## The Number of Symbol Comparisons in QuickSort & QuickSelect

I. Overview ~~~ Philippe Flajolet
 II. Average-Case Analysis ~~~ Brigitte Vallée
 III. Distributions ~~~ Jim Fill

- I.Algorithms & analysis
- 2. Cost measures
- 3. Sources (data model)
- 4. Results: average-case & distributional



# 1.QuickSort & QuickSelect



QuickSort 
$$(n, A)$$
: sorts the array  $A$   
Choose a pivot;  
 $(k, A_-, A_+) := \text{Partition}(A)$ ;  
QuickSort  $(k - 1, A_-)$ ;  
QuickSort  $(n - k, A_+)$ .



## Analyses of QuickSort

- Average-case: recurrences, then generating functions (GFs). Exchanges; Median-of-3, etc.
- Variance: multivariate GFs
- Distribution: MGFs & moments, Martingales, Contraction

Hoare; Knuth; Sedgewick [1960-1975] Hennequin, Régnier, Rösler [1989+] Fill & Janson [2000], Martinez...  $\begin{array}{l} \texttt{QuickSelect}\;(n,m,A)\text{: returns the value of the element of rank}\;\underline{m}\,\text{in}\;A.\\\\ \texttt{Choose a pivot;}\\ (k,A_-,A_+) \mathrel{\mathop:}= \texttt{Partition}(A)\text{;}\\\\ \texttt{If}\;m=k\;\text{then}\;\texttt{QuickSelect}\;\mathrel{\mathop:}=\;\texttt{pivot}\\\\ \texttt{else if}\;m<k\;\text{then}\;\texttt{QuickSelect}\;(k-1,m,A_-)\\\\\\ \texttt{else QuickSelect}\;(n-k,m-k,A_+)\text{;} \end{array}$ 



#### Various brands of QuickSelect:

| $\mathtt{QuickMin}(n)$            | := QuickSelect(1,n)                                     | finds the minimum                   |
|-----------------------------------|---------------------------------------------------------|-------------------------------------|
| $\mathtt{QuickMax}(n)$            | := QuickSelect(n, n)                                    | finds the maximum.                  |
| $\texttt{QuickQuant}_{\alpha}(n)$ | $:= \texttt{QuickSelect} (\lfloor \alpha n \rfloor, n)$ | finds the $\alpha$ -quantile        |
| $\mathtt{QuickMed}(n)$            | $:= \texttt{QuickSelect} (\lfloor n/2  floor, n)$       | finds the median                    |
| $\mathtt{QuickRand}(n)$           | $:= \texttt{QuickSelect}\ (m,n)$                        | for a rank $m \in [1n]_\mathcal{R}$ |

Average-case analyses

Mean number  $K_n$  of key comparisons



 $H_{\alpha} = \text{the entropy function} = \alpha |\log \alpha| + (1 - \alpha) |\log(1 - \alpha)|$ 

#### Knuth et al [ca 1970]

**Distributional analyses** 

- Quickselect: e.g., Dickman distribution Mahmoud-Modarres-Smythe, Grübel, Rösler, Hwang-Tsai, et al.
   perpetuities: 1+U1+U1U2+U1U2U3+... i.i.d. unif. [0, 1] (fixed rank; fixed quantile)
- Multiple Quickselect, ancestors, &c Lent-Mahmoud, Prodinger, et al.



# 2. Cost Measures



#### Sedgewick @ AofA-02(?): "actual complexity matters!"



- So far: number of key-comparisons
- But... keys are often "non-atomic" records!
- And...need common information-theoretic basis, to compare with radix methods, hashing, etc.

#### Alphabet: $\Sigma$

- Count all symbol comparisons in algorithms:
- comparing u and v has cost  $1 + \frac{\text{coincidence}}{(u,v)}$ .

coincidence=3; #comparisons=4.
 (γ)

$$\begin{split} A &= abbbbbbaaabab \quad B &= abbbbbbaabaa \quad C &= baabbbabbbba \quad D &= bbbaabbbbbaab \quad E &= bbabbaababbb \\ F &= abbbbbbbbbbbbbb \quad G &= bbaabbabbaba \quad H &= ababbbabbbab \quad I &= bbbaabbbbbbb \quad J &= abaabbbbbaabb \\ K &= bbbabbbbbbbbaa \quad L &= aaaabbabaabaa \quad M &= bbbaaabbbbbb \quad N &= abbbbbbbbabbaa \quad O &= abbabbabbbbb \quad P &= bbabbbbaaaabbb \\ \end{split}$$



## Under a wide range of *classical* STRING (WORD) MODELS:

It takes O(n.log n) <u>symbol comparisons</u> to "distinguish" <u>n</u> elements ---- in probability, on average

#### With high probability, the <u>common prefix</u> of any two words has length at most O(log *n*).

#### Many many people in the audience...

- Bernoulli, Markov, etc.
- Devroye's density model
- Vallée's dynamic sources...





Symbol comparisons

- QuickSort: [Janson & Fill 2002] binary source + density model.
  - $\sim Cn.log(n)^2$
- QuickSelect: [Fill-Nakama 2007-9] binary source for QuickMin/Max & QuickRand



CONSTANTS?

(cf also: Panholzer & Prodinger)



# 3. Sources

"A source models the way data (symbols) are produced."



### Axioms for SOURCES

- Totally ordered alphabet (usually finite) ∑
- Fundamental probabilities (p<sub>w</sub>) :=
   the probability of starting with w

• 
$$p_w \rightarrow 0$$
 as  $|w| \rightarrow \infty$ 

• Keys are invariably i.i.d.

[Later] + "regularity" conditions: <u>tameness</u>



**Property**: The Source is parameterized by [0, 1]: to an infinite word  $\underline{w}$ , there corresponds  $\underline{\alpha}$  such that  $M(\alpha)=w$ .

#### Notations:



## Fundamental constants of QuickStuffs will be all expressed in terms of fundamental probabilities

- Standard binary source (uniform: 1/2,1/2);
   Bernoulli sources such as 1/2, 1/6,1/3.
- <u>Density models</u>: Standard binary source with density f(x) or c.d.f F(x).
- Markov
- **Dynamical sources**

#### [Devroye 1986] [Vallée 2001; Clément-FI-Vallée 2001]



# 4. Results



(Le Savant Cosinus)

#### Average-case

**Theorem 1** For any tamed source, the mean number  $S_n$  of symbol comparisons used by QuickSort(n) satisfies

$$S_n \sim \frac{1}{h_S} n \log^2 n.$$

and involves the entropy  $h_{\mathcal{S}}$  of the source  $\mathcal{S}$ , defined as

$$h_{\mathcal{S}} := \lim_{k \to \infty} \left[ \frac{-1}{k} \sum_{w \in \Sigma^k} p_w \log p_w \right],$$

where  $p_w$  is the probability that a word begins with prefix w.

**Theorem 2** For any weakly tamed source, the mean number of symbol comparisons used by  $QuickQuant_{\alpha}(n)$  satisfies  $q_n^{(\alpha)} \sim \rho_{\mathcal{S}}(\alpha) n$ 

QuickMin, QuickRand

QuickVal @ @

#### **QUICKVAL(\alpha):** is dual to QuickSelect



- QuickVal(n, α) := rank of element whose parameter
   [corresponding to value ν] is α.
- QuickVal( $n, \alpha$ ) behaves "almost" as QuickSelect( $n\alpha$ ).

#### Distribution

**Theorem**: Assuming a suitable tameness condition, there exists a limiting distribution of the cost  $S_n/n$  of  $QuickQuant(\alpha)$ , which can be described explicitly