

# Process Convergence of the QuickSelect Residual

Jasper Ischebeck  
Universität Frankfurt

Bath 2024

# QuickSelect

- *Problem:* Given  $n$  distinct elements  $U_1; \dots; U_n$  and  $1 \leq k \leq n$ , find the element at rank  $k$ , i.e.  $U_{(k)}$  such that

$$\#\{i \mid U_i \leq U_{(k)}\} = k$$

- Sorting all values would take  $O(n \log n)$  time
- QuickSelect, derived from QuickSort, needs expected  $O(n)$  time

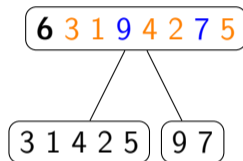
# QuickSelect – Example

Find the fourth-biggest element

6 3 1 9 4 2 7 5

# QuickSelect – Example

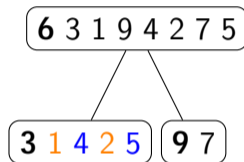
Find the fourth-biggest element



- Compare with 6 (the *pivot*) and sort into elements bigger and smaller

# QuickSelect – Example

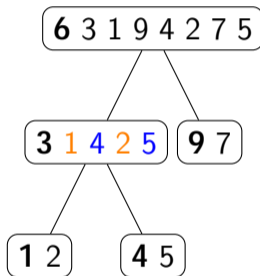
Find the fourth-biggest element



- Compare with 6 (the *pivot*) and sort into elements bigger and smaller
- New pivot 3

# QuickSelect – Example

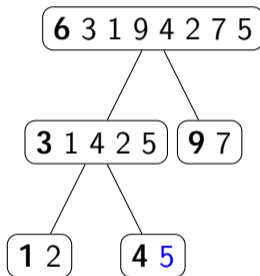
Find the fourth-biggest element



- Compare with 6 (the *pivot*) and sort into elements bigger and smaller
- New pivot 3

# QuickSelect – Example

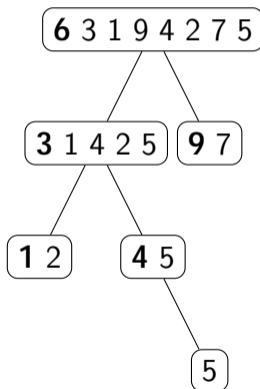
Find the fourth-biggest element



- Compare with 6 (the *pivot*) and sort into elements bigger and smaller
- New pivot 3
- Because 1, 2, 3 are smaller we look for the smaller of 4, 5

# QuickSelect – Example

Find the fourth-biggest element



- Compare with 6 (the *pivot*) and sort into elements bigger and smaller
- New pivot 3
- Because 1, 2, 3 are smaller we look for the smaller of 4, 5
- We find 4



# Analysis

- *Random model:* We assume  $U_1; \dots; U_n$  are independent and uniformly distributed on  $[0; 1]$ .
- Consider first the amount of **comparisons**
- But how to choose the rank for  $n \rightarrow \infty$ ?
- Uniformly at random on  $[1; n]$  (**grand average**)
- Rank  $btnc$  for some  $t \in [0; 1]$  (**QuickQuant process**)
- Let  $S_{t;n}$  be the amount of comparisons for rank  $btnc$
- $S_{t;n}/n$  converges for  $n \rightarrow \infty$  a.s. to a limit  $S_t$  (Grübel and Rösler 1996)

# Analysis

- *Random model:* We assume  $U_1; \dots; U_n$  are independent and uniformly distributed on  $[0; 1]$ .
- Consider first the amount of **comparisons**
- But how to choose the rank for  $n \rightarrow \infty$ ?
- Uniformly at random on  $[1; n]$  (**grand average**)
- Rank  $btnc$  for some  $t \in [0; 1]$  (**QuickQuant process**)
- Let  $S_{t;n}$  be the amount of comparisons for rank  $btnc$
- $S_{t;n}/n$  converges for  $n \rightarrow \infty$  a.s. to a limit  $S_t$  (Grübel and Rösler 1996)

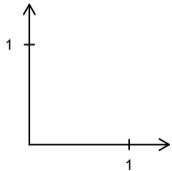
# Analysis

- *Random model:* We assume  $U_1; \dots; U_n$  are independent and uniformly distributed on  $[0; 1]$ .
- Consider first the amount of **comparisons**
- But how to choose the rank for  $n \rightarrow \infty$ ?
- Uniformly at random on  $[1; n]$  (**grand average**)
- Rank  $btnc$  for some  $t \in [0; 1]$  (**QuickQuant process**)
- Let  $S_{t;n}$  be the amount of comparisons for rank  $btnc$
- $S_{t;n}/n$  converges for  $n \rightarrow \infty$  a.s. to a limit  $S_t$  (Grübel and Rösler 1996)

# Analysis

- *Random model:* We assume  $U_1; \dots; U_n$  are independent and uniformly distributed on  $[0; 1]$ .
- Consider first the amount of **comparisons**
- But how to choose the rank for  $n \rightarrow \infty$ ?
- Uniformly at random on  $[1; n]$  (**grand average**)
- Rank  $btnc$  for some  $t \in [0; 1]$  (**QuickQuant process**)
- Let  $S_{t;n}$  be the amount of comparisons for rank  $btnc$
- $S_{t;n}/n$  converges for  $n \rightarrow \infty$  a.s. to a limit  $S_t$  (Grübel and Rösler 1996)

# The Limit Process



# The Limit Process

# The Limit Process

# The Limit Process



# The Limit Process

# The Limit Process

# The Limit Process

# The Limit Process

# The Limit Process

# The Limit Process

# The Limit Process

# The Limit Process



# The Limit Process

# The Limit Process

# The Limit Process

# The Limit Process

# The residual

What happens if we subtract the limit and rescale? We call this **residual**.

# The residual

What happens if we subtract the limit and rescale? We call this **residual**.

## Theorem (I., Neininger 2024+)

Let  $F_n := \frac{1}{n} \sum_{i=1}^n 1_{[U_i; 1]}$  be the empirical distribution function of  $U_1, \dots, U_n$ . The **residual process** converges

$$G_{t;n} := \frac{S_{t;n} - nS_{F_n^{-1}(t)}}{\sqrt{n}} \stackrel{d}{\rightarrow} G_{t;1} \quad \text{in } (D[0; 1]; d_{SK})$$

towards a mixed centred Gaussian process  $G_{t;1}$ .

What is  $D[0; 1]$ ? Why the empirical distribution function? And what are the covariances?

# Functional Results

- We study  $S_{t;n}$  resp.  $G_{t;n}$  as process to be able to consider all choices of random places simultaneously, .
- Many properties, e.g. the amount of comparisons at fixed and random places of the maximum can be written as functions of the process  $S_{t;n}$

## Theorem (Continuous mapping theorem)

For a measurable function

$$f: D \rightarrow \mathbb{R}^d \quad (X_n) \xrightarrow{d} (X) \implies f(X_n) \xrightarrow{d} f(X)$$

if  $X$  is a.s. not at a discontinuity of  $f$ .

- Process convergence  $S_{t;n}$  then implies convergence of these properties
- $S_{t;n}$  converges to  $S_t$  as process (Grubel and Rosler 1996)

# Functional Results

- We study  $S_{t;n}$  resp.  $G_{t;n}$  as process to be able to consider all choices of random places simultaneously, .
- Many properties, e.g. the amount of comparisons at fixed and random places of the maximum can be written as functions of the process  $S_{t;n}$

## Theorem (Continuous mapping theorem)

For a measurable function

$$X_n \xrightarrow{d} X \quad \Rightarrow \quad (X_n) \xrightarrow{d} (X)$$

if  $X$  is a.s. not at a discontinuity of  $f$ .

- Process convergence  $S_{t;n}$  then implies convergence of these properties
- $S_{t;n}$  converges to  $S_t$  as process (Grubel and Rosler 1996)



# Functional Results

- We study  $S_{t;n}$  resp.  $G_{t;n}$  as process to be able to consider all choices of random places simultaneously, .
- Many properties, e.g. the amount of comparisons at fixed and random places of the maximum can be written as functions of the process  $S_{t;n}$

## Theorem (Continuous mapping theorem)

For a measurable function

$$X_n \xrightarrow{d} X \quad \Rightarrow \quad f(X_n) \xrightarrow{d} f(X)$$

if  $X$  is a.s. not at a discontinuity of  $f$ .

- Process convergence  $S_{t;n}$  then implies convergence of these properties
- $S_{t;n}$  converges to  $S_t$  as process (Grubel and Rosler 1996)

# The space of càdlàg functions

- But first, we have to describe the space the process lives on
- $S_{t;n}$  is a right-continuous step function
- Functions that are right-continuous and have left limits are called càdlàg
- Write  $D[0; 1]$  for the space of càdlàg functions on  $[0; 1]$
- For measurability we need a metric accommodating for jumps not aligning, the Skorokhod metric

# The space of cadlag functions

- But first, we have to describe the space the process lives on
- $S_{t;n}$  is a right-continuous step function
- Functions that are right-continuous and have left limits are called cadlag
- Write  $D[0; 1]$  for the space of cadlag functions on  $[0; 1]$
- For measurability we need a metric accommodating for jumps not aligning, the Skorokhod metric

# The space of càdlàg functions

- The **Skorokhod metric** uses a monotonously growing bijection  $[0, 1] \rightarrow [0, 1]$  to align jumps.
- For  $f, g \in D[0, 1]$  define

$$d_{SK}(f; g) := \inf_k \|f \circ \gamma - g\|_{k_1} \quad \|\cdot\|_{k_1} = \|\text{id} \cdot \cdot\|_{k_1}$$

(id is the identity,  $\|\cdot\|_{k_1}$  the maximum,  $\|\cdot\|_{k_1}$  the uniform norm)

- Example: For  $a, b \in (0, 1)$ :  $d_{SK}(1_{[a;1)}; 1_{[b;1)}) = |b - a|$

# The space of càdlàg functions

- The **Skorokhod metric** uses a monotonously growing bijection  $[0, 1] \rightarrow [0, 1]$  to align jumps.
- For  $f, g \in D[0, 1]$  define

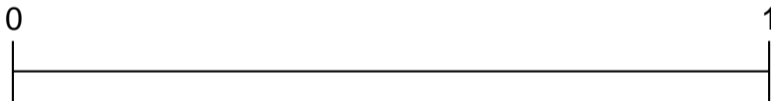
$$d_{SK}(f; g) := \inf \{k_f \circ g \circ k_1 \circ \text{id} \circ k_1\}$$

(id is the identity,  $\circ$  the maximum,  $k_1$  the uniform norm)

- Example: For  $a, b \in (0, 1)$ :  $d_{SK}(1_{[a;1)}; 1_{[b;1)}) = |b - a|$

# QuickVal

- Starting at the full interval  $[0, 1]$ , at step 0 the first pivot divides the interval in left and right
- The first value in each interval becomes the new pivot and splits the interval again
- Write  $l_k$  for the length of the interval in step  $k$  containing a value



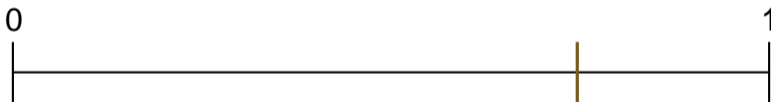
# QuickVal

- Starting at the full interval  $[0,1]$ , at step 0 the  $r$ st pivot divides the interval in left and right
- The  $r$ st value in each interval becomes the new pivot and splits the interval again
- Write  $l_{i,k}$  for the length of the interval in step  $k$  containing a value



# QuickVal

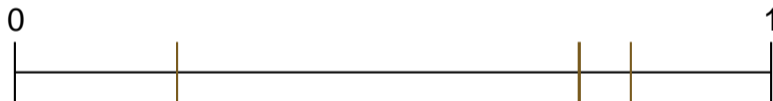
- Starting at the full interval  $[0, 1]$ , at step 0 the first pivot divides the interval in left and right
- The first value in each interval becomes the new pivot and splits the interval again
- Write  $l_{i,k}$  for the length of the interval in step  $k$  containing a value





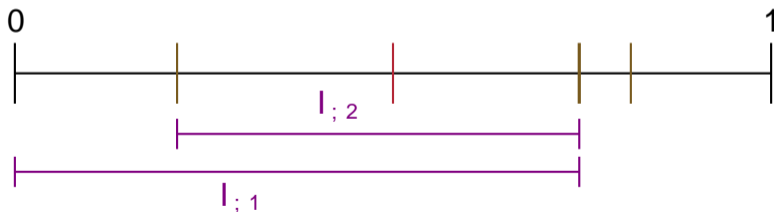
# QuickVal

- Starting at the full interval  $[0, 1]$ , at step 0 the first pivot divides the interval in left and right
- The first value in each interval becomes the new pivot and splits the interval again
- Write  $l_{i,k}$  for the length of the interval in step  $k$  containing a value



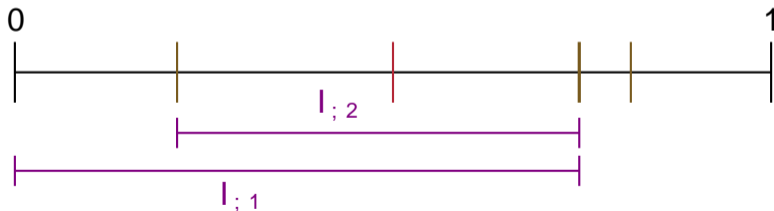
# QuickVal

- Starting at the full interval  $[0, 1]$ , at step 0 the first pivot divides the interval in left and right
- The first value in each interval becomes the new pivot and splits the interval again
- Write  $l; k$  for the length of the interval in step  $k$  containing a value



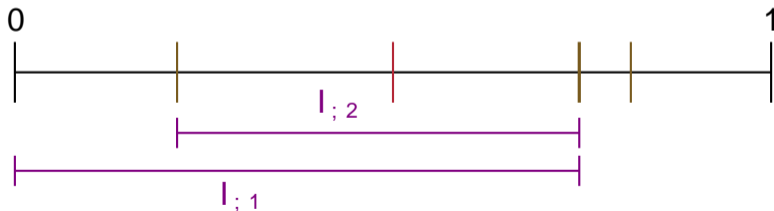
# QuickVal

- We always go into the interval where **the result** of our algorithm is
- The result is  $F_n^{-1}(t)$ , that is why the inverse empirical distribution function pops
- We should have indexed using the result instead of using  $i$ !
- The process  $S_{;n}$  of comparisons for result  $t$  is called **QuickVal**



# QuickVal

- We always go into the interval where **the result** of our algorithm is
- The result is  $F_n^{-1}(t)$ , that is why the inverse empirical distribution function pops
- We should have indexed using the result instead of using  $i$ !
- The process  $S_{;n}$  of comparisons for result is called **QuickVal**



# QuickVal

- Let  $S_{k;n}$  the amount of comparisons in step  $2 N_0$
- That is the amount of elements in the interval with length  $l_{k,n}$
- Conditional on the pivots  $S_{k;n}$  is thus almost  $B(n; l_{k,n})$ -distributed
- The error are the  $k$  pivots, so we write  $S_{k;n} = B(n; l_{k,n}) + O(k)$ .

# QuickVal

- Let  $S_{k;n}$  the amount of comparisons in step  $N_0$
- That is the amount of elements in the interval with length  $l_{k,n}$
- Conditional on the pivots  $S_{k;n}$  is thus almost  $B(n; l_{k,n})$ -distributed
- The error are the  $k$  pivots, so we write  $S_{k;n} = B(n; l_{k,n}) + O(k)$ .

# The limit process

- Since  $S_{k;n} = B(n; I_{k;n}) + O(k)$ ,

$$\frac{S_{k;n}}{n} \rightarrow I_{k;n} \text{ a.s.}$$

- Indeed, our limit process  $S$  is given by

$$S := \sum_{k \geq 1} X_{k;n} I_{k;n}$$

and  $n^{-1} S_{n; \cdot} \rightarrow S$  both a.s. (Fill and Nakama 2013) and  $L^2$  (Fill and Matterer 2014).

# The limit process

- Since  $S_{k;n} = B(n; I_{k;n}) + O(k)$ ,

$$\frac{S_{k;n}}{n} \rightarrow I_{k;n} \text{ a.s.}$$

- Indeed, our limit process  $S$  is given by

$$S := \sum_{k \in \mathbb{N}} X_{k;n} I_{k;n}$$

and  $n^{-1} S_{n; \cdot} \rightarrow S$  both a.s. (Fill and Nakama 2013) and  $L^2$  (Fill and Matterer 2014).



# Variances

What about the variation around this limit, the residual?

- Conditional on  $l_{i;k}$ ,

$$\frac{S_{i;k;n}}{p} - \frac{n l_{i;k}}{\bar{n}} \stackrel{!}{\sim} N(0; l_{i;k}(1 - l_{i;k}))$$

- As for the sum over all steps (Matterer 2015)

$$\frac{S_{i;n}}{p} - \frac{nS}{\bar{n}} \stackrel{!}{\sim} N(0; \dots)$$

- If an element is in  $l_{i;k}$ , then it is also in  $l_{i;l}$  for all  $l \leq k$ , so

$$S_{i;n} = \sum_{k;l \geq N_0} l_{i;k-1} - l_{i;k} l_{i;l}$$

# Variances

What about the variation around this limit, the residual?

- Conditional on  $l; k$ ,

$$\frac{S_{; k; n}}{p \frac{n l; k}{\bar{n}}} \stackrel{!}{\sim} N(0; l; k(1 - l; k))$$

- As for the sum over all steps (Matterer 2015)

$$\frac{S_{; n}}{p \frac{n S}{\bar{n}}} \stackrel{!}{\sim} N(0; ; )$$

- If an element is in  $l; k$ , then it is also in  $l; l$  for all  $l \leq k$ , so

$$S_{; k} = \sum_{l; l \leq k} S_{; l}$$

# Variances

What about the variation around this limit, the residual?

- Conditional on  $l_{i;k}$ ,

$$\frac{S_{i;k;n}}{p} - \frac{n l_{i;k}}{\bar{n}} \stackrel{!}{=} N(0; l_{i;k}(1 - l_{i;k}))$$

- As for the sum over all steps (Matterer 2015)

$$\frac{S_{i;n}}{p} - \frac{nS}{\bar{n}} \stackrel{!}{=} N(0; \dots)$$

- If an element is in  $l_{i;k}$ , then it is also in  $l_{i;l}$  for all  $l \leq k$ , so

$$S_{i;n} = \sum_{k;l \geq 2N_0} l_{i;k-1} - l_{i;k} l_{i;l}$$

# Covariances

- For  $\alpha \in [0; 1]$ , let  $J(\alpha)$  be the (random) last index where  $\alpha$  and  $\beta$  are in the same interval
- An element is in  $I_{\alpha; k}$  and  $I_{\alpha; l}$  for  $k, l \in \mathbb{N}_0$  at the same time if and only if  $k \leq J(\alpha)$  and the element is in  $I_{\alpha; l}$
- Combining interval and interval length, let us write

$$|I_{\alpha; k} \cap I_{\alpha; l}|$$

for the length of the intersection of  $I_{\alpha; k}$  and  $I_{\alpha; l}$

- Then, the covariances can be written as

$$\rho_{k,l} = \sum_{\alpha \in [0; 1]} |I_{\alpha; k} \cap I_{\alpha; l}| \cdot |I_{\alpha; k}| \cdot |I_{\alpha; l}|$$

# Main theorem

We can now restate our main theorem with **QuickVal**

## Theorem (I., Neininger 2024+)

The **residual process** converges

$$\frac{S_{;n} - nS}{\sqrt{n}} \xrightarrow{d} G_{;1} \quad \text{in } (D[0;1]; d_{SK})$$

towards a mixed centred Gaussian process  $G_{;1}$  with covariances

$$G_{;1}(j) \cdot G_{;1}(k) = \int_0^{\min(j,k)} \mathbb{1}_{[0;1]}(t) dt$$

# Proof Sketch

- Split for some  $K = K_n \ll N$

$$\frac{S_{;n}}{p \bar{n}} = \sum_{k=0}^{K} \frac{S_{;k;n}}{p \bar{n}} + \sum_{k=K+1}^{b:5 \log nc} \frac{S_{;k;n}}{p \bar{n}} + \sum_{k=d:5 \log ne} \frac{S_{;k;n}}{p \bar{n}}$$

- The first steps have only finitely many values, use the (Multivariate) Central Limit Theorem
- The following steps have few elements, and should be small, use Chernov bound
- In the last steps there are no elements, and is falling geometrically.

# Proof Sketch

- Split for some  $K = K_n \ll N$

$$\frac{S_{;n}}{p \bar{n}} = \sum_{k=0}^{K} \frac{S_{;k;n}}{p \bar{n}} + \sum_{k=K+1}^{b:5 \log nc} \frac{S_{;k;n}}{p \bar{n}} + \sum_{k=d:5 \log ne} \frac{S_{;k;n}}{p \bar{n}}$$

- The first steps have only finitely many values, use the (Multivariate) Central Limit Theorem
- The following steps have few elements, and should be small, use Chernov bound
- In the last steps there are no elements, and is falling geometrically.

# Proof Sketch

- Split for some  $K = K_n \ll N$

$$\frac{S_{;n}}{p \bar{n}} = \sum_{k=0}^{K} \frac{S_{;k;n}}{p \bar{n}} + \sum_{k=K+1}^{b:5 \log n c} \frac{S_{;k;n}}{p \bar{n}} + \sum_{k=d:5 \log n e} \frac{S_{;k;n}}{p \bar{n}}$$

- The first steps have only finitely many values, use the (Multivariate) Central Limit Theorem
- The following steps have few elements, and should be small, use Chernov bound
- In the last steps there are no elements, and is falling geometrically.



# Proof Sketch

- Split for some  $K = K_n \leq N$

$$\frac{S_{;n}}{p \bar{n}} = \sum_{k=0}^{K} \frac{S_{;k;n}}{p \bar{n}} + \sum_{k=K+1}^{b:5 \log n c} \frac{S_{;k;n}}{p \bar{n}} + \sum_{k=d:5 \log n e} \frac{S_{;k;n}}{p \bar{n}}$$

- The first steps have only finitely many values, use the (Multivariate) Central Limit Theorem
- The following steps have few elements, and should be small, use Chernov bound
- In the last steps there are no elements, and is falling geometrically.

# Swaps

- The step to partition into smaller and bigger elements is usually done by rearranging
- Depending on machine, swapping positions can be significantly slower than comparing
- We consider 2 schemes: Hoare and Lomuto

# Swaps

- The step to partition into smaller and bigger elements is usually done by rearranging
- Depending on machine, swapping positions can be significantly slower than comparing
- We consider 2 schemes: **Hoare** and **Lomuto**

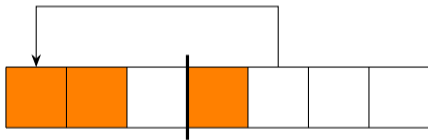
# Lomuto and Hoare { Visualisation

Lomuto



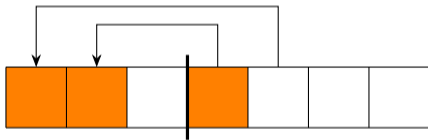
# Lomuto and Hoare { Visualisation

Lomuto



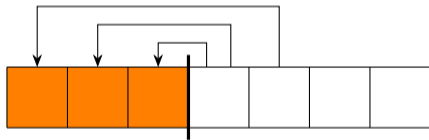
# Lomuto and Hoare { Visualisation

Lomuto



# Lomuto and Hoare { Visualisation

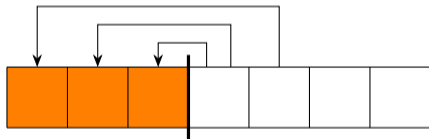
## Lomuto



**Cost:** Amount of smaller elements

# Lomuto and Hoare { Visualisation

Lomuto



**Cost:** Amount of smaller elements

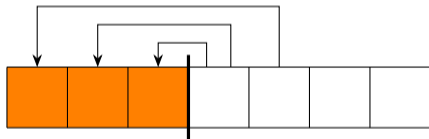
Hoare





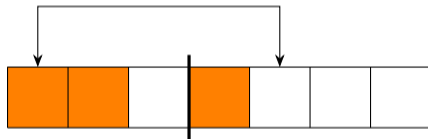
# Lomuto and Hoare { Visualisation

Lomuto



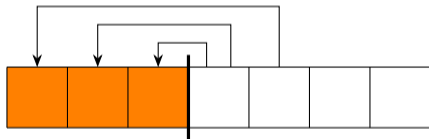
**Cost:** Amount of smaller elements

Hoare



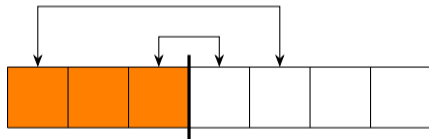
# Lomuto and Hoare { Visualisation

Lomuto



**Cost:** Amount of smaller elements

Hoare



**Cost:** Amount of smaller elements  
on the right side  
(hypergeometric)

# Swap Residual

- Given the interval sizes, the amount  $W_{j;k;n}$  of swaps for either **Hoare** or **Lomuto** at a fixed step  $k$  are asymptotically normal around some mean  $W_{j;k}$ , a function of interval sizes
- Using this limit, we again define a residual process

$$\overset{\times}{\rightarrow} \frac{W_{j;k;n} - nW_{j;k}}{\sqrt{n}}$$

$k=0$

# Swap Residual

## Theorem (I., Neininger 2024+)

For the *Hoare* scheme the *residual process for swaps* converges

$$\overset{\times}{\underset{k=0}{\rightarrow}} \frac{W_{;k;n} \rho \frac{nW_{;k}}{\bar{n}}}{\rho} \overset{!}{\rightarrow} G_{;1}^{\text{swap}} \text{ in } (D[0;1]; d_{\text{SK}})$$

towards a mixed centred Gaussian process  $G_{;1}^{\text{swap}}$ . The random covariances are functions of the interval sizes. The same holds for *Lomuto*, with other covariances.

# Variable comparison costs

- Other model: Comparison times depend on the items compared
- E.g. the time needed to compare strings is proportional to the length of their common prefix
- The same holds for decimals (number of bit comparisons)
- Typically, closer elements take longer to compare

# Variable comparison costs

- Let  $(u; v)$  be cost to compare pivot  $u$  to item  $v$
- Let  $V$  be uniformly distributed on  $[0; 1]$
- We call  $(u; V)$   $\epsilon$ -tame for  $\epsilon > 0$  if there exists a  $C > 0$  so that for all  $x \in [0; u] \subseteq [0; 1]$

$$P((u; V) > x) \leq Cx^{-\epsilon}$$

- Sufficient that  $(u; V)$  has a  $\epsilon$ -th moment, uniformly bounded in  $u$ .
- This covers bit comparisons, where  $(u; V)$  has exponential tails, so  $(u; V)$  is  $\epsilon$ -tame for all  $\epsilon > 0$

# Variable comparison costs

- Let  $(u; v)$  be cost to compare pivot  $u$  to item  $v$
- Let  $V$  be uniformly distributed on  $[0; 1]$
- We call  $(u; V)$  **"-tame** for  $" > 0$  if there exists a  $C > 0$  so that for all  $x \in [0; u] \subseteq [0; 1]$

$$P((u; V) > x) \leq Cx^{-" - 1}$$

- Sufficient that  $(u; V)$  has a  $" - 1$ -th moment, uniformly bounded in  $u$ .
- This covers bit comparisons, where  $(u; V)$  has exponential tails, so  $(u; V)$  is  $"$ -tame for all  $" > 0$

# Variable comparison costs

- Let  $c(u; v)$  be cost to compare pivot  $u$  to item  $v$
- Let  $V$  be uniformly distributed on  $[0; 1]$
- We call  $c$  **"-tame** for  $\epsilon > 0$  if there exists a  $C > 0$  so that for all  $x \in [0; u] \subseteq [0; 1]$

$$P(c(u; V) > x) \leq Cx^{-\epsilon}$$

- Sufficient that  $c(u; V)$  has a  $\epsilon^{-1}$ -th moment, uniformly bounded in  $u$ .
- This covers bit comparisons, where  $c(u; V)$  has exponential tails, so  $c$  is  $\epsilon$ -tame for all  $\epsilon > 0$



# Residual

The residual can be defined as before, and then

## Theorem (Matterer 2015)

*When  $\beta$  is  $\epsilon$ -tame for  $\epsilon < \frac{1}{2}$ , the residual at a fixed point converges in distribution to a centred mixed normal.*

## Theorem (I., Neininger 2024+)

*When  $\beta$  is  $\epsilon$ -tame for  $\epsilon < \frac{1}{4}$ , the residual process converges in distribution on  $(D[0;1]; d_{SK})$  to a centred mixed Gaussian process with covariances given as functions of  $\beta$  and the interval sizes.*

# Residual

The residual can be defined as before, and then

## Theorem (Matterer 2015)

*When  $\beta$  is  $\epsilon$ -tame for  $\epsilon < \frac{1}{2}$ , the residual at a fixed point converges in distribution to a centred mixed normal.*

## Theorem (I., Neininger 2024+)

*When  $\beta$  is  $\epsilon$ -tame for  $\epsilon < \frac{1}{4}$ , the residual process converges in distribution on  $(D[0;1]; d_{SK})$  to a centred mixed Gaussian process with covariances given as functions of  $\beta$  and the interval sizes.*

