

Distribution and Variance of the Optimal Multi-Pivot-Quicksort

Daniel Krenn

(joint work in progress with *Cecilia Holmgren*)



UPPSALA
UNIVERSITET

June 25, 2019



This presentation is licensed under a Creative Commons
Attribution-NonCommercial-ShareAlike 3.0 Unported License.

FWF
Der Wissenschaftsfonds.

Supported by the
Austrian Science Fund (FWF),
project P28466.

Quicksort & Quickselect



Quicksort & Quickselect



- choose a pivot element p

Quicksort & Quickselect



- choose a pivot element p
- partition into
 - small elements
 - large elements



Quicksort & Quickselect



- choose a pivot element p
- partition into
 - small elements
 - large elements



- proceed recursively

Dual Pivot Quicksort & Quickselect



Dual Pivot Quicksort & Quickselect

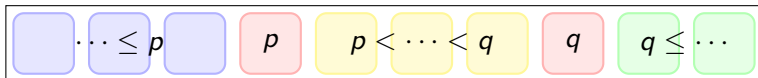


- choose pivot elements p and q

Dual Pivot Quicksort & Quickselect



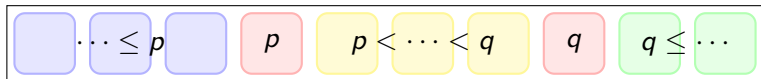
- choose pivot elements p and q
- partition into
 - small elements
 - medium elements
 - large elements



Dual Pivot Quicksort & Quickselect



- choose pivot elements p and q
- partition into
 - small elements
 - medium elements
 - large elements



- proceed recursively

Partitioning: How many Key Comparisons?

- “classical” quicksort



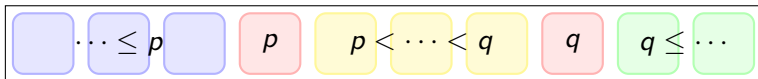
- always $\rightsquigarrow n - 1$ comparisons

Partitioning: How many Key Comparisons?

- “classical” quicksort

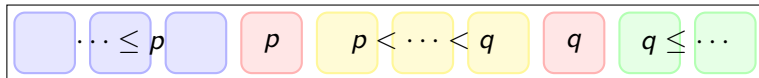


- always $\rightsquigarrow n - 1$ comparisons
- dual pivot quicksort



- always at most $\rightsquigarrow 2(n - 2)$ comparisons
- “Yaroslavskiy” $\rightsquigarrow (1.583 \dots)n + O(\log n)$ comparisons
- optimal $\rightsquigarrow 1.5n + 0.25 \log n + O(1)$ comparisons

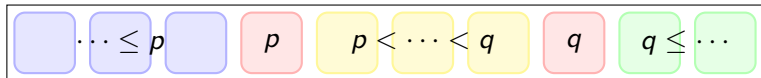
Average Number of Key Comparisons



- partitioning
 - “classical” $\rightsquigarrow n - 1$

- quicksort
 - “classical” $\rightsquigarrow 2n \log n - (2.84 \dots)n + O(\log n)$

Average Number of Key Comparisons



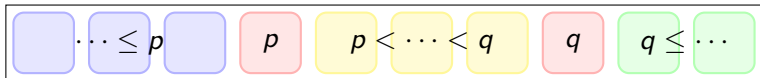
- partitioning

- “classical” $\rightsquigarrow n - 1$

- quicksort

- “classical” $\rightsquigarrow 2n \log n - (2.84 \dots)n + O(\log n)$
- “Yaroslavskiy–Bentley–Bloch” $\rightsquigarrow 1.9n \log n - (2.46 \dots)n + O(\log n)$
[Wild–Nebel 2012]

Average Number of Key Comparisons



- partitioning

- “classical” $\rightsquigarrow n - 1$

- “optimal dual pivot” $\rightsquigarrow 1.5n + 0.25 \log n + O(1)$

[Aumüller–Dietzfelbinger 2014,
Aumüller–Dietzfelbinger–Heuberger–K–Prodingen 2016]

- quicksort

- “classical” $\rightsquigarrow 2n \log n - (2.84 \dots)n + O(\log n)$

- “Yaroslavskiy–Bentley–Bloch” $\rightsquigarrow 1.9n \log n - (2.46 \dots)n + O(\log n)$

[Wild–Nebel 2012]

- “optimal dual pivot” $\rightsquigarrow 1.8n \log n - (2.38 \dots)n + O(\log n)$

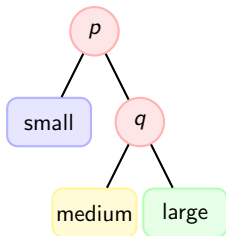
[Aumüller–Dietzfelbinger 2014,
Aumüller–Dietzfelbinger–Heuberger–K–Prodingen 2016]

Optimal Partitioning Strategy “Count”

- comparison of element with pivots:
 - seen more **small elements** \rightsquigarrow smaller pivot p first
 - seen more **large elements** \rightsquigarrow larger pivot q first
 - equality \rightsquigarrow choice: smaller pivot p first

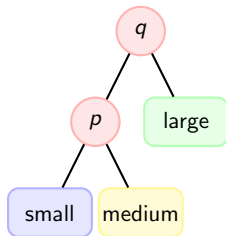
Optimal Partitioning Strategy “Count”

- comparison of element with pivots:
 - seen more **small elements** \rightsquigarrow smaller pivot p first
 - seen more **large elements** \rightsquigarrow larger pivot q first
 - equality \rightsquigarrow choice: smaller pivot p first
- comparison trees:



choose this tree if

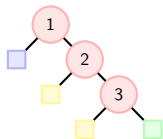
$$\#small \geq \#large$$



choose this tree if

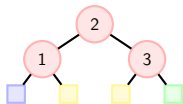
$$\#small < \#large$$

Optimal Partitioning Strategy with Three Pivots



$$s_1 \geq s_3$$

$$s_0 \geq s_2 + s_3 + 1$$

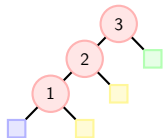


$$s_2 + s_3 + 1 \geq s_0$$

$$s_1 + s_2 + 1 \geq s_0$$

$$s_1 + s_2 + 1 \geq s_3$$

$$s_0 + s_1 + 1 \geq s_3$$



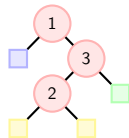
$$s_3 \geq s_0 + s_1 + 1$$

$$s_2 \geq s_0$$

$$s_0 \geq s_1 + s_2 + 1$$

$$s_3 \geq s_1$$

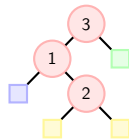
$$s_0 \geq s_3$$



$$s_0 \geq s_2$$

$$s_3 \geq s_0$$

$$s_3 \geq s_1 + s_2 + 1$$



comparison trees & polyhedra minimizing $\ell_t(s) = \sum_{i=0}^d h_i(t)(s_i + 1)$

Optimal Dual-Pivot Quicksort

Theorem (ADHKP 2016+2018)

average number of key comparisons
in dual pivot quicksort

with the **optimal** partitioning strategy "Count" is

$$\frac{9}{5}nH_n - \frac{1}{5}nH_n^{\text{alt}} - \frac{89}{25}n + \frac{67}{40}H_n - \frac{3}{40}H_n^{\text{alt}} - \frac{83}{800} + \frac{(-1)^n}{10} + \dots$$

$$= \frac{9}{5}n \log n + An + B \log n + C + \frac{D}{n} + \frac{E}{n^2} + \frac{(-1)^n F + G}{n^3} + O\left(\frac{1}{n^4}\right)$$

asymptotically as $n \rightarrow \infty$

- harmonic numbers

- $H_n = \sum_{i=1}^n 1/i$
- $H_n^{\text{alt}} = \sum_{i=1}^n (-1)^i / i$

- constant of linear term

$$A = \frac{9}{5}\gamma + \frac{1}{5} \log 2 - \frac{89}{25} = -2.3823823670652 \dots$$

- explicit constants B, C, \dots

Optimal Multi-Pivot Quicksort

Theorem (Heuberger–K 2019)



- *average number of key comparisons
in trial pivot quicksort
with the **optimal** partitioning strategy is*

$$\frac{133}{78}n \log n + An + B \log n + O(1)$$

- $A = \frac{133}{78}\gamma - \frac{2}{117}\sqrt{3}\pi + \frac{4}{39}\log 3 + \frac{3}{26}\log 2 - \frac{6761}{2028} = -2.24995\dots$
- $B = \frac{707}{468}$
- *average number of key comparisons
in quadral pivot quicksort*

*with the **optimal** partitioning strategy is*

$$\frac{9536}{5775}n \log n + An + B \log n + O(1)$$

- $A = \text{skipped} = -2.20515\dots$
- $B = \frac{48823}{34650}$

Limiting Distribution: Brief History

Quantity

number of key comparisons

- classical quicksort
 - convergence in law to some limiting distribution [Régnier 1989]
 - cumulants [Hennequin 1989]
 - distribution implicitly characterized by stochastic fixed point equation [Rösler 1991]
 - many more properties . . .
- Yaroslavskiy–Bentley–Bloch dual-pivot quicksort
 - distribution & variance [Wild–Nebel–Neininger 2015]
- optimal (“Count”) dual-pivot quicksort
 - distribution & variance [Neininger–Straub 2018]



Optimal Multi-pivot Quicksort Distribution

- random variable $C_n =$ number of key comparisons for sorting a list of n elements

Theorem (Holmgren–K 2019+)

- $\frac{1}{n}(C_n - \mathbb{E}(C_n)) \rightarrow C^*$ as $n \rightarrow \infty$ in distribution
- C^* whose distribution is *unique fixed point* of

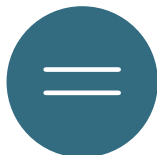
$$Z \stackrel{\mathcal{D}}{=} \sum_{i=0}^d (D_i Z^{(i)} + c_0 D_i \log D_i) + \sum_{t \in \mathcal{T}} [D \in \mathcal{C}_t^\infty] \ell_t^\infty(D)$$

with

- $\ell_t^\infty(x) = \sum_{i=0}^d x_i h_i(t)$
- non-negative random variables $D = (D_0, \dots, D_d)$ uniformly distributed among $D_0 + \dots + D_d = 1$
- asymptotic polyhedron \mathcal{C}_t^∞ for tree t

Towards Fixed Point Equation

- distributional recurrence $C_n \stackrel{\mathcal{D}}{=} \sum_{i=0}^d C_{l_i}^{(i)} + P_n$ with
 - $\sum_{i=0}^d l_i = n - d$
 - independent random variables $(l_0, \dots, l_d, P_n), (C_n^{(0)})_{n \geq 0}, \dots, (C_n^{(d)})_{n \geq 0}$



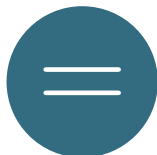
Towards Fixed Point Equation

- distributional recurrence $C_n \stackrel{\mathcal{D}}{=} \sum_{i=0}^d C_{l_i}^{(i)} + P_n$ with
 - $\sum_{i=0}^d l_i = n - d$
 - independent random variables $(l_0, \dots, l_d, P_n), (C_n^{(0)})_{n \geq 0}, \dots, (C_n^{(d)})_{n \geq 0}$

$$\begin{aligned}
 \bullet \quad C_n^* &= \frac{1}{n}(C_n - \mathbb{E}(C_n)) \\
 &= \underbrace{\sum_{i=0}^d \frac{l_i}{n} C_{l_i}^{*,(i)}}_{\rightarrow D_i} + \underbrace{\frac{1}{n} \left(-\mathbb{E}(C_n) + \sum_{i=0}^d \mathbb{E}(C_{l_i} | l_i) \right)}_{\substack{= c_0 \sum_{i=0}^d \frac{l_i}{n} \log \frac{l_i}{n} + o(1) \\ \rightarrow c_0 \sum_{i=0}^d D_i \log D_i}} + \underbrace{\frac{P_n}{n}}_{\rightarrow P^*}
 \end{aligned}$$

because

- $\mathbb{E}(C_n) = c_0 n \log n + c_1 n + o(n)$
- $\mathbb{E}(C_{l_i} | l_i) = c_0 l_i \log l_i + c_1 l_i + o(n)$



- contraction method

Distribution of Partitioning Cost

- number of key comparisons P_n for partitioning a list of n elements (by d pivot elements) into $d + 1$ sublists

Theorem (Holmgren–K 2019+)

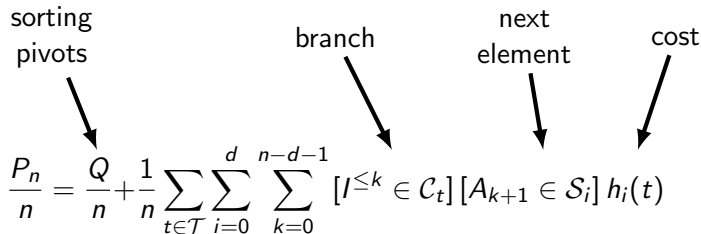
$$\frac{P_n}{n} \longrightarrow P^* = \sum_{t \in \mathcal{T}} [D \in \mathcal{C}_t^\infty] \ell_t^\infty(D)$$

as $n \rightarrow \infty$ in distribution with

- $\ell_t^\infty(x) = \sum_{i=0}^d x_i h_i(t)$
- non-negative random variables $D = (D_0, \dots, D_d)$ uniformly distributed among $D_0 + \dots + D_d = 1$
- asymptotic polyhedron \mathcal{C}_t^∞ for tree t

Distribution of Partitioning Cost: Idea

sorting pivots branch next element cost


$$\frac{P_n}{n} = \frac{Q}{n} + \frac{1}{n} \sum_{t \in \mathcal{T}} \sum_{i=0}^d \sum_{k=0}^{n-d-1} [I^{\leq k} \in \mathcal{C}_t] [A_{k+1} \in \mathcal{S}_i] h_i(t)$$

Distribution of Partitioning Cost: Idea

sorting pivots branch next element cost

$$\frac{P_n}{n} = \frac{Q}{n} + \frac{1}{n} \sum_{t \in \mathcal{T}} \sum_{i=0}^d \sum_{k=0}^{n-d-1} [I^{\leq k} \in \mathcal{C}_t] [A_{k+1} \in \mathcal{S}_i] h_i(t)$$

$$\rightsquigarrow 0 + \sum_{t \in \mathcal{T}} \sum_{i=0}^d \dots [I \in \mathcal{C}_t] \dots \frac{l_i}{n} \dots h_i(t)$$

$$\rightarrow \sum_{t \in \mathcal{T}} \dots [D \in \mathcal{C}_t^\infty] \dots \underbrace{\sum_{i=0}^d D_i}_{= \ell_t^\infty(D)} h_i(t)$$

Towards Variance

- $\frac{1}{n}(C_n - \mathbb{E}(C_n)) \longrightarrow C^*$ as $n \rightarrow \infty$

for variance
need to compute $\mathbb{E}((C^*)^2)$



Towards Variance

- $\frac{1}{n}(C_n - \mathbb{E}(C_n)) \longrightarrow C^*$ as $n \rightarrow \infty$

for variance
need to compute $\mathbb{E}((C^*)^2)$



- use fixed point equation:

$$\mathbb{E}((C^*)^2) = \mathbb{E}\left(\left(\sum_{i=0}^d D_i C^{*,(i)} + c_0 \sum_{i=0}^d D_i \log D_i + P^*\right)^2\right)$$

Towards Variance

- $\frac{1}{n}(C_n - \mathbb{E}(C_n)) \longrightarrow C^*$ as $n \rightarrow \infty$

for variance
need to compute $\mathbb{E}((C^*)^2)$



- use fixed point equation:

$$\mathbb{E}((C^*)^2) = \mathbb{E}\left(\left(\underbrace{\sum_{i=0}^d D_i}_{\downarrow} \underbrace{C^{*,(i)}}_{\downarrow} + c_0 \sum_{i=0}^d \underbrace{D_i}_{\downarrow} \log \underbrace{D_i}_{\downarrow} + \underbrace{P^*}_{\downarrow}\right)^2\right)$$

$$\mathbb{E}((C^*)^2) = \left(1 - \sum_{i=0}^d \mathbb{E}(D_i^2)\right)^{-1} \mathbb{E}\left(\left(c_0 \sum_{i=0}^d D_i \log D_i + P^*\right)^2\right)$$

Computing Expected Values

... for computing second moments

Integral

$$\mathbb{E}(f(D)) = \frac{1}{v} \int_{\mathcal{C}} f(x) d\lambda_d(x)$$

with volume $v = \lambda_d(\mathcal{C}) = \int_{\mathcal{C}} d\lambda_d(x)$

- function f :
 - polynomial
 - one logarithmic factor
 - more logarithmic factors
- polytope \mathcal{C}
(d -dimensional in
($d + 1$)-dimensional space):
 - standard d -simplex
 - multi-pivot quicksort polyhedra

Computing: Polynomial Functions

Integral

$$\int_{\mathcal{C}} f(x) \, d\lambda_d(x)$$

- multivariate polynomial f over rationals \mathbb{Q}
- full-dimensional polytope \mathcal{C} with rational vertices
 - software package *LattE integrale*
[Baldoni–Berline–DeLoera–Koeppel–Vergne 2011,
DeLoera–Dutra–Koeppel–Moreinis–Pinto–Wu 2013]
 - interfaced by the mathematics software *SageMath*
 - result is again rational



Computing: Polynomial Functions

Integral

$$\int_{\mathcal{C}} f(x) d\lambda_d(x)$$

- multivariate polynomial f over rationals \mathbb{Q}
- full-dimensional polytope \mathcal{C} with rational vertices
 - software package *LattE integrale*
[Baldoni–Berline–DeLoera–Koepppe–Vergne 2011,
DeLoera–Dutra–Koepppe–Moreinis–Pinto–Wu 2013]
 - interfaced by the mathematics software *SageMath*
 - result is again rational
- not full-dimensional polytope \mathcal{C} with rational vertices
 - project down to its affine hull
 - integrate over this new polytope



Computing: Logarithms

Integral

$$\int_{\mathcal{C}} f(x) d\lambda_d(x)$$

- $f(x) = p(x)(\ell(x) \log(\ell(x)))^m$
- multivariate polynomial p
- linear multivariate polynomial ℓ
- positive integer m

Computing: Logarithms

Integral

$$\int_{\mathcal{C}} f(x) d\lambda_d(x)$$

- $f(x) = p(x)(\ell(x) \log(\ell(x)))^m$
- multivariate polynomial p
- linear multivariate polynomial ℓ
- positive integer m

- higher-dimensional integration by parts

$$\begin{aligned} & \int_{\mathcal{C}} p(x)(\ell(x) \log(\ell(x)))^m d\lambda_{\delta}(x) \\ &= \sum_{\substack{(\delta-1)\text{-dimensional} \\ \text{face } \mathcal{B} \text{ of } \mathcal{C}}} \frac{n_j(\mathcal{B})}{b} \int_{\mathcal{B}} q(x)(\ell(x) \log(\ell(x)))^m d\lambda_{\delta-1}(x) \\ & \quad - \int_{\mathcal{C}} m q(x)(\ell(x) \log(\ell(x)))^{m-1} d\lambda_{\delta}(x) \end{aligned}$$

with $q(x)$ such that $\frac{\partial}{\partial \ell(x)} (q(x)\ell(x)^m) = p(x)\ell(x)^m$

Computing: More Logarithms

- What about products of logarithms?



$$\mathbb{E}\left(\left(\sum_{i=0}^d D_i \log D_i\right)^2\right)$$
$$= ???$$

Computing: More Logarithms

- What about products of logarithms?



$$\mathbb{E}\left(\left(\sum_{i=0}^d D_i \log D_i\right)^2\right)$$

$$= (H_{d+1}^{(2)} - 1) + (H_{d+1} - 1)^2 + \frac{d}{d+2} \left(1 - \frac{\pi^2}{6}\right)$$

with

- harmonic numbers $H_m = \sum_{i=1}^m 1/i$
- $H_m^{(2)} = \sum_{i=1}^m 1/i^2$

Variance



Theorem (Holmgren–K 2019+)

*variance of the
number of key comparisons C_n
for sorting a list of n elements
with the optimal d -pivot
quicksort algorithm
is $\sim \sigma_d^2 n^2$
as $n \rightarrow \infty$*

- $\sigma_1^2 = 7 - \frac{2}{3}\pi^2 = 0.4202637326071\dots$
- $\sigma_2^2 = \frac{1609}{300} - \frac{27}{50}\pi^2 + \frac{3}{10}\log 2 = 0.2416911109130\dots$
- $\sigma_3^2 = \frac{3051169}{657072} - \frac{17689}{36504}\pi^2 + \frac{1463}{2808}\log 2 - \frac{665}{8424}\log 3 = 0.1354122412131\dots$
- $\sigma_4^2 = \frac{54171532769}{5002593750} - \frac{45467648}{100051875}\pi^2 + \frac{594228688}{487265625}\log 2 + \frac{798044}{2165625}\log 3 - \frac{65057423}{97453125}\log 5$