuous time is that,

$$(3.3) \quad \text{for each fixed } k \text{ and } t, \text{ the variables } S_w(t) \text{ with } w \in \Sigma^k \text{ are independent.}$$

A further key observation, clear after a moment's thought, is this: For each $w \in \Sigma^*$, as stochastic processes,

$$(3.4) \quad (S_w(t) : t \in [0, \infty)) \text{ is a probabilistic replica of } (K(p_w t) : t \in [0, \infty)).$$

We define corresponding normalized variables as follows:

$$Y_k(t) := \frac{S_k(t) - \mathbf{E} S_k(t)}{t}, \quad Y_w(t) := \frac{S_w(t) - \mathbf{E} S_w(t)}{t},$$

with the normalized variable $Y(t)$ corresponding to $S(t)$ defined at (3.1). Then

$$Y(t) = \sum_{k=0}^{\infty} Y_k(t), \quad Y_k(t) = \sum_{w \in \Sigma^k} Y_w(t) \quad (k = 0, 1, \dots).$$

To complete the proof we then need only find random variables $Y_k(\infty)$ such that the hypotheses of Lemma 2.8 are satisfied for some $p' \in (p, \infty)$.

But, for each $w \in \Sigma^*$, the existence of an almost-sure limit, call it $Y_w(\infty)$, for $Y_w(t)$ as $t \rightarrow \infty$ follows from (3.4) and Lemma 3.3; indeed, we see that $Y_w(\infty)$ has the same distribution as $p_w T$, with T as in Lemma 3.3. Taking the finite sum over $w \in \Sigma^k$, we see that $Y_k(\infty)$ can be defined as $\sum_{w \in \Sigma^k} Y_w(\infty)$ to meet hypothesis (i) of Lemma 2.8.

To verify the remaining hypotheses we choose $t_0 = 1$ and need to bound the L^q -norm of $Y_k(t)$ for k a nonnegative integer, $t \in [1, \infty)$, and $q \in \{p, p'\}$. According to Lemma 3.4 to follow, for any real $q \in [2, \infty)$ there exists a constant c'_q such that

$$\|Y_k(t)\|_q \leq c'_q \left(\sum_{w \in \Sigma^k} p_w^2 \right)^{1/q}$$

for such k and t . Thus hypotheses (ii) and [for any $p' \in (p, \infty)$] (ii') of Lemma 2.8 hold, and the assumption (3.2) implies that (iii) does, as well. \square

Lemma 3.4. *Adopt the notation in the above proof of Theorem 3.1. Then for every real $q \in [2, \infty)$, there exists a constant $c'_q < \infty$ such that*

$$\|Y_k(t)\|_q \leq c'_q \left(\sum_{w \in \Sigma^k} p_w^2 \right)^{1/q}$$

for every nonnegative integer k and every $t \in [1, \infty)$.

Proof. Fix $q \in [2, \infty)$. The first step is to use (as did Fill and Janson [8, proof of Proposition 5.7]) Rosenthal's inequality, relying on the fact [recall (3.3)] that $S_k(t)$ is the *independent* sum of $S_w(t)$ with $w \in \Sigma^k$. According to Rosenthal's inequality [19, Theorem 3] (see also, e.g., [10, Theorem 3.9.1]) there exists a constant b_q (depending *only* on q) such that

$$\begin{aligned} t^q \|Y_k(t)\|_q^q &= \|S_k(t) - \mathbf{E} S_k(t)\|_q^q \\ &\leq b_q \max \left\{ \sum_{w \in \Sigma^k} \|S_w(t) - \mathbf{E} S_w(t)\|_q^q, \left[\sum_{w \in \Sigma^k} \|S_w(t) - \mathbf{E} S_w(t)\|_2^2 \right]^{q/2} \right\}. \end{aligned}$$

Utilizing (3.4) and Lemma 2.7 together with the assumptions $t \geq 1$ and $q \geq 2$ we therefore find

$$\begin{aligned} \|Y_k(t)\|_q^q &\leq b_q \max \left\{ \sum_{w \in A_k(t)} c_q^q p_w^q + \sum_{w \in B_k(t)} (2c_q)^q p_w^2, \left(\sum_{w \in \Sigma^k} c_2^2 p_w^2 \right)^{q/2} \right\} \\ &\leq b_q \max \left\{ (2c_q)^q \sum_{w \in \Sigma^k} p_w^2, c_2^q \left(\sum_{w \in \Sigma^k} p_w^2 \right)^{q/2} \right\} \\ &\leq (c'_q)^q \sum_{w \in \Sigma^k} p_w^2, \end{aligned}$$

where $A_k(t)$ and $B_k(t)$ are the intersections of those Σ^k with $\{w : p_w t \geq 1\}$ and $\{w : p_w t < 1\}$, respectively, and

$$c'_q := b_q^{1/q} \max\{2c_q, c_2\}.$$

This completes the proof of Lemma 3.4. \square

3.2. Identification of the limit variable Y . In a later draft of this paper (or possibly elsewhere) we hope to expand carefully on the following ideas concerning explicit identification of the limit variable Y appearing in Theorem 3.1. We would like to gain enough understanding of Y , for example, that moments of any given order—or at least the variance—could be computed explicitly in terms of the fundamental probabilities p_w , $w \in \Sigma^*$, of the source.

Recall from Theorem 3.1 and its proof that the limiting variable Y satisfies $Y = \sum_{k=0}^{\infty} Y_k$, where $Y_k \equiv Y_k(\infty) = \sum_{w \in \Sigma^k} Y_w$ and $Y_w \equiv Y_w(\infty)$ is a probabilistic replica of $p_w T$. So it ought to be possible to identify Y by identifying explicitly the random variable T in Régnier's [18] theorem (recall our Lemma 2.5).

But T can indeed be identified. Since we have now reduced to counting key comparisons, there is no loss of generality in assuming that the iid keys are uniformly distributed over $(0, 1)$. Construct an (almost surely complete) infinite rooted binary search tree in the usual way by starting with an empty tree and inserting each key as it is generated. Label the nodes in the natural binary way: the root gets an empty label, its left (respectively, right) child is labeled 0 (resp., 1), the left child of node 0 is labeled 00, etc. Let U_θ denote the key inserted at node θ . Let L_θ (resp., R_θ) denote the largest key smaller than U_θ (resp., smallest key larger than U_θ) inserted at any ancestor of θ , with the exceptions $L_\theta := 0$ and $R_\theta := 1$ if the specified

ancestor keys don't exist. Then one can prove that T is the limit as $\ell \rightarrow \infty$, both almost surely and in L^p for any finite p , of

$$T_\ell := \sum_{|\theta| \leq \ell} (R_\theta - L_\theta)g(U_\theta)$$

where $g(u) = 1 + 2u \ln u + 2(1-u) \ln(1-u)$ and $|\theta|$ is the length of the label θ . This result can be viewed as a marriage and extension of the **QuickSort** limit theorems of Régnier [18] and Rösler [20], since the former established the existence of T (but not its form) and the latter argued (cf. his Section 5) that $T_\ell \xrightarrow{\mathcal{L}} T$.

Acknowledgment. We thank Svante Janson for excellent suggestions that led to improvements to our main theorem.

REFERENCES

- [1] K. L. Chung. *A Course in Probability Theory*. Academic Press, London, 3rd edition, 2001.
- [2] J. Clément, P. Flajolet, and B. Vallée. Dynamical sources in information theory: a general analysis of trie structures. *Algorithmica*, 29(1-2):307–369, 2001. Average-case analysis of algorithms (Princeton, NJ, 1998).
- [3] J. A. Fill and T. Nakama. Analysis of the expected number of bit comparisons required by Quickselect. To appear in *Algorithmica*, 2009.
- [4] James Allen Fill and Svante Janson. A characterization of the set of fixed points of the Quicksort transformation. *Electron. Comm. Probab.*, 5:77–84 (electronic), 2000.
- [5] James Allen Fill and Svante Janson. Smoothness and decay properties of the limiting Quicksort density function. In *Mathematics and computer science (Versailles, 2000)*, Trends Math., pages 53–64. Birkhäuser, Basel, 2000.
- [6] James Allen Fill and Svante Janson. Quicksort asymptotics. *J. Algorithms*, 44(1):4–28, 2002. Analysis of algorithms.
- [7] James Allen Fill and Svante Janson. The number of bit comparisons used by Quicksort: an average-case analysis. In *Proceedings of the Fifteenth Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 300–307 (electronic), New York, 2004. ACM.
- [8] James Allen Fill and Svante Janson. The number of bit comparisons used by Quicksort: an average-case analysis (2010.03.24 draft of full paper). Available from <http://www.ams.jhu.edu/~fill/papers/BitsQuickfulldraft.pdf>, 2010.
- [9] James Allen Fill and Také Nakama. Distributional convergence for the number of symbol comparisons used by QuickSelect (draft). Available from <http://www.ams.jhu.edu/~fill/papers/QSelectdistndraft.pdf>, 2010.
- [10] Allan Gut. *Probability: a graduate course*. Springer Texts in Statistics. Springer, New York, 2005.
- [11] C. A. R. Hoare. Quicksort. *Comput. J.*, 5:10–15, 1962.
- [12] C. R. Hoare. Find (algorithm 65). *Communications of the ACM*, 4:321–322, 1961.
- [13] Olav Kallenberg. *Foundations of modern probability*. Probability and its Applications (New York). Springer-Verlag, New York, 1997.
- [14] C. Knessl and W. Szpankowski. Quicksort algorithm again revisited. *Discrete Mathematics and Theoretical Computer Science*, 3:43–64, 1999.
- [15] D. E. Knuth. *The Art of Computer Programming. Volume 3: Sorting and Searching*. Addison-Wesley, Reading, Massachusetts, 1998.
- [16] Takéhiko Nakama. *Analysis of Execution Costs for QuickSelect*. Ph.D. dissertation, The Johns Hopkins University, Department of Applied Mathematics and Statistics, August 2009. Available from <http://www.ams.jhu.edu/~fill/papers/NakamaDissertation.pdf>.
- [17] R. Neininger and L. Rüschenhof. Rates of convergence for Quickselect. *Journal of Algorithms*, 44:51–62, 2002.
- [18] M. Régnier. A limiting distribution of Quicksort. *RAIRO Informatique Théorique et Applications*, 23:335–343, 1989.
- [19] Haskell P. Rosenthal. On the subspaces of L^p ($p > 2$) spanned by sequences of independent random variables. *Israel J. Math.*, 8:273–303, 1970.

- [20] U. Rösler. A limit theorem for Quicksort. *RAIRO Informatique Théorique et Applications*, 25:85–100, 1991.
- [21] U. Rösler and L. Rüschendorf. The contraction method for recursive algorithms. *Algorithmica*, 29(1):3–33, 2001.
- [22] B. Vallée, J. Clément, J. A. Fill, and P. Flajolet. The number of symbol comparisons in Quicksort and Quickselect. In S. Albers et al., editor, *36th International Colloquium on Automata, Languages and Programming (ICALP 2009), Part I, LNCS 5555*, pages 750–763. Springer-Verlag, 2009.

DEPARTMENT OF APPLIED MATHEMATICS AND STATISTICS, THE JOHNS HOPKINS UNIVERSITY,
34TH AND CHARLES STREETS, BALTIMORE, MD 21218-2682 USA

E-mail address: jimfill@jhu.edu