### Median-of-k Quicksort Is Optimal For Many Equal Keys

Sebastian Wild

wild@cs.uni-kl.de

TECHNISCHE UNIVERSITÄT KAISERSLAUTERN

originates from joint work with Martin Aumüller, Martin Dietzfelbinger, Conrado Martínez, and Markus Nebel

#### AofA 2017

28th International Meeting on Probabilistic, Combinatorial and Asymptotic Methods for the Analysis of Algorithms

Sebastian Wild

**Quicksort Is Optimal For Many Equal Keys** 

2017-06-19

0 / 16





# **Quicksort and Search Trees**



**Saturated Fringe-Balanced Trees** 



**Back to Multiset Permutations** 

• Extensive literature and results on Quicksort

Type of	result	Analysis Techniqu	les	Algorithm	variants	Cost Mea	sures

- **but:** most results consider **random permutations** as input!
- partly justified: we can (should!) randomize Quicksort,
- Catch: Elements with equal keys won't go away!

< (1) →

• Extensive literature and results on Quicksort



- **but:** most results consider **random permutations** as input!
- partly justified: we can (should!) randomize Quicksort,
- Catch: Elements with equal keys won't go away!

• Extensive literature and results on Quicksort

Type of result	Analysis Techniques	Algorithm variants	Cost Measures
expected costs			
variance			

- **but:** most results consider **random permutations** as input!
- partly justified: we can (should!) randomize Quicksort,
  - → every input appears randomly ordered
- Catch: Elements with equal keys won't go away!

• Extensive literature and results on Quicksort

Type of result	Analysis Techniques	Algorithm variants	Cost Measures
expected costs			
variance			
tail inequalities			

- but: most results consider random permutations as input! not done in libraries
- partly justified: we can (should!) randomize Quicksort,
  - ••• every input appears randomly ordered
- Catch: Elements with equal keys won't go away!

< ∂ >

• Extensive literature and results on Quicksort

Type of result	Analysis Techniques	Algorithm variants	Cost Measures
expected costs			
variance			
tail inequalities			
limit distributions			

- but: most results consider random permutations as input!
- partly justified: we can (should!) randomize Quicksort,
- Catch: Elements with equal keys won't go away!

<∂>

• Extensive literature and results on Quicksort

Type of result	Analysis Techniques	Algorithm variants	Cost Measures
expected costs			
variance			
tail inequalities			
limit distributions			
semi-)local limit laws			

- but: most results consider random permutations as input!
- partly justified: we can (should!) randomize Quicksort,
- Catch: Elements with equal keys won't go away!

• Extensive literature and results on Quicksort

Type of result	Analysis Techniques	Algorithm variants	Cost Measures
expected costs	telescoping recurrences		
variance			
tail inequalities			
limit distributions			
semi-)local limit laws			

- but: most results consider random permutations as input! not done in libraries
- partly justified: we can (should!) randomize Quicksort,
- Catch: Elements with equal keys won't go away!

<∂>

• Extensive literature and results on Quicksort

Type of result	Analysis Techniques	Algorithm variants	Cost Measures
expected costs variance tail inequalities	telescoping recurrences singularity analysis		
semi-)local limit laws			

- but: most results consider random permutations as input! not done in libraries
- partly justified: we can (should!) randomize Quicksort,
   every input appears randomly ordered
- Catch: Elements with equal keys won't go away!

• Extensive literature and results on Quicksort

Type of result	Analysis Techniques	Algorithm variants	Cost Measures
expected costs variance tail inequalities limit distributions semi-)local limit laws	telescoping recurrences singularity analysis Euler differential eq.		

- but: most results consider random permutations as input! not done in libraries
- partly justified: we can (should!) randomize Quicksort,
   every input appears randomly ordered
- Catch: Elements with equal keys won't go away!

• Extensive literature and results on Quicksort

Type of result	Analysis Techniques	Algorithm variants	Cost Measures
expected costs variance tail inequalities limit distributions semi-)local limit laws	telescoping recurrences singularity analysis Euler differential eq. Martingales		

- but: most results consider random permutations as input!
   ...ot done in libraries...
- partly justified: we can (should!) randomize Quicksort,
   every input appears randomly ordered
- Catch: Elements with equal keys won't go away!

• Extensive literature and results on Quicksort

Type of result	Analysis Techniques	Algorithm variants	Cost Measures
expected costs t variance tail inequalities limit distributions semi-)local limit laws	telescoping recurrences singularity analysis Euler differential eq. Martingales contraction method		

- **but:** most results consider **random permutations** as input! , not done in libraries...
- partly justified: we can (should!) randomize Quicksort,
   every input appears randomly ordered
- Catch: Elements with equal keys won't go away!

<∂>

• Extensive literature and results on Quicksort

Type of result	Analysis Techniques	Algorithm variants	Cost Measures
expected costs variance tail inequalities limit distributions semi-)local limit laws	telescoping recurrences singularity analysis Euler differential eq. Martingales contraction method branching processes		

- but: most results consider random permutations as input!
- partly justified: we can (should!) randomize Quicksort,
   every input appears randomly ordered
- Catch: Elements with equal keys won't go away!

• Extensive literature and results on Quicksort

Type of result	Analysis Techniques	Algorithm variants	Cost Measures
expected costs variance tail inequalities limit distributions (semi-)local limit laws	telescoping recurrences singularity analysis Euler differential eq. Martingales contraction method branching processes continuous master theorem		

- but: most results consider **random permutations** as input!
- partly justified: we can (should!) randomize Quicksort,
- Catch: Elements with equal keys won't go away!

• Extensive literature and results on Quicksort

Type of result	Analysis Techniques	Algorithm variants	Cost Measures
expected costs variance tail inequalities limit distributions (semi-)local limit laws	telescoping recurrences singularity analysis Euler differential eq. Martingales contraction method branching processes continuous master theorem	pivot sampling	

- but: most results consider random permutations as input!
   partly justified: we can (should!) randomize Quicksort,
  - → **every** input appears randomly ordered
- Catch: Elements with equal keys won't go away!

• Extensive literature and results on Quicksort

Type of result	Analysis Techniques	Algorithm variants	Cost Measures
expected costs variance tail inequalities limit distributions (semi-)local limit laws	telescoping recurrences singularity analysis Euler differential eq. Martingales contraction method branching processes continuous master theorem	pivot sampling Insertionsort cutoff	

- but: most results consider random permutations as input! , not done in libraries...
   partly justified: we can (should!) randomize Quicksort,
   every input appears randomly ordered
- Catch: Elements with equal keys won't go away!

• Extensive literature and results on Quicksort

Type of result	Analysis Techniques	Algorithm variants	Cost Measures
expected costs variance tail inequalities limit distributions (semi-)local limit laws	telescoping recurrences singularity analysis Euler differential eq. Martingales contraction method branching processes continuous master theorem	pivot sampling Insertionsort cutoff multiway partitioning	

- but: most results consider random permutations as input! *indome in libraries*...
   partly justified: we can (should!) randomize Quicksort, *wery input appears randomly ordered*
- Catch: Elements with equal keys won't go away!

• Extensive literature and results on Quicksort

Type of result	Analysis Techniques	Algorithm variants	Cost Measures
expected costs variance tail inequalities limit distributions (semi-)local limit laws	telescoping recurrences singularity analysis Euler differential eq. Martingales contraction method branching processes continuous master theorem	pivot sampling Insertionsort cutoff multiway partitioning Quickselect	

- but: most results consider random permutations as input!
   partly justified: we can (should!) randomize Quicksort,
   → every input appears randomly ordered
- Catch: Elements with equal keys won't go away!

• Extensive literature and results on Quicksort

Type of result	Analysis Techniques	Algorithm variants	Cost Measures
expected costs variance tail inequalities limit distributions (semi-)local limit laws	telescoping recurrences singularity analysis Euler differential eq. Martingales contraction method branching processes continuous master theorem	pivot sampling Insertionsort cutoff multiway partitioning Quickselect constant space	

- but: most results consider random permutations as input!
   partly justified: we can (should!) randomize Quicksort,
   every input appears randomly ordered
- Catch: Elements with equal keys won't go away!

• Extensive literature and results on Quicksort

Type of result	Analysis Techniques	Algorithm variants	Cost Measures
expected costs variance tail inequalities limit distributions (semi-)local limit laws	telescoping recurrences singularity analysis Euler differential eq. Martingales contraction method branching processes continuous master theorem	pivot sampling Insertionsort cutoff multiway partitioning Quickselect constant space	key comparisons

- Catch: Elements with equal keys won't go away!

• Extensive literature and results on Quicksort

Type of result	Analysis Techniques	Algorithm variants	Cost Measures
expected costs variance tail inequalities limit distributions (semi-)local limit laws	telescoping recurrences singularity analysis Euler differential eq. Martingales contraction method branching processes continuous master theorem	pivot sampling Insertionsort cutoff multiway partitioning Quickselect constant space	key comparisons symbol comparisons

- but: most results consider random permutations as input! *indone in libraries*...
   partly justified: we can (should!) randomize Quicksort, *every input appears randomly ordered*
- Catch: Elements with equal keys won't go away!

• Extensive literature and results on Quicksort

Type of result	Analysis Techniques	Algorithm variants	Cost Measures
expected costs variance tail inequalities limit distributions (semi-)local limit laws	telescoping recurrences singularity analysis Euler differential eq. Martingales contraction method branching processes continuous master theorem	pivot sampling Insertionsort cutoff multiway partitioning Quickselect constant space	key comparisons symbol comparisons swaps

- but: most results consider random permutations as input!
   partly justified: we can (should!) randomize Quicksort,
   every input appears randomly ordered
- Catch: Elements with equal keys won't go away!

• Extensive literature and results on Quicksort

Type of result	Analysis Techniques	Algorithm variants	Cost Measures
expected costs variance tail inequalities limit distributions (semi-)local limit laws	telescoping recurrences singularity analysis Euler differential eq. Martingales contraction method branching processes continuous master theorem	pivot sampling Insertionsort cutoff multiway partitioning Quickselect constant space	key comparisons symbol comparisons swaps scanned elements

- Catch: Elements with equal keys won't go away!

• Extensive literature and results on Quicksort

Type of result	Analysis Techniques	Algorithm variants	Cost Measures
expected costs variance tail inequalities limit distributions (semi-)local limit laws	telescoping recurrences singularity analysis Euler differential eq. Martingales contraction method branching processes	pivot sampling Insertionsort cutoff multiway partitioning Quickselect constant space	key comparisons symbol comparisons swaps scanned elements branch misses
	continuous master theorem		branch misses

- but: most results consider random permutations as input!
   partly justified: we can (should!) randomize Quicksort,
   every input appears randomly ordered
- Catch: Elements with equal keys won't go away!

• Extensive literature and results on Quicksort

Type of result	Analysis Techniques	Algorithm variants	Cost Measures
expected costs variance tail inequalities limit distributions (semi-)local limit laws	telescoping recurrences singularity analysis Euler differential eq. Martingales contraction method branching processes	pivot sampling Insertionsort cutoff multiway partitioning Quickselect constant space	key comparisons symbol comparisons swaps scanned elements branch misses
	continuous master theorem		branch misses

- but: most results consider random permutations as input!
   partly justified: we can (should!) randomize Quicksort,
   every input appears randomly ordered
- Catch: Elements with equal keys won't go away!

• Extensive literature and results on Quicksort

Type of result	Analysis Techniques	Algorithm variants	Cost Measures
expected costs	telescoping recurrences singularity analysis	pivot sampling	key comparisons
tail inequalities	Euler differential eq. Martingales	multiway partitioning	swaps
limit distributions	contraction method	Quickselect	scanned elements
(semi-)local limit laws	continuous master theorem	constant space	branch misses

- but: most results consider random permutations as input!
- partly justified: we can (should!) randomize Quicksort,
   *wery* input appears randomly ordered
- Catch: Elements with equal keys won't go away!

<∂>

• Extensive literature and results on Quicksort

Type of result	Analysis Techniques	Algorithm variants	Cost Measures
expected costs variance tail inequalities	telescoping recurrences singularity analysis Euler differential eq. Martingales	pivot sampling Insertionsort cutoff multiway partitioning	key comparisons symbol comparisons swaps
limit distributions (semi-)local limit laws	contraction method branching processes continuous master theorem	Quickselect constant space	scanned elements branch misses

- **but:** most results consider **random permutations** as input!
- partly justified: we can (should!) randomize Quicksort,
   every input appears randomly ordered
- Catch: Elements with equal keys won't go away!

• Extensive literature and results on Quicksort

Type of result	Analysis Techniques	Algorithm variants	Cost Measures
expected costs variance tail inequalities limit distributions (semi-)local limit laws	telescoping recurrences singularity analysis Euler differential eq. Martingales contraction method branching processes	pivot sampling Insertionsort cutoff multiway partitioning Quickselect constant space	key comparisons symbol comparisons swaps scanned elements branch misses
	continuous musici meorem		

- **but:** most results consider **random permutations** as input!
- partly justified: we can (should!) randomize Quicksort,
   every input appears randomly ordered
- Catch: Elements with equal keys won't go away!

#### Assumptions:



#### Multiset Model

Random permutation  $U_1, \ldots, U_n$  of **fixed multiset**  $x_1, \ldots, x_u$  number of occurrences of values  $1, \ldots, u$ 

**Discrete i.i.d. Model:**   $U_1, \ldots, U_n$  i.i.d. with  $Pr[U_1 = v] = q_v$  $\vec{q} = (q_1, \ldots, q_u)$  a **fixed** universe **distribution** 

### 2 fat-pivot partitioning

- all duplicates of pivots removed
- subproblems of same type, (restricted to a sub-universe)
- 3 Cost: # ternary comparisons

#### Median-of-(2t+1) Quicksort:

• median-of-(2t + 1





2017-06-19

2 / 16

< (1) →

### Assumptions: **1** Input:

#### **Multiset Model:**

Random permutation  $U_1, \ldots, U_n$  of fixed multiset  $x_1, \ldots, x_n$  number of occurrences of values  $1, \ldots, n$ 

(A)

- **Cost:** # ternary comparisons

#### *Median-of-(2t+1) Quicksort:*



2017-06-19

2/16

#### Assumptions:



#### 2 fat-pivot partitioning

- all duplicates of pivots removed
- subproblems of same type, (restricted to a sub-universe)
- 3 Cost: # ternary comparisons

#### Median-of-(2t+1) Quicksort:

median-of-(2t + 1)



< (1) →

#### Assumptions:



#### 2 fat-pivot partitioning

- all duplicates of pivots removed
- subproblems of same type, (restricted to a sub-universe)
- 3 Cost: # ternary comparisons

#### Median-of-(2t+1) Quicksort:

• median-of-(2t + 1



#### ot today

#### Assumptions:



- **Cost:** # ternary comparisons

*Median-of-(2t+1) Ouicksort:* 

2017-06-19

#### Assumptions:



#### 2 fat-pivot partitioning

- all duplicates of pivots removed
- → subproblems of same type, (restricted to a sub-universe
- 3 Cost: # ternary comparisons

Median-of-(2t+1) Quicksort:

• median-of-(2t + 1

recursive call recursive call rse)



2017-06-19

#### Assumptions:

 Input: (A) Multiset Model: Random permutation U<sub>1</sub>,..., U<sub>n</sub> of fixed multiset x<sub>1</sub>,..., x<sub>u</sub> number of occurrences of values 1,..., u profile x of input

 (B) Discrete i.i.d. Model: U<sub>1</sub>,..., U<sub>n</sub> i.i.d. with Pr[U<sub>1</sub> = v] = q<sub>v</sub> → random profile X = Mult(n, q̄) q̄ = (q<sub>1</sub>,..., q<sub>u</sub>) a fixed universe distribution

#### 2 fat-pivot partitioning

• all duplicates of pivots removed



- → subproblems of **same** type, (restricted to a sub-universe)
- 3 Cost: # ternary comparisons

Median-of-(2t+1) Quicksort:

• median-of-(2t + 1

not toda



2017-06-19
### Assumptions:

 Input: (A) Multiset Model: Random permutation U<sub>1</sub>,..., U<sub>n</sub> of fixed multiset x<sub>1</sub>,..., x<sub>u</sub> number of occurrences of values 1,..., u profile x of input

 (B) Discrete i.i.d. Model: U<sub>1</sub>,..., U<sub>n</sub> i.i.d. with Pr[U<sub>1</sub> = v] = q<sub>v</sub> → random profile X = Mult(n, q̄) q̄ = (q<sub>1</sub>,..., q<sub>u</sub>) a fixed universe distribution

- 2 fat-pivot partitioning
  - all duplicates of pivots removed





- →→ subproblems of same type, (restricted to a sub-universe)
- **3** Cost: # ternary comparisons

Median-of-(2t+1) Quicksort:

median-of-(2t + 1)

not today



2017-06-19

P P P

### Assumptions:

• Input: (A) • Input: (A) • Multiset Model: Random permutation  $U_1, ..., U_n$  of fixed multiset  $x_1, ..., x_u$  number of occurrences of values 1, ..., u• profile  $\vec{x}$  of input Discrete *i.i.d.* Model:  $U_1, ..., U_n$  i.i.d. with  $\Pr[U_1 = v] = q_v$  ··· random profile  $\vec{X} \stackrel{2}{=} Mult(n, \vec{q})$  $\vec{q} = (q_1, ..., q_u)$  a fixed universe distribution

- 2 fat-pivot partitioning
  - all duplicates of pivots removed



**3** Cost: # ternary comparisons

Median-of-(2t+1) Quicksort:

median-of-(2t + 1

 P
 P
 P

 recursive call
 recursive call

### Assumptions:

• Input: (A) • Input: (A) • Multiset Model: Random permutation  $U_1, ..., U_n$  of fixed multiset  $x_1, ..., x_u$  number of occurrences of values 1, ..., u• profile  $\vec{x}$  of input **Discrete i.i.d. Model:**   $U_1, ..., U_n$  i.i.d. with  $\Pr[U_1 = v] = q_v \longleftarrow$  random profile  $\vec{X} \stackrel{a}{=} Mult(n, \vec{q})$  $\vec{q} = (q_1, ..., q_u)$  a fixed universe distribution

- 2 fat-pivot partitioning
  - all duplicates of pivots removed



- → subproblems of same type, (restricted to a sub-universe)
- **3** Cost: # ternary comparisons

### Median-of-(2t+1) Quicksort:

- median-of-(2t + 1)
- (extension to asymmetric sampling possible)



### Assumptions:

• Input: (A) • Input: (A) • Multiset Model: Random permutation  $U_1, ..., U_n$  of fixed multiset  $x_1, ..., x_u$  number of occurrences of values 1, ..., u• profile  $\vec{x}$  of input **Discrete i.i.d. Model:**   $U_1, ..., U_n$  i.i.d. with  $\Pr[U_1 = v] = q_v \longleftarrow$  random profile  $\vec{X} \stackrel{2}{=} Mult(n, \vec{q})$  $\vec{q} = (q_1, ..., q_u)$  a fixed universe distribution

- 2 fat-pivot partitioning
  - all duplicates of pivots removed



- →→ subproblems of same type, (restricted to a sub-universe)
- **3** Cost: # ternary comparisons

### Median-of-(2t+1) Quicksort:

- median-of-(2t + 1)
- (extension to asymmetric sampling possible)



### Assumptions:

• Input: (A) • Input: (A) • Multiset Model: Random permutation  $U_1, ..., U_n$  of fixed multiset  $x_1, ..., x_u$  number of occurrences of values 1, ..., u• profile  $\vec{x}$  of input **Discrete i.i.d. Model:**   $U_1, ..., U_n$  i.i.d. with  $\Pr[U_1 = v] = q_v \longleftarrow$  random profile  $\vec{X} \stackrel{a}{=} Mult(n, \vec{q})$  $\vec{q} = (q_1, ..., q_u)$  a fixed universe distribution

- 2 fat-pivot partitioning
  - all duplicates of pivots removed



- → subproblems of same type, (restricted to a sub-universe)
- 3 Cost: # ternary comparisons

### Median-of-(2t+1) Quicksort:

- median-of-(2t + 1)
- (extension to asymmetric sampling possible)





Quicksort Is Optimal For Many Equal Keys

not today

2 / 16

### Rather little is known!

- Sedgewick 1977: Quicksort on Equal Keys
- Sedgewick & Bentley 2002: Quicksort is Optimal (Talk at Knuthfest)

A bit more on BSTs:

- Burge 1976: An Analysis of BSTs Formed from Sequences of Nondistinct Keys
- Kemp 1996: BSTs constructed from nondistinct keys with/without specified probabilities
- Archibald & Clément 2006: Average depth in a BST with repeated keys

This is basically **all** literature on analysis of Quicksort with equal keys!

< (1) →

### Rather little is known!

- Sedgewick 1977: Quicksort on Equal Keys
- Sedgewick & Bentley 2002: Quicksort is Optimal (Talk at Knuthfest)

A bit more on BSTs:

- Burge 1976: An Analysis of BSTs Formed from Sequences of Nondistinct Keys
- Kemp 1996: BSTs constructed from nondistinct keys with/without specified probabilities
- Archibald & Clément 2006: Average depth in a BST with repeated keys

This is basically **all** literature on analysis of Quicksort with equal keys!

SIAM J. COMPUT. Vol. 6, No. 2, June 1977

#### QUICKSORT WITH EQUAL KEYS\*

ROBERT SEDGEWICK<sup>†</sup>

Abstract. This paper considers the problem of implementing and analyzing a Quicksort program when equal keys are likely to be present in the file to be sorted. Upper and lower bounds are derived on the average number of comparisons needed by any Quicksort program when equal keys are present. It is shown that, of the three strategies which have been suggested for dealing with equal keys, the method of always stopping the scanning pointers on keys equal to the partitioning element performs best.

Key words. analysis of algorithms, equal keys, Quicksort, sorting

Quicksort Is Optimal For Many Equal Keys

< (1) →

Rather little is known!

- Sedgewick 1977: Quicksort on Equal Keys
- Sedgewick & Bentley 2002: Quicksort is Optimal (Talk at Knuthfest)

A bit more on BSTs:

- Burge 1976: An Analysis of BSTs Formed from Sequences of Nondistinct Keys
- Kemp 1996: BSTs constructed from nondistinct keys with/without specified probabilities
- Archibald & Clément 2006: Average depth in a BST with repeated keys

This is basically **all** literature on analysis of Quicksort with equal keys!



< 67 ►

3/16

Rather little is known!

- Sedgewick 1977: Quicksort on Equal Keys
- Sedgewick & Bentley 2002: Quicksort is Optimal (Talk at Knuthfest)

A bit more on BSTs:

- Burge 1976: An Analysis of BSTs Formed from Sequences of Nondistinct Keys
- Kemp 1996: BSTs constructed from nondistinct keys with/without specified probabilities
- Archibald & Clément 2006: Average depth in a BST with repeated keys

This is basically **all** literature on analysis of Quicksort with equal keys!



<⊡>

3/16

Rather little is known!

- Sedgewick 1977: Quicksort on Equal Keys
- Sedgewick & Bentley 2002: Quicksort is Optimal (Talk at Knuthfest)

A bit more on BSTs:

- Burge 1976: An Analysis of BSTs Formed from Sequences of Nondistinct Keys
- Kemp 1996: BSTs constructed from nondistinct keys with/without specified probabilities
- Archibald & Clément 2006: Average depth in a BST with repeated keys

This is basically **all** literature on analysis of Quicksort with equal keys!



< A >

Rather little is known!

- Sedgewick 1977: Quicksort on Equal Keys
- Sedgewick & Bentley 2002: Quicksort is Optimal (Talk at Knuthfest)

A bit more on BSTs:

- Burge 1976: An Analysis of BSTs Formed from Sequences of Nondistinct Keys
- Kemp 1996: BSTs constructed from nondistinct keys with/without specified probabilities
- Archibald & Clément 2006: Average depth in a BST with repeated keys

### This is basically **all** literature on analysis of Quicksort with equal keys!



Rather little is known!

- Sedgewick 1977: Quicksort on Equal Keys
- Sedgewick & Bentley 2002: Quicksort is Optimal (Talk at Knuthfest)

A bit more on BSTs:

- Burge 1976: An Analysis of BSTs Formed from Sequences of Nondistinct Keys
- Kemp 1996: BSTs constructed from nondistinct keys with/without specified probabilities
- Archibald & Clément 2006: Average depth in a BST with repeated keys

This is basically **all** literature on analysis of Quicksort with equal keys!



Rather little is known!

- Sedgewick 1977: Quicksort on Equal Kevs
- Sedgewick & Bentley 2002: Quicksd

A bit more on BSTs:

- Burge 1976: An Analysis of BSTs Forme
- Kemp 1996: BSTs constructed from non
- Archibald & Clément 2006: Average depth in a BST with repeated keys

This is basically all literature on analysis of Quicksort with equal keys!



All these works only consider classic Quicksort:

- No sampling to choose pivots.
- (No multiway partitioning.)

**Classic Quicksort:** 

< @ ►

### **Classic Quicksort:**

# Analysis of Quicksort with equal keys 1. Define $C(x_1,...,x_n) = C(1,n)$ to be the mean # compares to sort the file $C(1,n) = N - 1 + \frac{1}{N} \sum_{1 \le j \le n} x_j(C(1,j-1) + C(j+1,n))$ 2. Multiply both sides by $N = x_1 + ... + x_n$ $NC(1,n) = N(N-1) + \sum_{1 \le j \le n} x_jC(1,j-1) + \sum_{1 \le j \le n} x_jC(j+1,n)$ 3. Subtract same equation for $x_2,...,x_n$ and let D(1,n) = C(1,n) - C(2,n) $(x_1 + ... + x_n)D(1,n) = x_1^2 - x_1 + 2x_1(x_2 + ... + x_n) + \sum_{j \ge j} x_jD(1,j-1)$

4. Subtract same equation for  $x_{1,...,x_{n-1}}$ 

$$(x_1 + ... + x_n)D(1, n) - (x_1 + ... + x_{n-1})D(1, n-1) = 2x_1x_n + x_nD(1, n-1)$$

### **Classic Quicksort:**

Analysis of Quicksort with equal keys 1. Define  $C(x_1,...,x_n) = C(1,n)$  to be the mean # compares to sor  $C(1,n) = N - 1 + \frac{1}{N} \sum_{1 < j < n} x_j (C(1, j-1) + C(j+1, n))$ 2. Multiply both sides by  $N = x_1 + ... + x_n$  $NC(1, n) = N(N-1) + \sum_{j} x_{j}C(1, j-1) + \sum_{j} x_{j}C(j+1, n)$ 1≤i≤n 1≤i≤n 3. Subtract same equation for  $x_2, ..., x_n$  and let D(1, n) = C(1, n) - C(1, n) $(x_1 + ... + x_n)D(1, n) = x_1^2 - x_1 + 2x_1(x_2 + ... + x_n) + \sum x_jD(1, j)$ 2≤j≤n 4. Subtract same equation for  $x_{1,...,x_{n-1}}$  $(x_1 + ... + x_n)D(1, n) - (x_1 + ... + x_{n-1})D(1, n-1) = 2x_1x_n + x_nD(1, n-1)$ 

#### Analysis of Quicksort with equal keys (cont.)

$$(x_1 + ... + x_n)D(1, n) - (x_1 + ... + x_{n-1})D(1, n-1) = 2x_1x_n + x_nD(1, n-1)$$

5. Simplify, divide both sides by  $N = x_1 + ... + x_n$ 

$$D(1,n) = D(1,n-1) + \frac{2x_1x_n}{x_1 + ... + x_n}$$

6. Telescope (twice)

$$C(1,n) = N - n + \sum_{1 \le k < j \le n} \frac{2x_k x_j}{x_k + \dots + x_j}$$

THEOREM. Quicksort (with 3-way partitioning, randomized) uses

$$N-n+2QN \quad \text{compares (where } Q = \sum_{1 \le k < j \le n} \frac{p_k p_j}{p_k + \ldots + p_j}, \text{ with } p_i = x_i / N)$$

to sort an  $(x_1, ..., x_n)$  – file, on the average .

< 17 >

Classic Quicksort: Expected comparisons expressible exactly.



#### Quicksort is optimal

The average number of compares per element C/N is always within a constant factor of the entropy H

lower bound: C > NH - N (information theory) upper bound: C < 2In2NH + N (Burge analysis, Melhorn bound)

No comparison-based algorithm can do better.

Conjecture: With sampling,  $C / N \rightarrow H$  as sample size increases.

#### Quicksort is optimal

The average number of compares per element C/N is always within a constant factor of the entropy H

lower bound: C > NH - N (information theory) upper bound: C < 2In2NH + N (Burge analysis, Melhorn bound)

No comparison-based algorithm can do better.

Conjecture: With sampling,  $C / N \rightarrow H$  as sample size increases.

#### Quicksort is optimal

The average number of compares per element C/N is always within a constant factor of the entropy H

lower bound: C > NH - N (information theory) upper bound: C < 2In2NH + N (Burge analysis, Melhorn bound)

No comparison-based algorithm can do better.

Conjecture: With sampling,  $C/N \rightarrow H$  as sample size increases.



< 🗗 ▶

#### Quicksort is optimal

The average number of compares per element C/N is always within a constant factor of the entropy H

lower bound: C > NH - N (information theory) upper bound: C < 2In2NH + N (Burge analysis, Melhorn bound)

No comparison-based algorithm can do better.

Conjecture: With sampling,  $C/N \rightarrow H$  as sample size increases.

\* subject to some assumptions



**Quicksort and Search Trees** 



Saturated Fringe-Balanced Trees



**Back to Multiset Permutations** 

### **Classic Fact:**

• Recursion Tree of Quicksort = Naturally grown BST from input

→ Comparisons in Quicksort = Comparisons to **built** BST

- = Comparisons to **search** input in **final** BST
- How about inputs with duplicates?

< 17 >

### **Classic Fact:**

- Recursion Tree of Quicksort = Naturally grown BST from input
  - → Comparisons in Quicksort = Comparisons to **built** BST
    - = Comparisons to search input in final BST
- How about inputs with duplicates?

< @ >

### **Classic Fact:**

- Recursion Tree of Quicksort = Naturally grown BST from input
  - → Comparisons in Quicksort = Comparisons to **built** BST
    - = Comparisons to **search** input in **final** BST
- How about inputs with duplicates?

### Classic Fact: (without duplicates)

- Recursion Tree of Quicksort = Naturally grown BST from input
  - → Comparisons in Quicksort = Comparisons to **built** BST
    - = Comparisons to **search** input in **final** BST
- How about inputs with duplicates?

### Classic Fact: (without duplicates)

- Recursion Tree of Quicksort = Naturally grown BST from input
  - → Comparisons in Quicksort = Comparisons to **built** BST
    - = Comparisons to **search** input in **final** BST
- How about inputs with duplicates?

Quicksort (Fat-Pivot)

4 2 1 3 3 5 4 4 3 5 2

### Classic Fact: (without duplicates)

- Recursion Tree of Quicksort = Naturally grown BST from input
  - → Comparisons in Quicksort = Comparisons to **built** BST
    - = Comparisons to **search** input in **final** BST
- How about inputs with duplicates?

Quicksort (Fat-Pivot)

### Classic Fact: (without duplicates)

- Recursion Tree of Quicksort = Naturally grown BST from input
  - → Comparisons in Quicksort = Comparisons to **built** BST
    - = Comparisons to **search** input in **final** BST
- How about inputs with duplicates?

#### Quicksort (Fat-Pivot)



### Classic Fact: (without duplicates)

- Recursion Tree of Quicksort = Naturally grown BST from input
  - → Comparisons in Quicksort = Comparisons to **built** BST
    - = Comparisons to **search** input in **final** BST
- How about inputs with duplicates?

Quicksort (Fat-Pivot)



### Classic Fact: (without duplicates)

- Recursion Tree of Quicksort = Naturally grown BST from input
  - → Comparisons in Quicksort = Comparisons to **built** BST
    - = Comparisons to **search** input in **final** BST
- How about inputs with duplicates?

Quicksort (Fat-Pivot)



### Classic Fact: (without duplicates)

- Recursion Tree of Quicksort = Naturally grown BST from input
  - → Comparisons in Quicksort = Comparisons to **built** BST
    - = Comparisons to **search** input in **final** BST
- How about inputs with duplicates?





### Classic Fact: (without duplicates)

- Recursion Tree of Quicksort = Naturally grown BST from input
  - → Comparisons in Quicksort = Comparisons to **built** BST
    - = Comparisons to **search** input in **final** BST
- How about inputs with duplicates?

#### Quicksort (Fat-Pivot)



16

### Classic Fact: (without duplicates)

- Recursion Tree of Quicksort = Naturally grown BST from input
  - → Comparisons in Quicksort = Comparisons to **built** BST
    - = Comparisons to **search** input in **final** BST
- How about inputs with duplicates?





### Classic Fact: (without duplicates)

- Recursion Tree of Quicksort = Naturally grown BST from input
  - → Comparisons in Quicksort = Comparisons to **built** BST
    - = Comparisons to **search** input in **final** BST
- How about inputs with duplicates?



#### Quicksort (Fat-Pivot)

### Classic Fact: (without duplicates)

- Recursion Tree of Quicksort = Naturally grown BST from input
  - → Comparisons in Quicksort = Comparisons to **built** BST
    - = Comparisons to **search** input in **final** BST
- How about inputs with duplicates?



#### Quicksort (Fat-Pivot)
### Classic Fact: (without duplicates)

- Recursion Tree of Quicksort = Naturally grown BST from input
  - → Comparisons in Quicksort = Comparisons to **built** BST
    - = Comparisons to **search** input in **final** BST
- How about inputs with duplicates?



#### Quicksort (Fat-Pivot)

**Binary Search Tree** 

4 2 1 3 3 5 4 4 3 5 2

### Classic Fact: (without duplicates)

- Recursion Tree of Quicksort = Naturally grown BST from input
  - → Comparisons in Quicksort = Comparisons to **built** BST
    - = Comparisons to **search** input in **final** BST
- How about inputs with duplicates?





Binary Search Tree



#### Classic Fact: (without duplicates)

- Recursion Tree of Quicksort = Naturally grown BST from input
  - → Comparisons in Quicksort = Comparisons to **built** BST
    - = Comparisons to **search** input in **final** BST
- How about inputs with duplicates?





**Binary Search Tree** 



#### Classic Fact: (without duplicates)

- Recursion Tree of Quicksort = Naturally grown BST from input
  - → Comparisons in Quicksort = Comparisons to **built** BST
    - = Comparisons to **search** input in **final** BST
- How about inputs with duplicates?

3 3

3

3

2

3 3 3 2

2



#### Quicksort (Fat-Pivot)

5 4 4

3 3

3 3

4

#### Classic Fact: (without duplicates)

- Recursion Tree of Quicksort = Naturally grown BST from input
  - → Comparisons in Quicksort = Comparisons to **built** BST
    - = Comparisons to **search** input in **final** BST
- How about inputs with duplicates?







#### Classic Fact: (without duplicates)

- Recursion Tree of Quicksort = Naturally grown BST from input
  - → Comparisons in Quicksort = Comparisons to **built** BST
    - = Comparisons to **search** input in **final** BST
- How about inputs with duplicates?



Quicksort (Fat-Pivot)

**Binary Search Tree** 



#### Classic Fact: (without duplicates)

- Recursion Tree of Quicksort = Naturally grown BST from input
  - → Comparisons in Quicksort = Comparisons to **built** BST

4

- = Comparisons to **search** input in **final** BST
- How about inputs with duplicates?

3 3 5 4 4

3

3



#### Quicksort (Fat-Pivot)

3 3

3 3

**Binary Search Tree** 

2

3 3 3 2

2

#### Classic Fact: (without duplicates)

- Recursion Tree of Quicksort = Naturally grown BST from input
  - → Comparisons in Quicksort = Comparisons to **built** BST

- = Comparisons to **search** input in **final** BST
- How about inputs with duplicates?



#### Quicksort (Fat-Pivot)

3 3

#### Classic Fact: (without duplicates)

- Recursion Tree of Quicksort = Naturally grown BST from input
  - → Comparisons in Quicksort = Comparisons to **built** BST

3

- = Comparisons to **search** input in **final** BST
- How about inputs with duplicates?

3

3 3 5 4 4

3

3



#### Quicksort (Fat-Pivot)

2

3 3

3 3

2

3 3

2

#### Classic Fact: (without duplicates)

- Recursion Tree of Quicksort = Naturally grown BST from input
  - → Comparisons in Quicksort = Comparisons to **built** BST
    - = Comparisons to **search** input in **final** BST
- How about inputs with duplicates?

3

3 3

3 3

3



#### Quicksort (Fat-Pivot)

3 3

2

3 3

2

#### Classic Fact: (without duplicates)

- Recursion Tree of Quicksort = Naturally grown BST from input
  - → Comparisons in Quicksort = Comparisons to **built** BST
    - = Comparisons to **search** input in **final** BST
- How about inputs with duplicates?



#### Classic Fact: (without duplicates)

- Recursion Tree of Quicksort = Naturally grown BST from input
  - → Comparisons in Quicksort = Comparisons to **built** BST
    - = Comparisons to **search** input in **final** BST
- How about inputs with duplicates?

Quicksort (Fat-Pivot)



→ Equivalence holds also with duplicates.



#### Classic Fact: (without duplicates)

- Recursion Tree of Quicksort = Naturally grown BST from input
  - → Comparisons in Quicksort = Comparisons to **built** BST
    - = Comparisons to **search** input in **final** BST
- How about inputs with duplicates?



**Quicksort** (Fat-Pivot)

**Binary Search Tree** 



→ Equivalence holds also with duplicates.

2017-06-19

#### Classic Fact: (without duplicates)

- Recursion Tree of Quicksort = Naturally grown BST from input
  - → Comparisons in Quicksort = Comparisons to **built** BST
    - = Comparisons to **search** input in **final** BST
- How about inputs with duplicates?



**Quicksort** (Fat-Pivot)



### $\rightsquigarrow\,$ Equivalence holds also with duplicates.

This was only basic Quicksort ... how about pivot sampling?

Sebastian Wild

#### Classic Fact: (without duplicates)

- Recursion Tree of Quicksort = Naturally grown BST from input
  - → Comparisons in Quicksort = Comparisons to **built** BST
    - = Comparisons to **search** input in **final** BST
- How about inputs with duplicates?



→ Equivalence holds also with duplicates.

This was only basic Quicksort ... how about pivot sampling?

Sebastian Wild

### k-Fringe-Balanced Search Trees:

- Leaves **buffer** k = 2t + 1 elements.
- If buffer is full, leaf is **split**  $\rightarrow$  new internal node with **chosen pivot**.

### k-Fringe-Balanced Search Trees:

- Leaves **buffer** k = 2t + 1 elements.
- If buffer is full, leaf is **split**  $\rightsquigarrow$  new internal node with **chosen pivot**.

#### k-Fringe-Balanced Search Trees:

- Leaves **buffer** k = 2t + 1 elements.
- If buffer is full, leaf is **split**  $\rightsquigarrow$  new internal node with **chosen pivot**.

4 2 1 3 3 5 4 4 3 5 2

### k-Fringe-Balanced Search Trees:

- Leaves **buffer** k = 2t + 1 elements.
- If buffer is full, leaf is **split**  $\rightsquigarrow$  new internal node with **chosen pivot**.

 $\label{eq:median-of-5} \mbox{ Quicksort} \qquad (t=2)$ 

4 2 1 3 3 5 4 4 3 5 2

4

### k-Fringe-Balanced Search Trees:

- Leaves **buffer** k = 2t + 1 elements.
- If buffer is full, leaf is **split**  $\rightsquigarrow$  new internal node with **chosen pivot**.

 Median-of-5 Quicksort
 (t = 2) 

 2
 1
 3
 5
 4
 4
 3
 5
 2

### k-Fringe-Balanced Search Trees:

- Leaves **buffer** k = 2t + 1 elements.
- If buffer is full, leaf is **split**  $\rightsquigarrow$  new internal node with **chosen pivot**.

 Median-of-5 Quicksort
 (t = 2) 

 4
 2
 1
 3
 5
 4
 4
 3
 5
 2

 2
 1
 2
 3
 3
 4
 5
 4
 4
 5

### k-Fringe-Balanced Search Trees:

- Leaves **buffer** k = 2t + 1 elements.
- If buffer is full, leaf is **split**  $\rightsquigarrow$  new internal node with **chosen pivot**.

 Median-of-5 Quicksort
 (t = 2) 

 4
 2
 1
 3
 5
 4
 4
 3
 5
 2

 2
 1
 2
 3
 3
 4
 5
 4
 4
 5
 2

 Insertionsort
 Insert

<∂>

### k-Fringe-Balanced Search Trees:

- Leaves **buffer** k = 2t + 1 elements.
- If buffer is full, leaf is **split**  $\rightsquigarrow$  new internal node with **chosen pivot**.

 Median-of-5 Quicksort
 (t = 2) 

 4
 2
 1
 3
 3
 5
 4
 4
 3
 5
 2

 2
 1
 2
 3
 3
 4
 5
 4
 4
 5
 2

### k-Fringe-Balanced Search Trees:

- Leaves **buffer** k = 2t + 1 elements.
- If buffer is full, leaf is **split**  $\rightsquigarrow$  new internal node with **chosen pivot**.

 Median-of-5 Quicksort
 (t = 2) 

 4
 2
 1
 3
 3
 5
 4
 4
 3
 5
 2

 2
 1
 2
 3
 3
 3
 4
 5
 4
 4
 5

### k-Fringe-Balanced Search Trees:

- Leaves **buffer** k = 2t + 1 elements.
- If buffer is full, leaf is **split**  $\rightsquigarrow$  new internal node with **chosen pivot**.

Median-of-5 Quicksort (t = 2)4 3 5 2 2 3 3 5 4 4 2 3 4 5 4 5 2 1 3 3 (4) 4 4 5 5

### k-Fringe-Balanced Search Trees:

- Leaves **buffer** k = 2t + 1 elements.
- If buffer is full, leaf is **split**  $\rightsquigarrow$  new internal node with **chosen pivot**.

Median-of-5 Quicksort (t = 2)5-Fringe-Balanced Tree 2 3 3 5 4 4 3 5 2 4 2 5 4 5 2 1 3 3 3 4 4 4 4 Insertionsort

#### k-Fringe-Balanced Search Trees:

- Leaves **buffer** k = 2t + 1 elements.
- If buffer is full, leaf is **split**  $\rightsquigarrow$  new internal node with **chosen pivot**.



< 🗗 ►

#### k-Fringe-Balanced Search Trees:

- Leaves **buffer** k = 2t + 1 elements.
- If buffer is full, leaf is **split**  $\rightsquigarrow$  new internal node with **chosen pivot**.



#### k-Fringe-Balanced Search Trees:

- Leaves **buffer** k = 2t + 1 elements.
- If buffer is full, leaf is **split**  $\rightsquigarrow$  new internal node with **chosen pivot**.



∢ 🗗 🕨

#### k-Fringe-Balanced Search Trees:

- Leaves **buffer** k = 2t + 1 elements.
- If buffer is full, leaf is **split**  $\rightsquigarrow$  new internal node with **chosen pivot**.





< Ø →

#### k-Fringe-Balanced Search Trees:

- Leaves **buffer** k = 2t + 1 elements.
- If buffer is full, leaf is **split**  $\rightsquigarrow$  new internal node with **chosen pivot**.

Median-of-5 Quicksort 3 2 3 5 4 4 3 5 2 2 3\_3 3 5 2 1 4 4 5 4 212 4 4 5 5 4 55 5-Fringe-Balanced Tree 1 3 3 5 4 4 3 5 2 4 2 1

(t = 2)

#### k-Fringe-Balanced Search Trees:

- Leaves **buffer** k = 2t + 1 elements.
- If buffer is full, leaf is **split**  $\rightsquigarrow$  new internal node with **chosen pivot**.

Median-of-5 Quicksort 2 3 3 5 4 4 3 5 2 2 3\_3 5 2 1 . 3 4 4 5 4 . 5 212 4 4 5 4 55 **5-Fringe-Balanced Tree** 3 3 5 4 4 3 5 2 4 2 1 3

(t = 2)

#### k-Fringe-Balanced Search Trees:

- Leaves **buffer** k = 2t + 1 elements.
- If buffer is full, leaf is **split**  $\rightsquigarrow$  new internal node with **chosen pivot**.





∢ 🗗 🕨

#### k-Fringe-Balanced Search Trees:

- Leaves **buffer** k = 2t + 1 elements.
- If buffer is full, leaf is **split**  $\rightsquigarrow$  new internal node with **chosen pivot**.



#### k-Fringe-Balanced Search Trees:

- Leaves **buffer** k = 2t + 1 elements.
- If buffer is full, leaf is **split**  $\rightsquigarrow$  new internal node with **chosen pivot**.

Median-of-5 Quicksort (t = 2)5-Fringe-Balanced Tree 3 2 3 5 4 4 3 5 2 5 4 4 3 5 2 2 3 5 2 1 3 3 4 4 5 4 212 4 5 5 4 4 55

#### k-Fringe-Balanced Search Trees:

- Leaves **buffer** k = 2t + 1 elements.
- If buffer is full, leaf is **split**  $\rightsquigarrow$  new internal node with **chosen pivot**.



16
#### k-Fringe-Balanced Search Trees:

- Leaves **buffer** k = 2t + 1 elements.
- If buffer is full, leaf is **split**  $\rightsquigarrow$  new internal node with **chosen pivot**.



16

#### k-Fringe-Balanced Search Trees:

- Leaves **buffer** k = 2t + 1 elements.
- If buffer is full, leaf is **split**  $\rightsquigarrow$  new internal node with **chosen pivot**.



6

#### k-Fringe-Balanced Search Trees:

- Leaves **buffer** k = 2t + 1 elements.
- If buffer is full, leaf is **split**  $\rightsquigarrow$  new internal node with **chosen pivot**.



6

#### k-Fringe-Balanced Search Trees:

- Leaves **buffer** k = 2t + 1 elements.
- If buffer is full, leaf is **split**  $\rightsquigarrow$  new internal node with **chosen pivot**.



16

#### k-Fringe-Balanced Search Trees:

- Leaves **buffer** k = 2t + 1 elements.
- If buffer is full, leaf is **split**  $\rightsquigarrow$  new internal node with **chosen pivot**.



<₽>

#### k-Fringe-Balanced Search Trees:

- Leaves **buffer** k = 2t + 1 elements.
- If buffer is full, leaf is **split**  $\rightsquigarrow$  new internal node with **chosen pivot**.



<∂>

#### k-Fringe-Balanced Search Trees:

- Leaves **buffer** k = 2t + 1 elements.
- If buffer is full, leaf is **split**  $\rightsquigarrow$  new internal node with **chosen pivot**.



#### k-Fringe-Balanced Search Trees:

- Leaves **buffer** k = 2t + 1 elements.
- If buffer is full, leaf is **split**  $\rightsquigarrow$  new internal node with **chosen pivot**.



.

### k-Fringe-Balanced Search Trees:

- Leaves **buffer** k = 2t + 1 elements.
- If buffer is full, leaf is **split**  $\rightsquigarrow$  new internal node with **chosen pivot**.



→ Correspondence extends to

- Pivot Sampling (any scheme, not only median)
- (s-way Partitioning ~> s-ary search trees) <--- not toda</li>

→ Analyze search trees instead of Quicksort.

### k-Fringe-Balanced Search Trees:

- Leaves **buffer** k = 2t + 1 elements.
- If buffer is full, leaf is **split**  $\rightsquigarrow$  new internal node with **chosen pivot**.



### → Correspondence extends to

- Pivot Sampling (any scheme, not only median)
- (s-way Partitioning ~> s-ary search trees) ~\_\_\_\_\_ not today

→ Analyze search trees instead of Quicksort.

### k-Fringe-Balanced Search Trees:

- Leaves **buffer** k = 2t + 1 elements.
- If buffer is full, leaf is **split**  $\rightsquigarrow$  new internal node with **chosen pivot**.



### → Correspondence extends to

- Pivot Sampling (any scheme, not only median)
- (s-way Partitioning ~> s-ary search trees) ~\_\_\_\_\_ not today

→ Analyze search trees instead of Quicksort.

### k-Fringe-Balanced Search Trees:

- Leaves **buffer** k = 2t + 1 elements.
- If buffer is full, leaf is **split**  $\rightsquigarrow$  new internal node with **chosen pivot**.



#### → Correspondence extends to

- Pivot Sampling (any scheme, not only median)

### → Analyze search trees instead of Quicksort.

### k-Fringe-Balanced Search Trees:

- Leaves **buffer** k = 2t + 1 elements.
- If buffer is full, leaf is **split**  $\rightsquigarrow$  new internal node with **chosen pivot**.



#### → Correspondence extends to

- Pivot Sampling (any scheme, not only median)
- (s-way Partitioning → s-ary search trees) ← not today

→ Analyze search trees instead of Quicksort.

### $\rightsquigarrow$ Quicksort costs = costs to **search** input $\vec{u} = (u_1, \dots, u_n)$ in **final** tree T.

- $\mathbb T$  fixed  $\, \rightsquigarrow \,$  search cost depends only on profile  $\vec X = (X_1, \ldots, X_u)$
- **but:** T also depends on  $\vec{U}$  (Recall: T is built from  $\vec{U}$ !)
- → direct analysis no simpler than for Quicksort

(k = 1)

4 5 5 5 2 3 5 4 2 2 2 5 5 5 4 2 5 1 5 4 2 5 1 5 4 2 3 4 2 2 1 2 4 2 2 5 5 2 3 4 2 3 4 1 3 2 2 1 4 4 4 1 2 3 3

**Observation:** T becomes **stationary** after each value was inserted!

- **Q-** Split input into **tree-growing** part and **searching** part:
  - **1** We built  $\mathcal{T}$  until it is stationary, **ignoring costs**.
  - 2 Determine costs of **searching** remaining elements.
  - $\rightsquigarrow$  profile  $\vec{X}_S$  of search part **independent** of  ${\mathbb T}$

 $\rightsquigarrow$  Quicksort costs = costs to **search** input  $\vec{u} = (u_1, \dots, u_n)$  in **final** tree T.

- ${\mathfrak T}$  fixed  $\, \rightsquigarrow \,$  search cost depends only on profile  $\vec{X} = (X_1, \ldots, X_u)$
- **but:** T also depends on  $\vec{U}$  (Recall: T is built from  $\vec{U}$ !)
- → direct analysis no simpler than for Quicksort

(k = 1)

4 5 5 5 2 3 5 4 2 2 2 5 5 5 4 2 5 5 5 4 2 5 1 5 4 2 5 1 5 4 2 3 4 2 2 1 2 4 2 2 5 5 2 3 4 2 3 4 1 3 2 2 1 4 4 4 1 2 3 3

**Observation:** T becomes **stationary** after each value was inserted!

- **Q-** Split input into **tree-growing** part and **searching** part:
  - **1** We built  $\mathcal{T}$  until it is stationary, **ignoring costs**.
  - 2 Determine costs of **searching** remaining elements.
  - $\rightsquigarrow$  profile  $\vec{X}_S$  of search part **independent** of  ${\mathbb T}$

 $\rightsquigarrow$  Quicksort costs = costs to **search** input  $\vec{u} = (u_1, \dots, u_n)$  in **final** tree T.

- ${\mathfrak T}$  fixed  $\,\rightsquigarrow\,$  search cost depends only on profile  $\vec{X}=(X_1,\ldots,X_u)$
- **but:** T also depends on  $\vec{u}$  (Recall: T is built from  $\vec{u}$ !)

→ direct analysis no simpler than for Quicksort

(k = 1)

4 5 5 5 2 3 5 4 2 2 2 5 5 5 4 2 5 1 5 4 2 5 1 5 4 2 3 4 2 2 1 2 4 2 2 5 5 2 3 4 2 3 4 1 3 2 2 1 4 4 4 1 2 3 3

**Observation:** T becomes **stationary** after each value was inserted!

- **Q-** Split input into **tree-growing** part and **searching** part:
  - **1** We built T until it is stationary, **ignoring costs**.
  - 2 Determine costs of **searching** remaining elements.
  - $\rightsquigarrow$  profile  $\vec{X}_S$  of search part **independent** of  ${\mathbb T}$

 $\rightsquigarrow$  Quicksort costs = costs to **search** input  $\vec{u} = (u_1, \dots, u_n)$  in **final** tree T.

- ${\mathfrak T}$  fixed  $\, \rightsquigarrow \,$  search cost depends only on profile  $\vec{X} = (X_1, \ldots, X_u)$
- **but:** T also depends on  $\vec{u}$  (Recall: T is built from  $\vec{u}$ !)
- → direct analysis no simpler than for Quicksort

(k = 1)

#### 4 5 5 5 2 3 5 4 2 2 2 5 5 5 4 2 5 1 5 4 2 5 1 5 4 2 3 4 2 3 4 2 2 1 2 4 2 5 5 2 3 4 2 3 4 1 3 2 2 1 4 4 4 1 2 3 3

**Observation:** T becomes **stationary** after each value was inserted!

- Split input into **tree-growing** part and **searching** part:
  - **1** We built  $\mathcal{T}$  until it is stationary, **ignoring costs**.
  - 2 Determine costs of **searching** remaining elements.
  - $\rightsquigarrow$  profile  $ec{X}_S$  of search part **independent** of  ${\mathbb T}$

 $\rightsquigarrow$  Quicksort costs = costs to **search** input  $\vec{u} = (u_1, \dots, u_n)$  in **final** tree T.

- ${\mathfrak T}$  fixed  $\, \rightsquigarrow \,$  search cost depends only on profile  $\vec{X} = (X_1, \ldots, X_u)$
- **but:** T also depends on  $\vec{u}$  (Recall: T is built from  $\vec{u}$ !)
- → direct analysis no simpler than for Quicksort

(k = 1)

#### 4 5 5 5 2 3 5 4 2 2 2 5 5 5 4 2 5 1 5 4 2 3 4 2 2 1 2 4 2 2 5 5 2 3 4 2 3 4 1 3 2 2 1 4 4 4 1 2 3 3

**Observation:** T becomes **stationary** after each value was inserted!

- Split input into **tree-growing** part and **searching** part:
  - **1** We built  $\mathcal{T}$  until it is stationary, **ignoring costs**.
  - 2 Determine costs of **searching** remaining elements.
  - $\rightsquigarrow$  profile  $\vec{X}_S$  of search part **independent** of  ${\mathbb T}$

 $\rightsquigarrow$  Quicksort costs = costs to **search** input  $\vec{u} = (u_1, \dots, u_n)$  in **final** tree T.

- ${\mathfrak T}$  fixed  $\, \rightsquigarrow \,$  search cost depends only on profile  $\vec{X} = (X_1, \ldots, X_u)$
- **but:** T also depends on  $\vec{u}$  (Recall: T is built from  $\vec{u}$ !)
- → direct analysis no simpler than for Quicksort

(k = 1)

#### 4 5 5 5 2 2 5 5 4 2 5 1 5 4 2 2 1 5 4 2 2 1 5 4 2 2 5 5 2 3 4 1 3 2 2 1 4 4 1 2 3 3

**Observation:** T becomes **stationary** after each value was inserted!

- Split input into **tree-growing** part and **searching** part:
  - **1** We built  $\mathcal{T}$  until it is stationary, **ignoring costs**.
  - 2 Determine costs of **searching** remaining elements.
  - $\rightsquigarrow$  profile  $\vec{X}_S$  of search part **independent** of  ${\mathbb T}$

٠

 $\rightsquigarrow$  Quicksort costs = costs to **search** input  $\vec{u} = (u_1, \dots, u_n)$  in **final** tree T.

- ${\mathfrak T}$  fixed  $\,\rightsquigarrow\,$  search cost depends only on profile  $\vec{X}=(X_1,\ldots,X_u)$
- **but:** T also depends on  $\vec{u}$  (Recall: T is *built* from  $\vec{u}$ !)
- → direct analysis no simpler than for Quicksort

(k = 1)

### 4 5 5 5 2 3 5 4 2 2 2 5 5 5 4 2 5 1 5 4 2 3 4 2 2 1 2 4 2 2 5 5 2 3 4 2 3 4 1 3 2 2 1 4 4 4 1 2 3 3

#### **Observation:** T becomes **stationary** after each value was inserted!

- Split input into **tree-growing** part and **searching** part:
  - **1** We built  $\mathcal{T}$  until it is stationary, **ignoring costs**.
  - 2 Determine costs of **searching** remaining elements.
  - $\rightsquigarrow$  profile  $ec{X}_S$  of search part **independent** of  ${\mathbb T}$

• •

 $\rightsquigarrow$  Quicksort costs = costs to **search** input  $\vec{u} = (u_1, \dots, u_n)$  in **final** tree T.

- ${\mathfrak T}$  fixed  $\,\rightsquigarrow\,$  search cost depends only on profile  $\vec{X}=(X_1,\ldots,X_u)$
- **but:** T also depends on  $\vec{u}$  (Recall: T is *built* from  $\vec{u}$ !)
- → direct analysis no simpler than for Quicksort

(k = 1)

### 4 5 5 5 2 3 5 4 2 2 2 5 5 5 4 2 5 1 5 4 2 3 4 2 2 1 2 4 2 2 5 5 2 3 4 2 3 4 1 3 2 2 1 4 4 4 1 2 3 3 · \ \

### **Observation:** T becomes **stationary** after each value was inserted!

- Split input into **tree-growing** part and **searching** part:
  - **1** We built  $\mathcal{T}$  until it is stationary, **ignoring costs**.
  - 2 Determine costs of **searching** remaining elements.
  - $\rightsquigarrow$  profile  $\vec{X}_S$  of search part **independent** of  ${\mathbb T}$

 $\rightsquigarrow$  Quicksort costs = costs to **search** input  $\vec{u} = (u_1, \dots, u_n)$  in **final** tree T.

- ${\mathfrak T}$  fixed  $\,\rightsquigarrow\,$  search cost depends only on profile  $\vec{X}=(X_1,\ldots,X_u)$
- **but:** T also depends on  $\vec{u}$  (Recall: T is *built* from  $\vec{u}$ !)
- → direct analysis no simpler than for Quicksort

(k = 1)

### 4 5 5 5 2 3 5 4 2 2 2 5 5 5 4 2 5 1 5 4 2 3 4 2 2 1 2 4 2 2 5 5 2 3 4 2 3 4 1 3 2 2 1 4 4 4 1 2 3 3 • \$

### **Observation:** T becomes **stationary** after each value was inserted!

- Split input into **tree-growing** part and **searching** part:
  - **1** We built T until it is stationary, **ignoring costs**.
  - 2 Determine costs of **searching** remaining elements.
  - $\rightsquigarrow$  profile  $\vec{X}_S$  of search part **independent** of  ${\mathbb T}$

 $\rightsquigarrow$  Quicksort costs = costs to **search** input  $\vec{u} = (u_1, \dots, u_n)$  in **final** tree T.

- ${\mathfrak T}$  fixed  $\, \rightsquigarrow \,$  search cost depends only on profile  $\vec{X} = (X_1, \ldots, X_u)$
- **but:** T also depends on  $\vec{u}$  (Recall: T is built from  $\vec{u}$ !)
- → direct analysis no simpler than for Quicksort

(k = 1)

### 4 5 5 5 2 3 5 4 2 2 2 5 5 5 4 2 5 1 5 4 2 3 4 2 2 1 2 4 2 2 5 5 2 3 4 2 3 4 1 3 2 2 1 4 4 4 1 2 3 3 • \ \ \ \

### **Observation:** T becomes **stationary** after each value was inserted!

- Split input into **tree-growing** part and **searching** part:
  - **1** We built  $\mathcal{T}$  until it is stationary, **ignoring costs**.
  - 2 Determine costs of **searching** remaining elements.
  - $\rightsquigarrow$  profile  $\vec{X}_S$  of search part **independent** of  ${\mathbb T}$

 $\rightsquigarrow$  Quicksort costs = costs to **search** input  $\vec{u} = (u_1, \dots, u_n)$  in **final** tree T.

- ${\mathfrak T}$  fixed  $\,\rightsquigarrow\,$  search cost depends only on profile  $\vec{X}=(X_1,\ldots,X_u)$
- **but:** T also depends on  $\vec{u}$  (Recall: T is *built* from  $\vec{u}$ !)
- → direct analysis no simpler than for Quicksort

(k = 1)

### 4 5 5 5 2 3 5 4 2 2 2 5 5 5 4 2 5 1 5 4 2 3 4 2 2 1 2 4 2 2 5 5 2 3 4 2 3 4 1 3 2 2 1 4 4 4 1 2 3 3 • \ \ \ \ \ \ \

### **Observation:** T becomes **stationary** after each value was inserted!

- Split input into **tree-growing** part and **searching** part:
  - **1** We built  $\mathcal{T}$  until it is stationary, **ignoring costs**.
  - 2 Determine costs of **searching** remaining elements.
  - $\rightsquigarrow$  profile  $\vec{X}_S$  of search part **independent** of  ${\mathbb T}$

 $\rightsquigarrow$  Quicksort costs = costs to **search** input  $\vec{u} = (u_1, \dots, u_n)$  in **final** tree T.

- ${\mathfrak T}$  fixed  $\,\rightsquigarrow\,$  search cost depends only on profile  $\vec{X}=(X_1,\ldots,X_u)$
- **but:** T also depends on  $\vec{u}$  (Recall: T is built from  $\vec{u}$ !)
- → direct analysis no simpler than for Quicksort

(k = 1)

### 4 5 5 5 2 3 5 4 2 2 2 5 5 5 4 2 5 1 5 4 2 3 4 2 2 1 2 4 2 2 5 5 2 3 4 2 3 4 1 3 2 2 1 4 4 4 1 2 3 3 • \ \ \ \ \ \ \ \ \ \ \ \ \ \ \

### **Observation:** T becomes **stationary** after each value was inserted!

- Split input into **tree-growing** part and **searching** part:
  - **1** We built T until it is stationary, **ignoring costs**.
  - 2 Determine costs of **searching** remaining elements.
  - $\rightsquigarrow$  profile  $\vec{X}_S$  of search part **independent** of  ${\mathbb T}$

 $\rightsquigarrow$  Quicksort costs = costs to **search** input  $\vec{u} = (u_1, \dots, u_n)$  in **final** tree T.

- ${\mathfrak T}$  fixed  $\, \rightsquigarrow \,$  search cost depends only on profile  $\vec{X} = (X_1, \ldots, X_u)$
- **but:** T also depends on  $\vec{u}$  (Recall: T is built from  $\vec{u}$ !)
- → direct analysis no simpler than for Quicksort

(k = 1)

### 4 5 5 5 2 3 5 4 2 2 2 5 5 5 4 2 5 1 5 4 2 3 4 2 2 1 2 4 2 2 5 5 2 3 4 2 3 4 1 3 2 2 1 4 4 4 1 2 3 3 • V V V A A A

### **Observation:** T becomes **stationary** after each value was inserted!

- Split input into **tree-growing** part and **searching** part:
  - **1** We built T until it is stationary, **ignoring costs**.
  - 2 Determine costs of **searching** remaining elements.
  - $\rightsquigarrow$  profile  $ec{X}_S$  of search part **independent** of  ${\mathbb T}$

 $\rightsquigarrow$  Quicksort costs = costs to **search** input  $\vec{u} = (u_1, \dots, u_n)$  in **final** tree T.

- ${\mathfrak T}$  fixed  $\, \rightsquigarrow \,$  search cost depends only on profile  $\vec{X} = (X_1, \ldots, X_u)$
- **but:** T also depends on  $\vec{u}$  (Recall: T is built from  $\vec{u}$ !)
- → direct analysis no simpler than for Quicksort

(k = 1)

### 

### **Observation:** T becomes **stationary** after each value was inserted!

- Split input into **tree-growing** part and **searching** part:
  - We built T until it is stationary, **ignoring costs**.
  - 2 Determine costs of **searching** remaining elements.
  - $\rightsquigarrow$  profile  $\vec{X}_S$  of search part **independent** of  ${\mathbb T}$

 $\rightsquigarrow$  Quicksort costs = costs to **search** input  $\vec{u} = (u_1, \dots, u_n)$  in **final** tree T.

- ${\mathfrak T}$  fixed  $\, \rightsquigarrow \,$  search cost depends only on profile  $\vec{X} = (X_1, \ldots, X_u)$
- **but:** T also depends on  $\vec{u}$  (Recall: T is built from  $\vec{u}$ !)
- → direct analysis no simpler than for Quicksort

(k = 1)

### 

### **Observation:** T becomes **stationary** after each value was inserted!

- Split input into **tree-growing** part and **searching** part:
  - **1** We built T until it is stationary, **ignoring costs**.
  - 2 Determine costs of **searching** remaining elements.
  - $\rightsquigarrow$  profile  $\vec{X}_S$  of search part **independent** of  ${\mathbb T}$

 $\rightsquigarrow$  Quicksort costs = costs to **search** input  $\vec{u} = (u_1, \dots, u_n)$  in **final** tree T.

- ${\mathfrak T}$  fixed  $\, \rightsquigarrow \,$  search cost depends only on profile  $\vec{X} = (X_1, \ldots, X_u)$
- **but:** T also depends on  $\vec{u}$  (Recall: T is *built* from  $\vec{u}$ !)
- → direct analysis no simpler than for Quicksort

(k = 1)

4 5 5 5 2 3 5 4 2 2 2 5 5 5 4 2 5 1 5 4 2 3 4 2 2 1 2 4 2 2 5 5 2 3 4 2 3 4 1 3 2 2 1 4 4 4 1 2 3 3

### **Observation:** T becomes **stationary** after each value was inserted!

- Split input into **tree-growing** part and **searching** part:
  - We built T until it is stationary, **ignoring costs**.
  - 2 Determine costs of **searching** remaining elements.
  - $\rightsquigarrow$  profile  $ec{X}_S$  of search part **independent** of  ${\mathbb T}$

 $\rightsquigarrow$  Quicksort costs = costs to **search** input  $\vec{u} = (u_1, \dots, u_n)$  in **final** tree T.

- ${\mathfrak T}$  fixed  $\rightsquigarrow$  search cost depends only on profile  $\vec{X}=(X_1,\ldots,X_u)$
- **but:** T also depends on  $\vec{U}$  (Recall: T is *built* from  $\vec{U}$ !)
- → direct analysis no simpler than for Quicksort

(k = 1)

4 5 5 5 2 3 5 4 2 2 2 5 5 5 4 2 5 1 5 4 2 3 4 2 2 1 2 4 2 2 5 5 2 3 4 2 3 4 1 3 2 2 1 4 4 4 1 2 3 3

### **Observation:** T becomes **stationary** after each value was inserted!

- Split input into **tree-growing** part and **searching** part:
  - **1** We built T until it is stationary, **ignoring costs**.
  - 2 Determine costs of **searching** remaining elements.
  - $\rightsquigarrow$  profile  $\vec{X}_S$  of search part **independent** of  ${\mathbb T}$

 $\rightsquigarrow$  Quicksort costs = costs to **search** input  $\vec{u} = (u_1, \dots, u_n)$  in **final** tree T.

- ${\mathfrak T}$  fixed  $\, \rightsquigarrow \,$  search cost depends only on profile  $\vec{X} = (X_1, \ldots, X_u)$
- **but:** T also depends on  $\vec{U}$  (Recall: T is *built* from  $\vec{U}$ !)
- → direct analysis no simpler than for Quicksort

(k = 1)

4 5 5 5 2 3 5 4 2 2 2 5 5 5 4 2 5 1 5 4 2 3 4 2 2 1 2 4 2 2 5 5 2 3 4 2 3 4 1 3 2 2 1 4 4 4 1 2 3 3

#### **Observation:** T becomes **stationary** after each value was inserted!

Fringe-balanced:  $\rightarrow$  stationary after each value inserted  $\mathbf{k} = 2\mathbf{t} + 1$  times (up to k duplicates in buffer)

• Split input into **tree-growing** part and **searching** part:

- **1** We built  $\mathcal{T}$  until it is stationary, **ignoring costs**.
- 2 Determine costs of **searching** remaining elements.
- $\rightsquigarrow$  profile  $ec{X}_S$  of search part **independent** of  ${\mathfrak T}$

 $\rightsquigarrow$  Quicksort costs = costs to **search** input  $\vec{u} = (u_1, \dots, u_n)$  in **final** tree T.

- ${\mathfrak T} \text{ fixed } \rightsquigarrow \text{ search cost depends only on } {\textbf{profile }} \vec{X} = (X_1, \ldots, X_u)$
- **but:** T also depends on  $\vec{U}$  (Recall: T is *built* from  $\vec{U}$ !)
- → direct analysis no simpler than for Quicksort

(k = 1)

**Observation:** T becomes **stationary** after each value was inserted!

Fringe-balanced:  $\rightarrow$  stationary after each value inserted k = 2t + 1 times (up to k duplicates in buffer)

Split input into **tree-growing** part and **searching** part:

- We built T until it is stationary, ignoring costs.
- 2 Determine costs of **searching** remaining elements.
- $\rightsquigarrow$  profile  $ec{X}_S$  of search part **independent** of  $\mathfrak T$

(日)

 $\rightsquigarrow$  Quicksort costs = costs to **search** input  $\vec{u} = (u_1, \dots, u_n)$  in **final** tree T.

- ${\mathfrak T} \text{ fixed } \rightsquigarrow \text{ search cost depends only on } {\textbf{profile }} \vec{X} = (X_1, \ldots, X_u)$
- **but:** T also depends on  $\vec{U}$  (Recall: T is *built* from  $\vec{U}$ !)
- → direct analysis no simpler than for Quicksort

(k = 1)



**Observation:** T becomes **stationary** after each value was inserted!<sup>\*</sup>

 $\rightarrow$  stationary after each value inserted k = 2t + 1 times (up to k duplicates in buffer)

- Split input into **tree-growing** part and **searching** part:
  - **1** We built T until it is stationary, **ignoring costs**.
  - 2 Determine costs of **searching** remaining elements.
  - $\rightsquigarrow$  profile  $\vec{X}_S$  of search part **independent** of  $\mathfrak{T}$

 $\rightsquigarrow$  Quicksort costs = costs to **search** input  $\vec{u} = (u_1, \dots, u_n)$  in **final** tree T.

- ${\mathfrak T} \text{ fixed } \rightsquigarrow \text{ search cost depends only on } {\textbf{profile }} \vec{X} = (X_1, \ldots, X_u)$
- **but:** T also depends on  $\vec{U}$  (Recall: T is *built* from  $\vec{U}$ !)
- → direct analysis no simpler than for Quicksort

(k = 1)



 $\rightsquigarrow$  Quicksort costs = costs to **search** input  $\vec{u} = (u_1, \dots, u_n)$  in **final** tree T.

- ${\mathfrak T} \text{ fixed } \rightsquigarrow \text{ search cost depends only on } {\textbf{profile }} \vec{X} = (X_1, \ldots, X_u)$
- **but:** T also depends on  $\vec{u}$  (Recall: T is *built* from  $\vec{u}$ !)
- → direct analysis no simpler than for Quicksort

(k = 1)



 $\rightsquigarrow$  Quicksort costs = costs to **search** input  $\vec{u} = (u_1, \dots, u_n)$  in **final** tree T.

- ${\mathfrak T}$  fixed  $\, \rightsquigarrow \,$  search cost depends only on profile  $\vec{X} = (X_1, \ldots, X_u)$
- **but:** T also depends on  $\vec{U}$  (Recall: T is *built* from  $\vec{U}$ !)
- → direct analysis no simpler than for Quicksort

(k = 1)


# Tree-growing and searching

 $\rightsquigarrow$  Quicksort costs = costs to **search** input  $\vec{u} = (u_1, \dots, u_n)$  in **final** tree T.

- ${\mathfrak T} \text{ fixed } \rightsquigarrow \text{ search cost depends only on } {\textbf{profile }} \vec{X} = (X_1, \ldots, X_u)$
- **but:** T also depends on  $\vec{u}$  (Recall: T is *built* from  $\vec{u}$ !)
- → direct analysis no simpler than for Quicksort

(k = 1)



# Tree-growing and searching

 $\rightsquigarrow$  Quicksort costs = costs to **search** input  $\vec{u} = (u_1, \dots, u_n)$  in **final** tree T.

- ${\mathfrak T} \text{ fixed } \rightsquigarrow \text{ search cost depends only on } {\textbf{profile }} \vec{X} = (X_1, \ldots, X_u)$
- **but:** T also depends on  $\vec{U}$  (Recall: T is *built* from  $\vec{U}$ !)
- → direct analysis no simpler than for Quicksort

(k = 1)



**Goal:** ignore tree-growing for analysis.

 $\rightsquigarrow$  Allow only first  $n_T$  elements for tree growing.

< 17 >

Goal: ignore tree-growing for analysis.



 $\rightsquigarrow$  Allow only first  $n_T$  elements for tree growing.

< @ >



16

< @ >



**Problem:** if a value occurs < k times in first  $n_T$  elements, T not complete





6









• Costs to grow  $\mathcal{T}$ :



• never more than  $\leq n_T \cdot u$ 





6



< (1) →



< (1) →



- $\alpha(\vec{q}) =$  expected search cost in random  $\mathcal{T}$
- ullet error term covers: tree-growing, sorting samples, degenerate inputs, inputs with high  ${\mathcal T}$

Sebastian Wild



- $\alpha(\vec{q}) =$  **expected search cost** in random  $\Upsilon$
- ullet error term covers: tree-growing, sorting samples, degenerate inputs, inputs with high  ${\mathfrak T}$

Sebastian Wild



- $\alpha(\vec{q}) =$  **expected search cost** in random  $\Upsilon$
- error term covers: tree-growing, sorting samples, degenerate inputs, inputs with high  ${\mathbb T}$

Sebastian Wild

2017-06-19



**Back to Multiset Permutations** 

**Recall:** 
$$\alpha(\vec{q}) = \sum_{\nu=1}^{u} q_{\nu} \cdot depth_{\mathcal{T}}(\nu)$$
  $\mathcal{T}$  from inserting i. i. d.  $\mathcal{D}(\vec{q})$  el

*Warmup:* Ordinary BSTs (t = 0)

Allen & Munro 1978 Self-Organizing Binary Search Tree  $\alpha(\vec{q}) = 2\mathcal{H}_Q(\vec{q}) + 1 \text{ with}$   $\mathcal{H}_Q(\vec{q}) = \sum_{1 \leq i < j \leq u} \frac{q_i q_j}{q_i + \dots + q_j}$ 

• Proof sketch:

Sum prob. that i is ancestor of j over all i, j ancestor  $\iff$  i **first** inserted key among i,..., j

**Recall:** 
$$\alpha(\vec{q}) = \sum_{\nu=1}^{u} q_{\nu} \cdot depth_{\mathcal{T}}(\nu)$$
  $\mathcal{T}$  from inserting i. i. d.  $\mathcal{D}(\vec{q})$  elements until

*Warmup:* Ordinary BSTs (t = 0)

Allen & Munro 1978 Self-Organizing Binary Search Tree  $\alpha(\vec{q}) = 2\mathcal{H}_Q(\vec{q}) + 1 \text{ with}$   $\mathcal{H}_Q(\vec{q}) = \sum_{1 \leq i < j \leq u} \frac{q_i q_j}{q_i + \dots + q_j}$ 

• Proof sketch:

Sum prob. that i is ancestor of j over all i, j ancestor  $\iff$  i first inserted key among i,..., j

< 17 >

saturation

**Recall:** 
$$\alpha(\vec{q}) = \sum_{\nu=1}^{u} q_{\nu} \cdot depth_{\mathcal{T}}(\nu)$$

distribution with prob. weights  $q_1, \ldots, q_u$ T from inserting i.i.d.  $\hat{\mathcal{D}}(\vec{q})$  elements **until saturation** 

*Warmup:* Ordinary BSTs (t = 0)

Allen & Munro 1978: Self-Organizing Binary Search Trees  $\alpha(\vec{q}) = 2\mathcal{H}_Q(\vec{q}) + 1 \text{ with}$   $\mathcal{H}_Q(\vec{q}) = \sum_{\substack{1 \leq i < j \leq u}} \frac{q_i q_j}{q_i + \dots + q_j}$ 

• Proof sketch:

Sum prob. that i is ancestor of j over all i, j ancestor  $\iff$  i **first** inserted key among i, ..., j

**Recall:** 
$$\alpha(\vec{q}) = \sum_{\nu=1}^{u} q_{\nu} \cdot depth_{\mathcal{T}}(\nu)$$

distribution with prob. weights  $q_1, \ldots, q_u$ T from inserting i.i.d.  $\hat{\mathcal{D}}(\vec{q})$  elements **until saturation** 

#### *Warmup:* Ordinary BSTs (t = 0)

• Proof sketch:

Sum prob. that i is ancestor of j over all i, j ancestor  $\iff$  i **first** inserted key among i,..., j

**Recall:** 
$$\alpha(\vec{q}) = \sum_{\nu=1}^{u} q_{\nu} \cdot depth_{\mathcal{T}}(\nu)$$
  $\mathcal{T}$  from inserting i.i.d.  $\mathcal{D}(\vec{q})$  elements until saturation

#### *Warmup:* Ordinary BSTs (t = 0)

 Proof sketch:
 Sum prob. that i is ancestor of j over all i, j ancestor ⇔ i first inserted key among i,..., j

**Recall:** 
$$\alpha(\vec{q}) = \sum_{\nu=1}^{u} q_{\nu} \cdot depth_{\mathcal{T}}(\nu)$$
  $\mathcal{T}$  from inserting i.i.d.  $\mathcal{D}(\vec{q})$  elements until saturation

#### *Warmup:* Ordinary BSTs (t = 0)

#### • Proof sketch:

Sum prob. that i is ancestor of j over all i, j ancestor  $\iff$  i **first** inserted key among i,..., j

< 67 ▶

**Recall:** 
$$\alpha(\vec{q}) = \sum_{\nu=1}^{u} q_{\nu} \cdot depth_{\mathfrak{T}}(\nu)$$

#### *Warmup:* Ordinary BSTs (t = 0)

 $\begin{array}{lll} & \text{Allen \& Munro 1978:} \\ & \text{Self-Organizing Binary Search Trees} \\ & \alpha(\vec{q}) &= 2\mathcal{H}_Q(\vec{q}) + 1 \text{ with} \\ & \mathcal{H}_Q(\vec{q}) = \sum_{1 \leqslant i < j \leqslant u} \frac{q_i q_j}{q_i + \cdots + q_j} \end{array}$ 

• Proof sketch:

Sum prob. that i is ancestor of j over all i, j ancestor  $\iff$  i **first** inserted key among i,..., j

#### distribution with prob. weights $q_1, \ldots, q_u$ T from inserting i.i.d. $\hat{\mathcal{D}}(\vec{q})$ elements **until saturation**

#### Fringe-balanced trees (t > 1)



 $\rightsquigarrow\,$  prob. that i inserted first among i, . . . j  $\ref{eq:selected_selecte$ 

< (1) →

**Recall:** 
$$\alpha(\vec{q}) = \sum_{\nu=1}^{u} q_{\nu} \cdot depth_{\mathfrak{T}}(\nu)$$

#### *Warmup:* Ordinary BSTs (t = 0)

 $\begin{array}{lll} & \text{Allen \& Munro 1978:} \\ & \text{Self-Organizing Binary Search Trees} \\ & \alpha(\vec{q}) &= 2\mathcal{H}_Q(\vec{q}) + 1 \text{ with} \\ & \mathcal{H}_Q(\vec{q}) = \sum_{1 \leqslant i < j \leqslant u} \frac{q_i q_j}{q_i + \cdots + q_j} \end{array}$ 

• Proof sketch:

Sum prob. that i is ancestor of j over all i, j ancestor  $\iff$  i **first** inserted key among i,..., j

#### distribution with prob. weights $q_1, \ldots, q_u$ T from inserting i.i.d. $\hat{\mathcal{D}}(\vec{q})$ elements **until saturation**

#### Fringe-balanced trees (t > 1)

• probability of given value in root:



 $\rightsquigarrow\,$  prob. that i inserted first among i, . . . j  $\ref{eq:selected_selecte$ 

**Recall:** 
$$\alpha(\vec{q}) = \sum_{\nu=1}^{u} q_{\nu} \cdot depth_{\mathcal{T}}(\nu)$$

#### *Warmup:* Ordinary BSTs (t = 0)

**Old result:** Allen & Munro 1978: Self-Organizing Binary Search Trees  $\alpha(\vec{q}) = 2\mathcal{H}_{O}(\vec{q}) + 1$  with  $\mathcal{H}_Q(\vec{q}) = \sum_{1 \le i < j \le u} \frac{q_i q_j}{q_i + \dots + q_j}$ 

Proof sketch:

Sum prob. that i is ancestor of j over all i, j ancestor  $\iff$  i **first** inserted key among i, ..., j

distribution with prob. weights  $q_1, \ldots, q_u$  $\mathfrak{T}$  from inserting i.i.d.  $\mathfrak{D}(\mathbf{q})$  elements until saturation

#### Fringe-balanced trees (t > 1)

• probability of given value in root:



 $\rightarrow$  prob. that i inserted first among i, ... j ??

**Recall:** 
$$\alpha(\vec{q}) = \sum_{\nu=1}^{u} q_{\nu} \cdot depth_{\mathfrak{T}}(\nu)$$

#### *Warmup:* Ordinary BSTs (t = 0)

• Proof sketch:

Sum prob. that i is ancestor of j over all i, j ancestor  $\iff$  i **first** inserted key among i,..., j

#### distribution with prob. weights $q_1, \ldots, q_u$ T from inserting i.i.d. $\hat{\mathcal{D}}(\vec{q})$ elements **until saturation**

#### Fringe-balanced trees (t > 1)

• probability of given value in root:





**Recall:** 
$$\alpha(\vec{q}) = \sum_{\nu=1}^{u} q_{\nu} \cdot depth_{\mathfrak{T}}(\nu)$$

#### *Warmup:* Ordinary BSTs (t = 0)

 $\begin{array}{lll} & \text{Allen \& Munro 1978:} \\ & \text{Self-Organizing Binary Search Trees} \\ & \alpha(\vec{q}) &= 2\mathcal{H}_Q(\vec{q}) + 1 \text{ with} \\ & \mathcal{H}_Q(\vec{q}) = \sum_{1 \leqslant i < j \leqslant u} \frac{q_i q_j}{q_i + \cdots + q_j} \end{array}$ 

• Proof sketch:

Sum prob. that i is ancestor of j over all i, j ancestor  $\iff$  i **first** inserted key among i,...,j

 $\bullet~$  In the same paper:  $\mathfrak{H}_Q(\vec{q}) < \mathfrak{H}_{\text{ln}}(\vec{q})$ 

#### distribution with prob. weights $q_1, \ldots, q_u$ T from inserting i.i.d. $\hat{\mathcal{D}}(\vec{q})$ elements until saturation

#### Fringe-balanced trees (t > 1)

• probability of given value in root:





**Recall:** 
$$\alpha(\vec{q}) = \sum_{\nu=1}^{u} q_{\nu} \cdot depth_{\mathfrak{T}}(\nu)$$

#### *Warmup:* Ordinary BSTs (t = 0)

• Proof sketch:

Sum prob. that i is ancestor of j over all i, j ancestor  $\iff$  i **first** inserted key among i,..., j

base e Shannon entropy

• In the same paper:  $\mathfrak{H}_Q(\vec{q}) < \mathfrak{H}_{ln}^{\checkmark}(\vec{q})$ 

#### distribution with prob. weights $q_1, \ldots, q_u$ T from inserting i.i.d. $\hat{\mathcal{D}}(\vec{q})$ elements until saturation

#### Fringe-balanced trees (t > 1)

• probability of given value in root:



→ prob. that i inserted first among i,...j ??
Pold approach does not work

**Recall:** 
$$\alpha(\vec{q}) = \sum_{\nu=1}^{u} q_{\nu} \cdot depth_{\mathfrak{T}}(\nu)$$

#### *Warmup:* Ordinary BSTs (t = 0)

• Proof sketch:

Sum prob. that i is ancestor of j over all i, j ancestor  $\iff$  i **first** inserted key among i,..., j

base 2 entropy

base e Shannon entropy

• In the same paper:  $\mathfrak{H}_Q(\vec{q}) < \mathfrak{H}_{ln}^{\not}(\vec{q})$ 

# $\rightsquigarrow \alpha(\vec{q}) < 2 \ln 2 \cdot \mathcal{H}_{ld}^{\checkmark}(\vec{q}) + 1,$

distribution with prob. weights  $q_1, \ldots, q_u$ T from inserting i.i.d.  $\hat{\mathcal{D}}(\vec{q})$  elements **until saturation** 

### Fringe-balanced trees (t > 1)

• probability of given value in root:



→ prob. that i inserted first among i,...j ??
 Pold approach does not work

**Recall:** 
$$\alpha(\vec{q}) = \sum_{\nu=1}^{u} q_{\nu} \cdot depth_{\mathfrak{T}}(\nu)$$

#### *Warmup:* Ordinary BSTs (t = 0)

• Proof sketch:

Sum prob. that i is ancestor of j over all i, j ancestor  $\iff$  i **first** inserted key among i,..., j

base e Shannon entropy

• In the same paper:  $\mathfrak{H}_Q(\vec{q}) < \mathfrak{H}_{ln}^{\checkmark}(\vec{q})$ 

base 2 entropy

 $\rightsquigarrow \alpha(\vec{q}) < 2 \ln 2 \cdot \mathcal{H}_{ld}^{\checkmark}(\vec{q}) + 1,$ only factor  $2 \ln 2 \approx 1.386$  from optimal! distribution with prob. weights  $q_1, \ldots, q_u$ T from inserting i.i.d.  $\mathcal{D}(\vec{q})$  elements until saturation

### Fringe-balanced trees (t > 1)

• probability of given value in root:





2017-06-19

**Recall:** 
$$\alpha(\vec{q}) = \sum_{\nu=1}^{u} q_{\nu} \cdot depth_{\mathfrak{T}}(\nu)$$

#### *Warmup:* Ordinary BSTs (t = 0)

• Proof sketch:

Sum prob. that i is ancestor of j over all i, j ancestor  $\iff$  i **first** inserted key among i,..., j

base 2 entropy

base e Shannon entropy

• In the same paper:  $\mathfrak{H}_Q(\vec{q}) < \mathfrak{H}_{ln}^{\checkmark}(\vec{q})$ 

# $\rightsquigarrow \alpha(\vec{q}) < 2 \ln 2 \cdot \mathcal{H}_{ld}^{\downarrow}(\vec{q}) + 1,$

only factor  $2 \ln 2 \approx 1.386$  from **optimal**!

distribution with prob. weights  $q_1, \ldots, q_u$  $\mathfrak{T}$  from inserting i.i.d.  $\mathfrak{D}(\vec{q})$  elements until saturation

### Fringe-balanced trees (t > 1)

• probability of given value in root:





 $\rightsquigarrow\,$  Try to generalize this! —

< (1) →
• One of the defining properties of Shannon entropy: *aggregation* 



< @ >

• One of the defining properties of Shannon entropy: *aggregation* 



• One of the defining properties of Shannon entropy: *aggregation* 





• One of the defining properties of Shannon entropy: *aggregation* 





• One of the defining properties of Shannon entropy: *aggregation* 





< ₽ >

• One of the defining properties of Shannon entropy: *aggregation* 



• One of the defining properties of Shannon entropy: *aggregation* 





< @ >

• One of the defining properties of Shannon entropy: *aggregation* 







First partitioning step / Root of BST: Split into  $\langle P \rangle$ , =P,  $\geq P$ 

$$\Rightarrow \mathcal{H}(\vec{q}) = \mathcal{H}(V_1, H, V_2) + \sum_{i=1}^{N} V_j \cdot \mathcal{H}(Z_j) \quad \begin{array}{l} Z_1 = \left(\frac{q_1}{V_1}, \dots, \frac{q_{P-1}}{V_1}\right) \\ Z_2 = \left(\frac{q_{P+1}}{V_2}, \dots, \frac{q_u}{V_2}\right) \end{array}$$

< 17 >

• One of the defining properties of Shannon entropy: *aggregation* 







First partitioning step / Root of BST: Split into  $\langle P, =P, \rangle$ 

 $\rightarrow \mathcal{H}(\vec{q}) = \mathcal{H}(V_1, H, V_2) + \sum_{i=1}^{2} V_j \cdot \mathcal{H}(Z_j) \quad \begin{array}{l} Z_1 = \left(\frac{q_1}{V_1}, \dots, \frac{q_{P-1}}{V_1}\right) \\ Z_2 = \left(\frac{q_{P+1}}{V_2}, \dots, \frac{q_u}{V_2}\right) \end{array}$ 

2017-06-19

• One of the defining properties of Shannon entropy: *aggregation* 



operties  
ggregation  

$$\frac{1}{2} \frac{1}{3} \frac{1}{6}$$

$$\mathcal{H}(\frac{1}{2}, \frac{1}{3}, \frac{1}{6}) = \mathcal{H}(\frac{1}{2}, \frac{1}{2}) + \frac{1}{2} \cdot 0 + \frac{1}{2} \cdot \mathcal{H}(\frac{2}{3}, \frac{1}{3})$$
rst partitioning step / Root of BST: Split into  $\langle \mathsf{P}, \mathsf{P}, \mathsf{P}\rangle$   
 $\rightsquigarrow \mathcal{H}(\vec{\mathsf{q}}) = \mathcal{H}(\mathsf{V}_1, \mathsf{H}, \mathsf{V}_2) + \sum_{j=1}^{2} \mathsf{V}_j \cdot \mathcal{H}(\mathsf{Z}_j)$ 

$$\overset{Z_1}{=} (\frac{q_{1,1}}{V_1}, \dots, \frac{q_{p-1}}{V_1})$$

$$Z_2 = (\frac{q_{p+1}}{V_2}, \dots, \frac{q_{p}}{V_2})$$

< 🗗 ►

• One of the defining properties of Shannon entropy: *aggregation* 



$$\mathfrak{H}(\frac{1}{2}, \frac{1}{3}, \frac{1}{6}) = \mathfrak{H}(\frac{1}{2}, \frac{1}{2}) + \frac{1}{2} \cdot \mathfrak{H}(\frac{2}{3}, \frac{1}{3})$$

$$\mathfrak{H}(\frac{1}{2}, \frac{1}{3}, \frac{1}{6}) = \mathfrak{H}(\frac{1}{2}, \frac{1}{2}) + \frac{1}{2} \cdot \mathfrak{O} + \frac{1}{2} \cdot \mathfrak{H}(\frac{2}{3}, \frac{1}{3})$$

$$\mathfrak{H}(\frac{1}{2}, \frac{1}{3}, \frac{1}{6}) = \mathfrak{H}(V_1, \mathbb{H}, \mathbb{V}_2) + \sum_{j=1}^{2} \mathbb{V}_j \cdot \mathfrak{H}(\mathbb{Z}_j) \qquad Z_1 = (\frac{q_1}{\mathbb{V}_1}, \dots, \frac{q_{p-1}}{\mathbb{V}_1})$$

$$Z_2 = (\frac{q_{p+1}}{\mathbb{V}_2}, \dots, \frac{q_u}{\mathbb{V}_2})$$

2017-06-19

< 🗗 ►

One of the defining properties • of Shannon entropy: aggregation



 $\mathcal{H}(\frac{1}{2},\frac{1}{3},\frac{1}{6}) = \mathcal{H}(\frac{1}{2},\frac{1}{2}) + \frac{1}{2} \cdot \mathbf{0} + \frac{1}{2} \cdot \mathcal{H}(\frac{2}{3},\frac{1}{3})$ First partitioning step / Root of BST: Split into  $\langle P \rangle$ ,  $|=P \rangle$ ,  $|>P \rangle$  $\rightarrow \mathcal{H}(\vec{q}) = \mathcal{H}(V_1, H, V_2) + \sum_{i=1}^{2} V_j \cdot \mathcal{H}(Z_j) \qquad \begin{array}{l} Z_1 = \left(\frac{q_1}{V_1}, \dots, \frac{q_{P-1}}{V_1}\right) \\ Z_2 = \left(\frac{q_{P+1}}{V_2}, \dots, \frac{q_u}{V_2}\right) \end{array}$ • Recurrence for search costs:  $\alpha(\vec{q}) = 1 + \sum_{j=1}^{2} V_j \cdot \alpha(Z_j)$ 

 $\frac{1}{3}$ 

• One of the defining properties of Shannon entropy: *aggregation* 



• Recurrence for search costs:

 $\frac{1}{6}$ 

 $\frac{1}{2}$ 

÷

< @ >

 One of the defining properties of Shannon entropy: aggregation





< (1) →

 One of the defining properties of Shannon entropy: aggregation



- Recurrence for **search costs**:
- $\mathcal{H}(\frac{1}{2},\frac{1}{3},\frac{1}{6})$  $= \mathcal{H}(\frac{1}{2},\frac{1}{2}) + \frac{1}{2} \cdot \mathbf{0} + \frac{1}{2} \cdot \mathcal{H}(\frac{2}{3},\frac{1}{3})$ First partitioning step / Root of BST: Split into  $\langle P \rangle$ ,  $|=P \rangle$ ,  $|>P \rangle$  $\begin{array}{c|c} & \xrightarrow{P} & \xrightarrow{P} \\ \hline q_3 & \xrightarrow{q_4} q_5 & q_6 \end{array} & \longrightarrow & \mathcal{H}(\vec{q}) = & \mathcal{H}(V_1, H, V_2) + \sum_{j=1}^2 V_j \cdot \mathcal{H}(Z_j) & Z_1 = \left(\frac{q_1}{V_1}, \dots, \frac{q_{P-1}}{V_1}\right) \\ Z_2 = \left(\frac{q_{P+1}}{V_2}, \dots, \frac{q_u}{V_2}\right) \end{array}$ same shape!  $\alpha(\vec{q}) \stackrel{I}{=} 1 + \sum_{i=1}^{2} V_{i} \cdot \alpha(Z_{i}) \quad \rightsquigarrow \quad \underbrace{\mathcal{H}(V_{1}, H, V_{2}) \quad \text{vs. } 1?}$

#### Technical Issues

- 2  $\mathbb{E}[\mathcal{H}_{ln}(V_1, H, V_2)] \approx \mathbb{E}[\mathcal{H}_{ln}(D, 1-D)] = H_{k+1} H_{t+1}$  where  $D \stackrel{\mathbb{D}}{=} \text{Beta}(t+1, t+1)$ < (1) →

 $\frac{1}{3}$ 

 One of the defining properties of Shannon entropy: aggregation



 $\mathcal{H}(\frac{1}{2}, \frac{1}{3}, \frac{1}{6})$  $= \mathcal{H}(\frac{1}{2},\frac{1}{2}) + \frac{1}{2} \cdot \mathcal{O} + \frac{1}{2} \cdot \mathcal{H}(\frac{2}{3},\frac{1}{3})$ First partitioning step / Root of BST: Split into  $\langle P \rangle$ ,  $|=P \rangle$ ,  $|>P \rangle$  $\begin{array}{c|c} & \xrightarrow{P} & \xrightarrow{P} \\ \hline q_3 & \xrightarrow{q_4} q_5 & q_6 \end{array} & \longrightarrow & \mathcal{H}(\vec{q}) = & \mathcal{H}(V_1, H, V_2) + \sum_{j=1}^2 V_j \cdot \mathcal{H}(Z_j) & Z_1 = \left(\frac{q_1}{V_1}, \dots, \frac{q_{P-1}}{V_1}\right) \\ Z_2 = \left(\frac{q_{P+1}}{V_2}, \dots, \frac{q_u}{V_2}\right) \end{array}$ same shape! • Recurrence for search costs:  $\alpha(\vec{q}) \stackrel{\checkmark}{=} 1 + \sum_{j=1}^{2} V_j \cdot \alpha(Z_j) \rightsquigarrow (\mathcal{H}(V_1, H, V_2) \text{ vs. } 1?)$ 

#### **Technical Issues**

- **1** Pivot P is random  $\rightarrow$  take expectations over P (and thus V<sub>1,2</sub>, Z<sub>1,2</sub>).
- $\mathbb{E}[\mathcal{H}_{ln}(V_1,H,V_2)] \approx \mathbb{E}[\mathcal{H}_{ln}(D,1-D)] = H_{k+1} H_{t+1} \qquad \text{where } D \stackrel{\mathbb{P}}{=} \text{Beta}(t+1,t+1)$

 $\frac{1}{3}$ 

< (1) →

 One of the defining properties of Shannon entropy: aggregation



 $\mathcal{H}(\frac{1}{2}, \frac{1}{3}, \frac{1}{6})$  $= \mathcal{H}(\frac{1}{2},\frac{1}{2}) + \frac{1}{2} \cdot \mathbf{0} + \frac{1}{2} \cdot \mathcal{H}(\frac{2}{3},\frac{1}{3})$ First partitioning step / Root of BST: Split into  $|\langle P|, |=P|, |\rangle |P|$  $\begin{array}{c|c} & \xrightarrow{P} & \xrightarrow{P} \\ \hline q_3 & \xrightarrow{q_4} q_5 & q_6 \end{array} & & \\ & & \\ \end{array} & & \\ & & \\ \end{array} & \begin{array}{c} \mathcal{H}(\vec{q}) = \mathcal{H}(V_1, H, V_2) + \sum_{i=1}^2 V_j \cdot \mathcal{H}(Z_j) & \begin{array}{c} Z_1 = \left(\frac{q_1}{V_1}, \dots, \frac{q_{P-1}}{V_1}\right) \\ Z_2 = \left(\frac{q_{P+1}}{V_2}, \dots, \frac{q_u}{V_2}\right) \end{array}$ same shape! • Recurrence for search costs:  $\alpha(\vec{q}) \stackrel{\checkmark}{=} 1 + \sum_{i=1}^{2} V_{j} \cdot \alpha(Z_{j}) \rightsquigarrow (\mathcal{H}(V_{1}, H, V_{2}) \text{ vs. } 1?)$ 

#### **Technical Issues**

- **1** Pivot P is random  $\rightarrow$  take expectations over P (and thus V<sub>1,2</sub>, Z<sub>1,2</sub>).
- 2  $\mathbb{E}[\mathcal{H}_{ln}(V_1, H, V_2)] \approx \mathbb{E}[\mathcal{H}_{ln}(D, 1-D)] = H_{k+1} H_{t+1}$ where  $D \stackrel{\mathcal{D}}{=} Beta(t+1,t+1)$ < (1) →

 One of the defining properties of Shannon entropy: aggregation



- $\mathcal{H}(\frac{1}{2},\frac{1}{3},\frac{1}{4})$  $= \mathcal{H}(\frac{1}{2},\frac{1}{2}) + \frac{1}{2} \cdot \mathcal{O} + \frac{1}{2} \cdot \mathcal{H}(\frac{2}{3},\frac{1}{3})$ First partitioning step / Root of BST: Split into  $|\langle P|, |=P|, |\rangle |P|$  $\begin{array}{c|c} & \xrightarrow{P} & \xrightarrow{P} \\ \hline q_3 & \xrightarrow{q_4} q_5 & q_6 \end{array} & & \\ & & \\ \end{array} & & \\ & & \\ \end{array} & \begin{array}{c} \mathcal{H}(\vec{q}) = \mathcal{H}(V_1, H, V_2) + \sum_{i=1}^2 V_j \cdot \mathcal{H}(Z_j) & \begin{array}{c} Z_1 = \left(\frac{q_1}{V_1}, \dots, \frac{q_{P-1}}{V_1}\right) \\ Z_2 = \left(\frac{q_{P+1}}{V_2}, \dots, \frac{q_u}{V_2}\right) \end{array}$ same shape! • Recurrence for search costs:  $\alpha(\vec{q}) \stackrel{=}{=} 1 + \sum_{j=1}^{2} V_{j} \cdot \alpha(Z_{j}) \rightsquigarrow (\mathcal{H}(V_{1}, H, V_{2}) \text{ vs. } 1?)$

#### **Technical Issues**

- **1** Pivot P is random  $\rightarrow$  take expectations over P (and thus V<sub>1,2</sub>, Z<sub>1,2</sub>).
- 2  $\mathbb{E}[\mathcal{H}_{ln}(V_1, H, V_2)] \approx \mathbb{E}[\mathcal{H}_{ln}(D, 1-D)] = H_{k+1} H_{t+1}$ where  $D \stackrel{\mathcal{D}}{=} Beta(t+1,t+1)$ but not an inequality in either direction < (1) →

Quicksort Is Optimal For Many Equal Keys

#### **Entropy Bounds for Search Costs**

 $\rightsquigarrow \ \alpha(\vec{q}) = c \cdot \mathcal{H}(\vec{q})$  does not seem to hold for any constant c

 $\begin{array}{rcl} \alpha(\vec{q}) &\leqslant & c \ \cdot \ \mathcal{H}(\vec{q}) \ + \ d \\ \\ \alpha_{\vec{q}} &\geqslant & c' \cdot \ \mathcal{H}(\vec{q}) \ - \ d' \end{array}$ 

for family of constants (c, d) and (c', d').

• Always have  $c' < \alpha_k < c$  where  $\alpha_k = \frac{\ln 2}{H_{k+1} - H_{t+1}}$ 

• But we can show 5

 $\rightsquigarrow \ \alpha(\vec{q}) = c \cdot \mathcal{H}(\vec{q})$  does not seem to hold for any constant c

• But we can show 🔊

 $\begin{array}{lll} \alpha(\vec{q}) &\leqslant \ c \ \cdot \mbox{\boldmath $\mathcal{H}$}(\vec{q}) \ + \ d \\ \alpha_{\vec{q}} &\geqslant \ c' \cdot \mbox{\boldmath $\mathcal{H}$}(\vec{q}) \ - \ d' \end{array}$ 

for family of constants (c, d) and (c', d').

• Always have 
$$c' < \alpha_k < c$$
 where  $\alpha_k = \frac{\ln 2}{H_{k+1} - H_{t+1}}$ 

6

< 67 ▶

 $\alpha(\vec{q}) = c \cdot \mathcal{H}(\vec{q})$  does **not** seem to hold for any constant c  $\rightarrow$ 

• But we can show

 $\alpha(\vec{q}) \leq c \cdot \mathcal{H}(\vec{q}) + d$  $\alpha_{\vec{q}} \geq c' \cdot \mathcal{H}(\vec{q}) - d'$ 

for family of constants (c, d) and (c', d').

• Always have 
$$c' < \alpha_k < c$$
 where  

$$\alpha_k = \frac{\ln 2}{H_{k+1} - H_{t+1}}$$



< 17 >

 $\rightsquigarrow \ \alpha(\vec{q}) = c \cdot \mathcal{H}(\vec{q})$  does not seem to hold for any constant c

• But we can show

w zod

 $\begin{array}{lll} \alpha(\vec{q}) &\leqslant \ c \ \cdot \mbox{$\mathcal{H}$}(\vec{q}) \ + \ d \\ \alpha_{\vec{q}} &\geqslant \ c' \cdot \mbox{$\mathcal{H}$}(\vec{q}) \ - \ d' \end{array}$ 

for family of constants (c, d) and (c', d').

• Always have  $c' < \alpha_k < c$  where  $\alpha_k = \frac{ln 2}{H_{k+1} - H_{t+1}}$ 



< 17 >

 $\rightsquigarrow \ \alpha(\vec{q}) = c \cdot \mathcal{H}(\vec{q})$  does not seem to hold for any constant c

• But we can show

 $\begin{array}{lll} \alpha(\vec{q}) &\leqslant \ c \ \cdot \mbox{$\mathcal{H}$}(\vec{q}) \ + \ d \\ \\ \alpha_{\vec{q}} &\geqslant \ c' \cdot \mbox{$\mathcal{H}$}(\vec{q}) \ - \ d' \end{array}$ 

for family of constants (c, d) and (c', d').

• Always have  $c' < \alpha_k < c$  where  $\alpha_k = \frac{\ln 2}{H_{k+1} - H_{t+1}}$ 



 $\rightsquigarrow$  Asymptotically matching values for c and c'

$$\rightsquigarrow \left( \alpha(\vec{q}) = \alpha_k \cdot \mathcal{H}_{\mathsf{Id}}(\vec{q}) \pm O\left(\mathcal{H}(\vec{q})^{\frac{t+2}{t+3}} \log(\mathcal{H}(\vec{q}))\right) \right)$$

2017-06-19



#### Time to put the pieces together!

- Separation Theorem: Quicksort costs
  - in the i.i.d. model
  - with **"many duplicates"** ( $\Omega(n^{\varepsilon})$  duplicates of each value in expectation)

are given by (as  $n \to \infty$ , for any  $\delta \in (0, \varepsilon)$ )

 $\left[\mathbb{E}[C_{n,\vec{q}}] = \alpha(\vec{q}) \cdot n \pm O(n^{1-\delta})\right]$ 

Average search costs in saturated k-fringe-balanced trees: (as H → ∞)

 $\alpha(\vec{q}) = \alpha_k \cdot \mathcal{H} \pm O\left(\mathcal{H}^{\frac{t+2}{t+3}}\log \mathcal{H}\right)$ 

• 
$$\mathcal{H} = \mathcal{H}_{ld}(\vec{q})$$
  
•  $\alpha_k = \frac{\ln 2}{H_{k+1} - H_{t+1}}$ 

#### Quicksort Costs (i.i.d. model)

 $\sim \rightarrow$ 

Under the assumptions above, we have for any  $\delta \in (0, \frac{1}{t+3})$  $\mathbb{E}[C_{n,\vec{q}}] = \alpha_k \mathcal{H}_{\mathsf{Id}}(\vec{q}) \cdot n \pm O((\mathcal{H}(\vec{q})^{1-\delta}+1)n).$ 

< (1) →



#### Time to put the pieces together!

#### Separation Theorem: Quicksort costs

- in the i.i.d. model
- with "many duplicates"
   (Ω(n<sup>ε</sup>) duplicates of each value in expectation)

are given by  $(as n \to \infty, \text{ for any } \delta \in (0, \epsilon))$ 

 $\left(\mathbb{E}[C_{n,\vec{q}}] = \alpha(\vec{q}) \cdot n \pm O(n^{1-\delta})\right)$ 

Average search costs in saturated k-fringe-balanced trees:
(as ff => ∞)

 $\alpha(\vec{q}) \; = \; \alpha_k \cdot \mathcal{H} \; \pm \; O \big( \mathcal{H}^{\frac{t+2}{t+3}} \log \mathcal{H} \big)^{-1}$ 

• 
$$\mathcal{H} = \mathcal{H}_{\mathsf{ld}}(\vec{q})$$
  
•  $\alpha_{\mathsf{k}} = \frac{\ln 2}{\mathsf{H}_{\mathsf{k}+1} - \mathsf{H}_{\mathsf{t}+1}}$ 

#### Quicksort Costs (i.i.d. model)

 $\sim \rightarrow$ 

Under the assumptions above, we have for any  $\delta \in (0, \frac{1}{t+3})$  $\mathbb{E}[C_{n,\vec{q}}] = \alpha_k \mathcal{H}_{\mathsf{ld}}(\vec{q}) \cdot n \pm O((\mathcal{H}(\vec{q})^{1-\delta}+1)n).$ 



#### Time to put the pieces together!

- Separation Theorem: Quicksort costs
  - in the i.i.d. model
  - with "many duplicates"
     (Ω(n<sup>ε</sup>) duplicates of each value in expectation)

are given by  $(as n \to \infty, \text{ for any } \delta \in (0, \epsilon))$ 

 $\left(\mathbb{E}[C_{n,\vec{q}}] = \alpha(\vec{q}) \cdot n \pm O(n^{1-\delta})\right)$ 

Average search costs

in saturated k-fringe-balanced trees:

(as  $\mathcal{H} \to \infty$ )

 $\alpha(\vec{q}) \; = \; \alpha_k \cdot \mathcal{H} \; \pm \; O \big( \mathcal{H}^{\frac{t+2}{t+3}} \log \mathcal{H} \big)$ 

• 
$$\mathcal{H} = \mathcal{H}_{\mathsf{ld}}(\vec{q})$$
  
•  $\alpha_k = \frac{\ln 2}{H_{k+1} - H_{t+1}}$ 

 $\sim \rightarrow$ 

Under the assumptions above, we have for any  $\delta \in (0, \frac{1}{t+3})$  $\mathbb{E}[C_{n,\vec{q}}] = \alpha_k \mathcal{H}_{ld}(\vec{q}) \cdot n \pm O((\mathcal{H}(\vec{q})^{1-\delta} + 1)n).$ 



#### Time to put the pieces together!

- Separation Theorem: Quicksort costs
  - in the i.i.d. model
  - with "many duplicates"
     (Ω(n<sup>ε</sup>) duplicates of each value in expectation)

are given by  $(as n \to \infty, \text{ for any } \delta \in (0, \epsilon))$ 

 $\left(\mathbb{E}[C_{n,\vec{q}}] = \alpha(\vec{q}) \cdot n \pm O(n^{1-\delta})\right)$ 

# Average search costs

in saturated k-fringe-balanced trees:

(as  $\mathcal{H} \to \infty$ )

 $\alpha(\vec{q}) \; = \; \alpha_k \cdot \mathcal{H} \; \pm \; O \big( \mathcal{H}^{\frac{t+2}{t+3}} \log \mathcal{H} \big)$ 

• 
$$\mathcal{H} = \mathcal{H}_{\mathsf{ld}}(\vec{q})$$
  
•  $\alpha_k = \frac{\ln 2}{H_{k+1} - H_{t+1}}$ 

Quicksort Costs (i.i.d. model)Under the assumptions above, we have for any  $\delta \in (0, \frac{1}{t+3})$  $\mathbb{E}[C_{n,\vec{q}}] = \alpha_k \mathcal{H}_{\mathsf{ld}}(\vec{q}) \cdot n \pm O((\mathcal{H}(\vec{q})^{1-\delta}+1)n).$ 

 $\sim \rightarrow$ 





# **Quicksort and Search Trees**



**Saturated Fringe-Balanced Trees** 

**Back to Multiset Permutations** 

How about the multiset model? 🔀



- **1** Replace multiset model with profile  $\vec{x}$  by i.i.d. model with  $\vec{q} = \vec{x}/n$
- **2** Use Chernoff bounds to **bound difference** between costs.

→ Need Chernoff bound for **multinomial variables**.

How about the multiset model?

 $\sum$  - Many duplicates  $\rightarrow$  profile  $\vec{X}$  concentrated around  $\mathbb{E}[\vec{X}] = \vec{q}n$ 

- **1** Replace multiset model with profile  $\vec{x}$  by i.i.d. model with  $\vec{q} = \vec{x}/n$
- 2 Use Chernoff bounds to **bound difference** between costs.

→ Need Chernoff bound for **multinomial variables**.

How about the multiset model?

- $\bigvee$  Many duplicates  $\rightsquigarrow$  profile  $\vec{X}$  concentrated around  $\mathbb{E}[\vec{X}] = \vec{q}n$
- **①** Replace multiset model with profile  $\vec{x}$  by i.i.d. model with  $\vec{q} = \vec{x}/n$
- **2** Use Chernoff bounds to **bound difference** between costs.

→ Need Chernoff bound for **multinomial variables**.

How about the multiset model?

- $\bigvee$  Many duplicates  $\rightsquigarrow$  profile  $\vec{X}$  concentrated around  $\mathbb{E}[\vec{X}] = \vec{q}n$ 
  - **①** Replace multiset model with profile  $\vec{x}$  by i.i.d. model with  $\vec{q} = \vec{x}/n$
  - **2** Use Chernoff bounds to **bound difference** between costs.

→ Need Chernoff bound for **multinomial variables**.

How about the multiset model?

- $\bigvee$  Many duplicates  $\rightsquigarrow$  profile  $\vec{X}$  concentrated around  $\mathbb{E}[\vec{X}] = \vec{q}n$ 
  - **0** Replace multiset model with profile  $\vec{x}$  by i.i.d. model with  $\vec{q} = \vec{x}/n$
  - **2** Use Chernoff bounds to **bound difference** between costs.
  - ~ Need Chernoff bound for multinomial variables.

How about the multiset model? 🗙

- $\bigvee$  Many duplicates  $\rightsquigarrow$  profile  $\vec{X}$  concentrated around  $\mathbb{E}[\vec{X}] = \vec{q}n$
- **①** Replace multiset model with profile  $\vec{x}$  by i.i.d. model with  $\vec{q} = \vec{x}/n$
- **2** Use Chernoff bounds to **bound difference** between costs.

→ Need Chernoff bound for **multinomial variables**.

*The Annals of Statistics* 1983, Vol. 11, No. 3, 896–904

# THE EQUIVALENCE OF WEAK, STRONG AND COMPLETE CONVERGENCE IN $L_1$ FOR KERNEL DENSITY ESTIMATES<sup>1</sup>

By LUC DEVROYE

McGill University

Let f be a density on  $\mathbb{R}^d$ , and let  $f_n$  be the kernel estimate of f,

 $f_n(x) = (nh^d)^{-1} \sum_{i=1}^n K((x - X_i)/h)$ 

where  $h = h_n$  is a sequence of positive numbers, and K is an absolutely integrable function with  $\int K(x) dx = 1$ . Let  $J_n = \int |f_n(x) - f(x)| dx$ . We show

Quicksort Is Optimal For Many Equal Keys

How about the multiset model? 🗙

 $\dot{\mathbf{X}}$ - Many duplicates  $\rightsquigarrow$  profile  $\vec{X}$  concentrated around  $\mathbb{E}[\vec{X}] = \vec{q}n$ 

- **①** Replace multiset model with profile  $\vec{x}$  by i.i.d. model with  $\vec{q} = \vec{x}/n$
- **2** Use Chernoff bounds to **bound difference** between costs.

→ Need Chernoff bound for **multinomial variables**.

#### The Annals of Statistics

LEMMA 3. (A multinomial distribution inequality). Let  $(X_1, \dots, X_k)$  be a multinomial  $(n, p_1, \dots, p_k)$  random vector. For all  $\varepsilon \in (0, 1)$  and all k satisfying  $k/n \le \varepsilon^2/20$ , we have

```
P(\sum_{i=1}^{k} |X_i - E(X_i)| > n\varepsilon) \le 3 \exp(-n\varepsilon^2/25).
```

#### McGill University

Let f be a density on  $\mathbb{R}^d$ , and let  $f_n$  be the kernel estimate of f,

$$f_n(x) = (nh^d)^{-1} \sum_{i=1}^n K((x - X_i)/h)$$

where  $h = h_n$  is a sequence of positive numbers, and K is an absolutely integrable function with  $\int K(x) dx = 1$ . Let  $J_n = \int |f_n(x) - f(x)| dx$ . We show

Quicksort Is Optimal For Many Equal Keys

How about the multiset model?

 $\dot{\mathbf{V}}$ - Many duplicates  $\rightsquigarrow$  profile  $\vec{X}$  concentrated around  $\mathbb{E}[\vec{X}] = \vec{q}n$ 

- **1** Replace multiset model with profile  $\vec{x}$  by i.i.d. model with  $\vec{q} = \vec{x}/n$
- 2 Use Chernoff bounds to **bound difference** between costs.

→ Need Chernoff bound for **multinomial variables**.

#### The Annals of Statistics

LEMMA 3. (A multinomial distribution inequality). Let  $(X_1, \dots, X_k)$  be a multinomial  $(n, p_1, \dots, p_k)$  random vector. For all  $\varepsilon \in (0, 1)$  and all k satisfying  $k/n \le \varepsilon^2/20$ , we have

```
P(\sum_{i=1}^{k} |X_i - E(X_i)| > n\varepsilon) \le 3 \exp(-n\varepsilon^2/25).
```

McGill University


## Findings

- First analysis of median-of-k Quicksort on equal keys ... for "many duplicates".
  - Same relative speedup as for random permutations.  $\rightarrow$
- Partial Answer to **conjecture of Sedgewick & Bentley**:
- Not in this talk: For uniform  $\vec{q} = (\frac{1}{u}, \dots, \frac{1}{u})$  with  $u = O(n^{1-\epsilon})$



## Findings

- First analysis of median-of-k Quicksort on equal keys ... for "many duplicates".
  - → Same relative speedup as for random permutations.
- Partial Answer to conjecture of Sedgewick & Bentley:

Median-of-k Quicksort approaches lower bound for  $k \to \infty.$ 

- Not in this talk: For uniform  $\vec{q} = (\frac{1}{u}, \dots, \frac{1}{u})$  with  $u = O(n^{1-\varepsilon})$ 
  - better error bounds
  - extension for multiway partitioning

## **Open Problems**

- Get rid of "many duplicates" restriction;  $n^{1-\epsilon}$  seems (to me) best possible so that
  - inputs are non-degenerate w.h.p.
  - tree-building costs are still negligible
  - difference between i.i.d. model and multiset model is negligible

the entropy is a lower bound



## Findings

- First analysis of median-of-k Quicksort on equal keys ... for "many duplicates".
  - → Same relative speedup as for random permutations.
- Partial Answer to conjecture of Sedgewick & Bentley:

Median-of-k Quicksort approaches lower bound for  $k \to \infty.$ 

- Not in this talk: For uniform  $\vec{q} = (\frac{1}{u}, \dots, \frac{1}{u})$  with  $u = O(n^{1-\epsilon})$ 
  - better error bounds
  - extension for multiway partitioning

## **Open Problems**

- Get rid of "many duplicates" restriction;  $n^{1-\epsilon}$  seems (to me) best possible so that
  - inputs are non-degenerate w.h.p.
  - tree-building costs are still negligible
  - difference between i.i.d. model and multiset model is negligible.

the entropy is a lower boundary



## Findings

- First analysis of median-of-k Quicksort on equal keys ... for "many duplicates".
  - Same relative speedup as for random permutations.  $\sim \rightarrow$
- Partial Answer to **conjecture of Sedgewick & Bentley**:

Median-of-k Quicksort approaches lower bound for  $k \to \infty$ .

- Not in this talk: For uniform  $\vec{q} = (\frac{1}{n}, \dots, \frac{1}{n})$  with  $u = O(n^{1-\varepsilon})$  better error bounds
  - extension for multiway partitioning



## Findings

- First analysis of median-of-k Quicksort on equal keys ... for "many duplicates".
  - → Same relative speedup as for random permutations.
- Partial Answer to conjecture of Sedgewick & Bentley:

Median-of-k Quicksort approaches lower bound for  $k \to \infty.$ 

- Not in this talk: For uniform  $\vec{q} = (\frac{1}{u}, \dots, \frac{1}{u})$  with  $u = O(n^{1-\epsilon})$ 
  - better error bounds
  - extension for multiway partitioning

## **Open Problems**

- Get rid of "many duplicates" restriction;  $n^{1-\epsilon}$  seems (to me) best possible so that
  - inputs are non-degenerate w.h.p.:
  - tree-building costs are still negligible
  - difference between i.i.d. model and multiset model is negligible.

the entropy is a lower bound



## Findings

- First analysis of median-of-k Quicksort on equal keys ... for "many duplicates".
  - → Same relative speedup as for random permutations.
- Partial Answer to conjecture of Sedgewick & Bentley:

Median-of-k Quicksort approaches lower bound for  $k \to \infty.$ 

- Not in this talk: For uniform  $\vec{q} = (\frac{1}{u}, \dots, \frac{1}{u})$  with  $u = O(n^{1-\epsilon})$ 
  - better error bounds
  - extension for multiway partitioning

#### **Open Problems**

- Get rid of "many duplicates" restriction;  $n^{1-\varepsilon}$  seems (to me) best possible so that
  - inputs are non-degenerate w.h.p.
  - tree-building costs are still negligible
  - difference between i.i.d. model and multiset model is negligible
  - the entropy is a lower bound



## Findings

- First analysis of median-of-k Quicksort on equal keys ... for "many duplicates".
  - → Same relative speedup as for random permutations.
- Partial Answer to conjecture of Sedgewick & Bentley:

Median-of-k Quicksort approaches lower bound for  $k \to \infty.$ 

- Not in this talk: For uniform  $\vec{q} = (\frac{1}{u}, \dots, \frac{1}{u})$  with  $u = O(n^{1-\epsilon})$ 
  - better error bounds
  - extension for multiway partitioning

## **Open Problems**

- Get rid of "many duplicates" restriction;  $n^{1-\epsilon}$  seems (to me) best possible so that
  - inputs are non-degenerate w.h.p.
  - tree-building costs are still negligible
  - difference between i.i.d. model and multiset model is negligible
  - the entropy is a lower bound



## Findings

- First analysis of median-of-k Quicksort on equal keys ... for "many duplicates".
  - → Same relative speedup as for random permutations.
- Partial Answer to conjecture of Sedgewick & Bentley:

Median-of-k Quicksort approaches lower bound for  $k \to \infty.$ 

- Not in this talk: For uniform  $\vec{q} = (\frac{1}{u}, \dots, \frac{1}{u})$  with  $u = O(n^{1-\epsilon})$ 
  - better error bounds
  - extension for multiway partitioning

## **Open Problems**

- Get rid of "many duplicates" restriction;  $n^{1-\epsilon}$  seems (to me) best possible so that
  - inputs are non-degenerate w.h.p.
  - tree-building costs are still negligible
  - difference between i.i.d. model and multiset model is negligible
  - the entropy is a lower bound



## Findings

- First analysis of median-of-k Quicksort on equal keys ... for "many duplicates".
  - → Same relative speedup as for random permutations.
- Partial Answer to conjecture of Sedgewick & Bentley:

Median-of-k Quicksort approaches lower bound for  $k \to \infty.$ 

- Not in this talk: For uniform  $\vec{q} = (\frac{1}{u}, \dots, \frac{1}{u})$  with  $u = O(n^{1-\epsilon})$ 
  - better error bounds
  - extension for multiway partitioning

## **Open Problems**

- Get rid of "many duplicates" restriction;  $n^{1-\varepsilon}$  seems (to me) best possible so that
  - inputs are non-degenerate w.h.p.
  - tree-building costs are still negligible
  - difference between i.i.d. model and multiset model is negligible
  - the entropy is a lower bound



## Findings

- First analysis of median-of-k Quicksort on equal keys ... for "many duplicates".
  - → Same relative speedup as for random permutations.
- Partial Answer to conjecture of Sedgewick & Bentley:

Median-of-k Quicksort approaches lower bound for  $k \to \infty.$ 

- Not in this talk: For uniform  $\vec{q} = (\frac{1}{u}, \dots, \frac{1}{u})$  with  $u = O(n^{1-\epsilon})$ 
  - better error bounds
  - extension for multiway partitioning

## **Open Problems**

- Get rid of "many duplicates" restriction;  $n^{1-\epsilon}$  seems (to me) best possible so that
  - inputs are non-degenerate w.h.p.
  - tree-building costs are still negligible
  - difference between i.i.d. model and multiset model is negligible
  - the entropy is a lower bound





## Findings

- First analysis of median-of-k Quicksort on equal keys ... for "many duplicates".
  - → Same relative speedup as for random permutations.
- Partial Answer to conjecture of Sedgewick & Bentley:

Median-of-k Quicksort approaches lower bound for  $k \to \infty.$ 

- Not in this talk: For uniform  $\vec{q} = (\frac{1}{u}, \dots, \frac{1}{u})$  with  $u = O(n^{1-\epsilon})$ 
  - better error bounds
  - extension for multiway partitioning

## **Open Problems**

- Get rid of "many duplicates" restriction;  $n^{1-\varepsilon}$  seems (to me) best possible so that
  - inputs are non-degenerate w.h.p.
  - tree-building costs are still negligible
  - difference between i.i.d. model and multiset model is negligible
  - the entropy is a lower bound





Icons made by Freepik and Gregor Cresnar from www.flaticon.com.

	_
0.0000000	

<∂>