

# Wavelet Trees Meet Suffix Trees

Maxim Babenko<sup>1</sup>   Paweł Gawrychowski<sup>2→3</sup>  
**Tomasz Kociumaka**<sup>3</sup>   Tatiana Starikovskaya<sup>1→4</sup>

<sup>1</sup>National Research University Higher School of Economics, Moscow, Russia

<sup>2</sup>Max-Planck-Institut für Informatik, Saarbrücken, Germany

<sup>3</sup>University of Warsaw, Poland

<sup>4</sup>University of Bristol, UK

**SODA 2015**

San Diego, California, USA

January 4, 2015

# Rank and Select Queries

## General Setting

**Universe** totally ordered universe  $U$  (integers, strings, ...)

**Input** multiset  $A$  over universe  $U$

$rank_A(u)$  count elements in  $A$  not exceeding  $u$

$select_A(k)$  return the  $k$ -th smallest element in  $A$

# Rank and Select Queries

## General Setting

**Universe** totally ordered universe  $U$  (integers, strings, ...)

**Input** multiset  $A$  over universe  $U$

$rank_A(u)$  count elements in  $A$  not exceeding  $u$

$select_A(k)$  return the  $k$ -th smallest element in  $A$

Example: an integer multiset

$$A = \{0, 1, 3, 4, 5, 5, 8, 8, 9, 9\}$$

# Rank and Select Queries

## General Setting

**Universe** totally ordered universe  $U$  (integers, strings, ...)

**Input** multiset  $A$  over universe  $U$

$rank_A(u)$  count elements in  $A$  not exceeding  $u$

$select_A(k)$  return the  $k$ -th smallest element in  $A$

Example: an integer multiset

$$A = \{0, 1, 3, 4, 5, 5, 8, 8, 9, 9\}$$

$$rank_A(7) = 6$$

# Rank and Select Queries

## General Setting

**Universe** totally ordered universe  $U$  (integers, strings, ...)

**Input** multiset  $A$  over universe  $U$

$rank_A(u)$  count elements in  $A$  not exceeding  $u$

$select_A(k)$  return the  $k$ -th smallest element in  $A$

Example: an integer multiset

$$A = \{0, 1, 3, 4, 5, 5, 8, 8, 9, 9\}$$

$$rank_A(7) = 6$$

$$select_A(7) = 8$$

## General Setting

**Universe** totally ordered universe  $U$  (integers, strings, ...)

**Input** multiset  $A$  over universe  $U$

$rank_A(u)$  count elements in  $A$  not exceeding  $u$

$select_A(k)$  return the  $k$ -th smallest element in  $A$

Example: an integer multiset

$$A = \{0, 1, 3, 4, 5, 5, 8, 8, 9, 9\}$$

$$rank_A(7) = 6 \qquad select_A(7) = 8$$

Example: a set encoded as bitmask

$$A = \{2, 3, 5, 9, 10, 11, 13, 14, 15, 16, 17, 18, 20\}$$

0110100011101111101

## General Setting

**Universe** totally ordered universe  $U$  (integers, strings, ...)

**Input** multiset  $A$  over universe  $U$

$rank_A(u)$  count elements in  $A$  not exceeding  $u$

$select_A(k)$  return the  $k$ -th smallest element in  $A$

Example: an integer multiset

$$A = \{0, 1, 3, 4, 5, 5, 8, 8, 9, 9\}$$

$$rank_A(7) = 6 \quad select_A(7) = 8$$

Example: a set encoded as bitmask

$$A = \{2, 3, 5, 9, 10, 11, 13, 14, 15, 16, 17, 18, 20\}$$

01101000111011111101

$$rank_1(7) = 3$$

# Rank and Select Queries

## General Setting

**Universe** totally ordered universe  $U$  (integers, strings, ...)

**Input** multiset  $A$  over universe  $U$

$rank_A(u)$  count elements in  $A$  not exceeding  $u$

$select_A(k)$  return the  $k$ -th smallest element in  $A$

Example: an integer multiset

$$A = \{0, 1, 3, 4, 5, 5, 8, 8, 9, 9\}$$

$$rank_A(7) = 6 \quad select_A(7) = 8$$

Example: a set encoded as bitmask

$$A = \{2, 3, 5, 9, 10, 11, 13, 14, 15, 16, 17, 18, 20\}$$

$$01101000111011111101$$

$$rank_1(7) = 3 \quad select_1(7) = 13$$



## Range Rank and Select Queries

**Input** An integer sequence  $S = (S_1, \dots, S_n)$ .

**Queries** Rank and select on a range  $S_{l..r} = \{S_l, S_{l+1}, \dots, S_r\}$ .

## Range Rank and Select Queries

**Input** An integer sequence  $S = (S_1, \dots, S_n)$ .

**Queries** Rank and select on a range  $S_{l..r} = \{S_l, S_{l+1}, \dots, S_r\}$ .

Special case:

- Range Minimum Queries (RMQ).

## Range Rank and Select Queries

**Input** An integer sequence  $S = (S_1, \dots, S_n)$ .

**Queries** Rank and select on a range  $S_{l..r} = \{S_l, S_{l+1}, \dots, S_r\}$ .

Special case:

- Range Minimum Queries (RMQ).

Example:

$S$ : 

2	6	0	7	0	0	2	1	4	1	8	1	6	7	4	2	9	3	0	5
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

  
1 20

## Range Rank and Select Queries

**Input** An integer sequence  $S = (S_1, \dots, S_n)$ .

**Queries** Rank and select on a range  $S_{l..r} = \{S_l, S_{l+1}, \dots, S_r\}$ .

Special case:

- Range Minimum Queries (RMQ).

Example:

S:	2	6	0	7	0	0	2	1	4	1	8	1	6	7	4	2	9	3	0	5
	1			5								14								20

$$\text{rank}_{5..14}(6) =$$

## Range Rank and Select Queries

**Input** An integer sequence  $S = (S_1, \dots, S_n)$ .

**Queries** Rank and select on a range  $S_{l..r} = \{S_l, S_{l+1}, \dots, S_r\}$ .

Special case:

- Range Minimum Queries (RMQ).

Example:

S:	2	6	0	7	0	0	2	1	4	1	8	1	6	7	4	2	9	3	0	5
	1			5									14							20

$$\text{rank}_{5..14}(6) = 8$$

## Range Rank and Select Queries

**Input** An integer sequence  $S = (S_1, \dots, S_n)$ .

**Queries** Rank and select on a range  $S_{l..r} = \{S_l, S_{l+1}, \dots, S_r\}$ .

Special case:

- Range Minimum Queries (RMQ).

Example:

S:	2	6	0	7	0	0	2	1	4	1	8	1	6	7	4	2	9	3	0	5
	1						9									17		20		

$$\text{rank}_{5..14}(6) = 8$$

$$\text{select}_{9..17}(5) =$$

## Range Rank and Select Queries

**Input** An integer sequence  $S = (S_1, \dots, S_n)$ .

**Queries** Rank and select on a range  $S_{l..r} = \{S_l, S_{l+1}, \dots, S_r\}$ .

Special case:

- Range Minimum Queries (RMQ).

Example:

$S$ : 

2	6	0	7	0	0	2	1	4	1	8	1	6	7	4	2	9	3	0	5
1							9								17				20

$$\text{rank}_{5..14}(6) = 8$$

$$\text{select}_{9..17}(5) = 4.$$

# Range Rank and Selection: Results

Problem	Space	Query	Construction	Reference
Rank (negative)	$\mathcal{O}(n)$	$\mathcal{O}\left(\frac{\log n}{\log \log n}\right)$	$\mathcal{O}(n\sqrt{\log n})$	Chan & Pătrașcu SODA 2010
	$n \log^{\mathcal{O}(1)} n$	$\Omega\left(\frac{\log n}{\log \log n}\right)$	—	Pătrașcu STOC 2007



# Range Rank and Selection: Results

Problem	Space	Query	Construction	Reference
Rank (negative)	$\mathcal{O}(n)$	$\mathcal{O}\left(\frac{\log n}{\log \log n}\right)$	$\mathcal{O}(n\sqrt{\log n})$	Chan & Pătraşcu SODA 2010
	$n \log^{\mathcal{O}(1)} n$	$\Omega\left(\frac{\log n}{\log \log n}\right)$	–	Pătraşcu STOC 2007
Select (negative)	$\mathcal{O}(n)$	$\mathcal{O}(\log n)$	$\mathcal{O}(n\sqrt{\log n})$	Chan & Pătraşcu SODA 2010
	$\mathcal{O}(n)$	$\mathcal{O}\left(\frac{\log n}{\log \log n}\right)$	$\mathcal{O}(n \log n)$	Brodal et al. ICALP 2009
	$n \log^{\mathcal{O}(1)} n$	$\Omega\left(\frac{\log n}{\log \log n}\right)$	–	Jørgensen & Larsen SODA 2011

# Range Rank and Selection: Results

Problem	Space	Query	Construction	Reference
Rank (negative)	$\mathcal{O}(n)$	$\mathcal{O}\left(\frac{\log n}{\log \log n}\right)$	$\mathcal{O}(n\sqrt{\log n})$	Chan & Pătraşcu SODA 2010
	$n \log^{\mathcal{O}(1)} n$	$\Omega\left(\frac{\log n}{\log \log n}\right)$	–	Pătraşcu STOC 2007
Select (negative)	$\mathcal{O}(n)$	$\mathcal{O}(\log n)$	$\mathcal{O}(n\sqrt{\log n})$	Chan & Pătraşcu SODA 2010
	$\mathcal{O}(n)$	$\mathcal{O}\left(\frac{\log n}{\log \log n}\right)$	$\mathcal{O}(n \log n)$	Brodal et al. ICALP 2009
	$n \log^{\mathcal{O}(1)} n$	$\Omega\left(\frac{\log n}{\log \log n}\right)$	–	Jørgensen & Larsen SODA 2011
Select	$\mathcal{O}(n)$	$\mathcal{O}\left(\frac{\log n}{\log \log n}\right)$	$\mathcal{O}(n\sqrt{\log n})$	this work

# Range Rank and Selection: Results

Problem	Space	Query	Construction	Reference
Rank (negative)	$\mathcal{O}(n)$	$\mathcal{O}\left(\frac{\log n}{\log \log n}\right)$	$\mathcal{O}(n\sqrt{\log n})$	Chan & Pătrașcu SODA 2010
	$n \log^{\mathcal{O}(1)} n$	$\Omega\left(\frac{\log n}{\log \log n}\right)$	–	Pătrașcu STOC 2007
Select (negative)	$\mathcal{O}(n)$	$\mathcal{O}(\log n)$	$\mathcal{O}(n\sqrt{\log n})$	Chan & Pătrașcu SODA 2010
	$\mathcal{O}(n)$	$\mathcal{O}\left(\frac{\log n}{\log \log n}\right)$	$\mathcal{O}(n \log n)$	Brodal et al. ICALP 2009
	$n \log^{\mathcal{O}(1)} n$	$\Omega\left(\frac{\log n}{\log \log n}\right)$	–	Jørgensen & Larsen SODA 2011
Select	$\mathcal{O}(n)$	$\mathcal{O}\left(\frac{\log n}{\log \log n}\right)$	$\mathcal{O}(n\sqrt{\log n})$	this work

Theorem (independently: Munro, Nekrich, Vitter; SPIRE 2014)

*Wavelet trees can be constructed in  $\mathcal{O}(n\sqrt{\log n})$  time.*

## Suffix Rank and Selection Queries

**Universe** strings over  $\Sigma$  (denoted  $\Sigma^*$ )

**Input** a string  $T$  of length  $n$

**Queries** rank and select on the set of suffixes of  $T$  ( $Suf(T)$ ).

# Suffix Rank and Selection

## Suffix Rank and Selection Queries

**Universe** strings over  $\Sigma$  (denoted  $\Sigma^*$ )

**Input** a string  $T$  of length  $n$

**Queries** rank and select on the set of suffixes of  $T$  ( $Suf(T)$ ).

$T$	a	b	a	a	b	a	b	a	a	b	a
$i$	1	2	3	4	5	6	7	8	9	10	11

# Suffix Rank and Selection

## Suffix Rank and Selection Queries

**Universe** strings over  $\Sigma$  (denoted  $\Sigma^*$ )

**Input** a string  $T$  of length  $n$

**Queries** rank and select on the set of suffixes of  $T$  ( $Suf(T)$ ).

$T$	a	b	a	a	b	a	b	a	a	b	a
$i$	1	2	3	4	5	6	7	8	9	10	11
ISA	6	10	3	7	11	5	9	2	4	8	1

$i$	SA	
1	11	a
2	8	aaba
3	3	aababaaba
4	9	aba
5	6	abaaba
6	1	abaababaaba
7	4	ababaaba
8	10	ba
9	7	baaba
10	2	baababaaba
11	5	babaaba

# Suffix Rank and Selection

## Suffix Rank and Selection Queries

**Universe** strings over  $\Sigma$  (denoted  $\Sigma^*$ )

**Input** a string  $T$  of length  $n$

**Queries** rank and select on the set of suffixes of  $T$  ( $Suf(T)$ ).

$T$	a	b	a	a	b	a	b	a	a	b	a
$i$	1	2	3	4	5	6	7	8	9	10	11
ISA	6	10	3	7	11	5	9	2	4	8	1

$select_T(5) =$

$i$	SA	
1	11	a
2	8	aaba
3	3	aababaaba
4	9	aba
5	6	abaaba
6	1	abaababaaba
7	4	ababaaba
8	10	ba
9	7	baaba
10	2	baababaaba
11	5	babaaba

# Suffix Rank and Selection

## Suffix Rank and Selection Queries

**Universe** strings over  $\Sigma$  (denoted  $\Sigma^*$ )

**Input** a string  $T$  of length  $n$

**Queries** rank and select on the set of suffixes of  $T$  ( $Suf(T)$ ).

$T$	a	b	a	a	b	a	b	a	a	b	a
$i$	1	2	3	4	5	6	7	8	9	10	11
ISA	6	10	3	7	11	5	9	2	4	8	1

$$\text{select}_T(5) = T[\text{SA}[5]..] = T[6..] = \text{abaaba}$$

$i$	SA	
1	11	a
2	8	aaba
3	3	aababaaba
4	9	aba
5	6	abaaba
6	1	abaababaaba
7	4	ababaaba
8	10	ba
9	7	baaba
10	2	baababaaba
11	5	babaaba



# Suffix Rank and Selection

## Suffix Rank and Selection Queries

**Universe** strings over  $\Sigma$  (denoted  $\Sigma^*$ )

**Input** a string  $T$  of length  $n$

**Queries** rank and select on the set of suffixes of  $T$  ( $Suf(T)$ ).

$T$	a	b	a	a	b	a	b	a	a	b	a
$i$	1	2	3	4	5	6	7	8	9	10	11
ISA	6	10	3	7	11	5	9	2	4	8	1

$$\text{select}_T(5) = T[\text{SA}[5]..] = T[6..] = \text{abaaba}$$

$$\text{rank}_T(T[4..]) =$$

$i$	SA	
1	11	a
2	8	aaba
3	3	aababaaba
4	9	aba
5	6	abaaba
6	1	abaababaaba
7	4	ababaaba
8	10	ba
9	7	baaba
10	2	baababaaba
11	5	babaaba

# Suffix Rank and Selection

## Suffix Rank and Selection Queries

**Universe** strings over  $\Sigma$  (denoted  $\Sigma^*$ )

**Input** a string  $T$  of length  $n$

**Queries** rank and select on the set of suffixes of  $T$  ( $Suf(T)$ ).

$T$	a	b	a	a	b	a	b	a	a	b	a
$i$	1	2	3	4	5	6	7	8	9	10	11
ISA	6	10	3	7	11	5	9	2	4	8	1

$$\text{select}_T(5) = T[\text{SA}[5]..] = T[6..] = \text{abaaba}$$

$$\text{rank}_T(T[4..]) = \text{ISA}[4] = 7$$

$i$	SA	
1	11	a
2	8	aaba
3	3	aababaaba
4	9	aba
5	6	abaaba
6	1	abaababaaba
7	4	ababaaba
8	10	ba
9	7	baaba
10	2	baababaaba
11	5	babaaba

## Suffix Rank and Selection Queries

**Universe** strings over  $\Sigma$  (denoted  $\Sigma^*$ )

**Input** a string  $T$  of length  $n$

**Queries** rank and select on the set of suffixes of  $T$  ( $Suf(T)$ ).

$T$	a	b	a	a	b	a	b	a	a	b	a
$i$	1	2	3	4	5	6	7	8	9	10	11
ISA	6	10	3	7	11	5	9	2	4	8	1

$$\text{select}_T(5) = T[\text{SA}[5]..] = T[6..] = \text{abaaba}$$

$$\text{rank}_T(T[4..]) = \text{ISA}[4] = 7$$

$$\text{rank}_T(\text{aabb}) =$$

$i$	SA	
1	11	a
2	8	aaba
3	3	aababaaba
4	9	aba
5	6	abaaba
6	1	abaababaaba
7	4	ababaaba
8	10	ba
9	7	baaba
10	2	baababaaba
11	5	babaaba

# Suffix Rank and Selection

## Suffix Rank and Selection Queries

**Universe** strings over  $\Sigma$  (denoted  $\Sigma^*$ )

**Input** a string  $T$  of length  $n$

**Queries** rank and select on the set of suffixes of  $T$  ( $Suf(T)$ ).

$T$	a	b	a	a	b	a	b	a	a	b	a
$i$	1	2	3	4	5	6	7	8	9	10	11
ISA	6	10	3	7	11	5	9	2	4	8	1

$$\text{select}_T(5) = T[\text{SA}[5]..] = T[6..] = \text{abaaba}$$

$$\text{rank}_T(T[4..]) = \text{ISA}[4] = 7$$

$$\text{rank}_T(\text{aabb}) = 3$$

$i$		SA
1	11	a
2	8	aaba
3	3	aababaaba
4	9	aba
5	6	abaaba
6	1	abaababaaba
7	4	ababaaba
8	10	ba
9	7	baaba
10	2	baababaaba
11	5	babaaba

## Substring Suffix Rank and Selection Queries

**Universe** strings over  $\Sigma$  (denoted  $\Sigma^*$ )

**Input** a string  $T$  of length  $n$

**Queries** rank and select on suffixes of subwords  $T[\ell..r]$  of  $T$ .

# Substring Suffix Rank and Selection

## Substring Suffix Rank and Selection Queries

**Universe** strings over  $\Sigma$  (denoted  $\Sigma^*$ )

**Input** a string  $T$  of length  $n$

**Queries** rank and select on suffixes of subwords  $T[l..r]$  of  $T$ .

$T$ : 

a	b	a	a	b	a	b	a	a	b	a	a	b	a	b	a	a	b
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

  
118

# Substring Suffix Rank and Selection

## Substring Suffix Rank and Selection Queries

**Universe** strings over  $\Sigma$  (denoted  $\Sigma^*$ )

**Input** a string  $T$  of length  $n$

**Queries** rank and select on suffixes of subwords  $T[l..r]$  of  $T$ .

$T$ : 

a	b	a	a	b	a	b	a	a	b	a	a	b	a	b	a	a	b
1				5									14				18

$select_{T[5..14]}(2) =$

# Substring Suffix Rank and Selection

## Substring Suffix Rank and Selection Queries

**Universe** strings over  $\Sigma$  (denoted  $\Sigma^*$ )

**Input** a string  $T$  of length  $n$

**Queries** rank and select on suffixes of subwords  $T[l..r]$  of  $T$ .

$T$ : 

a	b	a	a	b	a	b	a	a	b	a	a	b	a	b	a	a	b
1			5						11		14						18

$$\text{select}_{T[5..14]}(2) = T[11..14] = \text{aaba}$$



# Substring Suffix Rank and Selection

## Substring Suffix Rank and Selection Queries

**Universe** strings over  $\Sigma$  (denoted  $\Sigma^*$ )

**Input** a string  $T$  of length  $n$

**Queries** rank and select on suffixes of subwords  $T[\ell..r]$  of  $T$ .

$T$ : 

a	b	a	a	b	a	b	a	a	b	a	b	a	a	b
1							9		12				17	18

$$\text{select}_{T[5..14]}(2) = T[11..14] = \text{aaba}$$

$$\text{rank}_{T[9..17]}(T[12..17]) =$$

# Substring Suffix Rank and Selection

## Substring Suffix Rank and Selection Queries

**Universe** strings over  $\Sigma$  (denoted  $\Sigma^*$ )

**Input** a string  $T$  of length  $n$

**Queries** rank and select on suffixes of subwords  $T[\ell..r]$  of  $T$ .

$T$ : 

a	b	a	a	b	a	b	a	a	b	a	b	a	a	b
1							9	11	12	14	16	17	18	

$$\text{select}_{T[5..14]}(2) = T[11..14] = \text{aaba}$$

$$\text{rank}_{T[9..17]}(T[12..17]) = 6$$

# Substring Suffix Rank and Selection

## Substring Suffix Rank and Selection Queries

**Universe** strings over  $\Sigma$  (denoted  $\Sigma^*$ )

**Input** a string  $T$  of length  $n$

**Queries** rank and select on suffixes of subwords  $T[\ell..r]$  of  $T$ .

$T$ : 

a	b	a	a	b	a	b	a	a	b	a	b	a	a	b
1						9							17	18

$$\text{select}_{T[5..14]}(2) = T[11..14] = \text{aaba}$$

$$\text{rank}_{T[9..17]}(T[12..17]) = 6$$

$$\text{rank}_{T[9..17]}(\text{abaab}) =$$

# Substring Suffix Rank and Selection

## Substring Suffix Rank and Selection Queries

**Universe** strings over  $\Sigma$  (denoted  $\Sigma^*$ )

**Input** a string  $T$  of length  $n$

**Queries** rank and select on suffixes of subwords  $T[\ell..r]$  of  $T$ .

$T$ : 

a	b	a	a	b	a	b	a	a	b	a	b	a	a	b
1							9	11		14		16	17	18

$$\text{select}_{T[5..14]}(2) = T[11..14] = \text{aaba}$$

$$\text{rank}_{T[9..17]}(T[12..17]) = 6$$

$$\text{rank}_{T[9..17]}(\text{abaab}) = 4$$

# Substring Suffix Rank and Selection

## Substring Suffix Rank and Selection Queries

**Universe** strings over  $\Sigma$  (denoted  $\Sigma^*$ )

**Input** a string  $T$  of length  $n$

**Queries** rank and select on suffixes of subwords  $T[l..r]$  of  $T$ .

$T$ : 

a	b	a	a	b	a	b	a	a	b	a	a	b	a	a	b
1				5				9	11			14	16	17	18

$$\text{select}_{T[5..14]}(2) = T[11..14] = \text{aaba}$$

$$\text{rank}_{T[9..17]}(T[12..17]) = 6$$

$$\text{rank}_{T[9..17]}(\text{abaab}) = 4$$

$$\text{rank}_{T[9..17]}(T[1..5]) = 4$$

# Substring Suffix Rank and Selection

## Substring Suffix Rank and Selection Queries

**Universe** strings over  $\Sigma$  (denoted  $\Sigma^*$ )

**Input** a string  $T$  of length  $n$

**Queries** rank and select on suffixes of subwords  $T[\ell..r]$  of  $T$ .

$T$ : 

a	b	a	a	b	a	b	a	a	b	a	a	b	a	b	a	a	b
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

  
118

$$\text{select}_{T[5..14]}(2) = T[11..14] = \text{aaba}$$

$$\text{rank}_{T[9..17]}(T[12..17]) = 6$$

$$\text{rank}_{T[9..17]}(\text{abaab}) = 4$$

$$\text{rank}_{T[9..17]}(T[1..5]) = 4$$

# Subword Suffix Rank and Selection: Results

Problem	Space	Query	Construction	Reference
Maximum	$\mathcal{O}(n)$	$\mathcal{O}(1)$	$\mathcal{O}(n)$	B.G.K.S. CPM 2014
Minimum	$\mathcal{O}(n)$	$\mathcal{O}(1)$ $\mathcal{O}(\log n)$	$\mathcal{O}(n \log n)$ $\mathcal{O}(n)$	B.G.K.S. CPM 2014

# Subword Suffix Rank and Selection: Results

Problem	Space	Query	Construction	Reference
Maximum	$\mathcal{O}(n)$	$\mathcal{O}(1)$	$\mathcal{O}(n)$	B.G.K.S. CPM 2014
Minimum	$\mathcal{O}(n)$	$\mathcal{O}(1)$ $\mathcal{O}(\log n)$	$\mathcal{O}(n \log n)$ $\mathcal{O}(n)$	B.G.K.S. CPM 2014
Rank & Select	$\mathcal{O}(n)$	$\mathcal{O}(\log n)$	$\mathcal{O}(n\sqrt{\log n})^*$	this work
(negative)	$n \log^{\mathcal{O}(1)} n$	$\Omega\left(\frac{\log n}{\log \log n}\right)$	—	



# Subword Suffix Rank and Selection: Results

Problem	Space	Query	Construction	Reference
Maximum	$\mathcal{O}(n)$	$\mathcal{O}(1)$	$\mathcal{O}(n)$	B.G.K.S. CPM 2014
Minimum	$\mathcal{O}(n)$	$\mathcal{O}(1)$ $\mathcal{O}(\log n)$	$\mathcal{O}(n \log n)$ $\mathcal{O}(n)$	B.G.K.S. CPM 2014
Rank & Select	$\mathcal{O}(n)$	$\mathcal{O}(\log n)$	$\mathcal{O}(n\sqrt{\log n})^*$	this work
(negative)	$n \log^{\mathcal{O}(1)} n$	$\Omega\left(\frac{\log n}{\log \log n}\right)$	—	

- \* Wavelet suffix trees: randomized construction:  $\mathcal{O}(n\sqrt{\log n})$  deterministic +  $\mathcal{O}(n)$  expected.

# Applications to Substring Compression

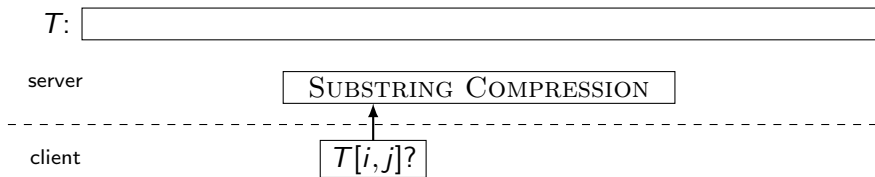
$T$ :

server

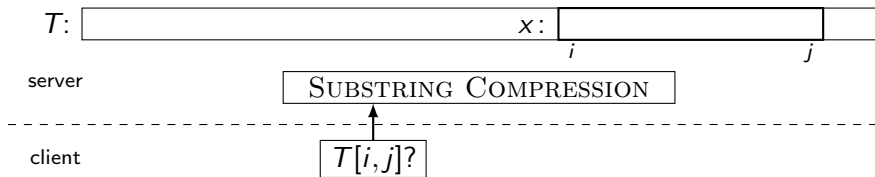
SUBSTRING COMPRESSION

-----  
client

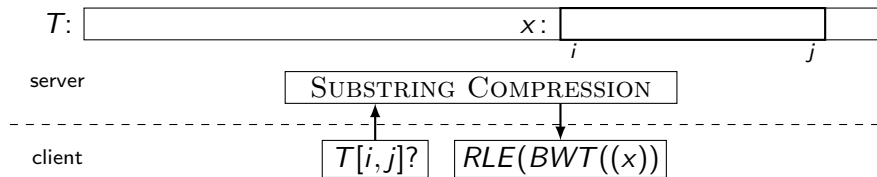
# Applications to Substring Compression



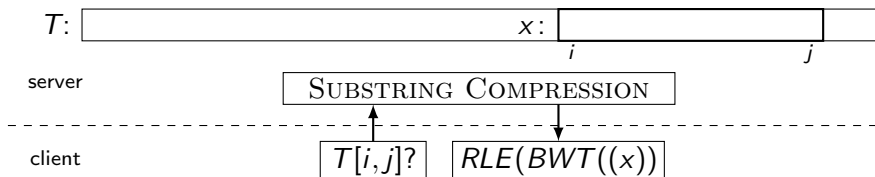
# Applications to Substring Compression



# Applications to Substring Compression



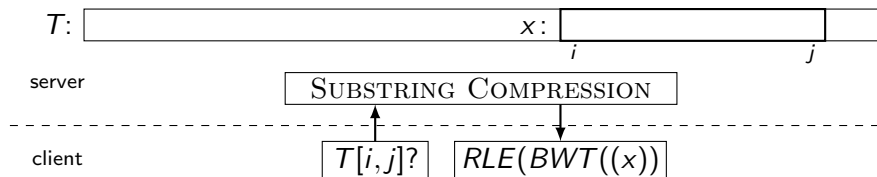
# Applications to Substring Compression



Burrows-Wheeler transform:

- often makes data easier to compress with simple methods:
  - run-length encoding.
- $RLE(BWT(\text{banana}\$)) = RLE(\text{annb}\$aa) = a^1 n^2 b^1 \$^1 a^2$ .

# Applications to Substring Compression



Burrows-Wheeler transform:

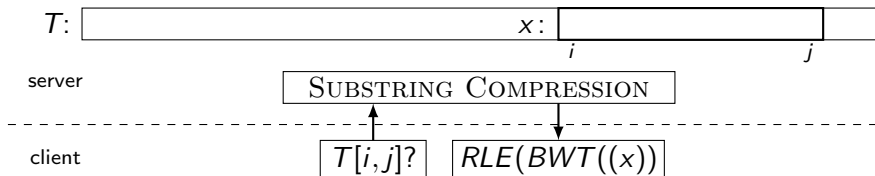
- often makes data easier to compress with simple methods:
  - run-length encoding.
- $RLE(BWT(\text{banana}\$)) = RLE(\text{annb}\$aa) = a^1 n^2 b^1 \$^1 a^2$ .

Substring compression: [Cormode & Muthukrishnan; SODA 2005]

LZ77  $\mathcal{O}(s \log^\epsilon n)$  [Keller et al.; Theor. Comp. Sci., 2014]

BWT+RLE  $\mathcal{O}(s \log n)$  [this work]

# Applications to Substring Compression



Burrows-Wheeler transform:

- often makes data easier to compress with simple methods:
  - run-length encoding.
- $RLE(BWT(\text{banana}\$)) = RLE(\text{annb}\$aa) = a^1 n^2 b^1 \$^1 a^2$ .

Substring compression: [Cormode & Muthukrishnan; SODA 2005]

LZ77  $\mathcal{O}(s \log^\epsilon n)$  [Keller et al.; Theor. Comp. Sci., 2014]

BWT+RLE  $\mathcal{O}(s \log n)$  [this work]

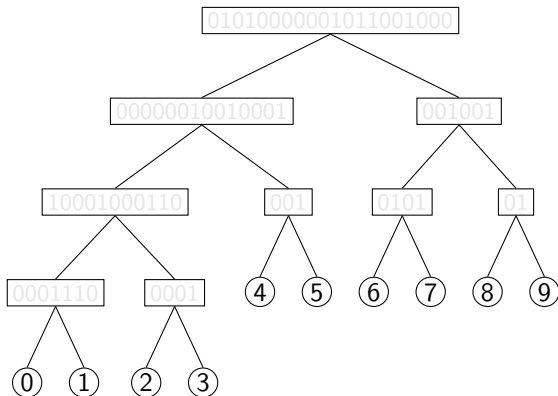
## Acknowledgement

Thanks to Djamel Belazzougui, Travis Gagie, and Simon Puglisi (University of Helsinki) for suggesting the of study BWT queries.



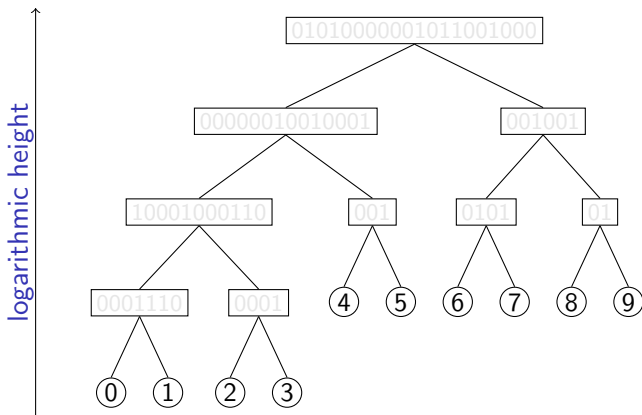
# Wavelet Trees [Grossi, Gupta, Vitter; SODA 2003]

Wavelet tree of  $S = \begin{matrix} 2 & 6 & 0 & 7 & 0 & 0 & 2 & 1 & 4 & 1 & 8 & 1 & 6 & 7 & 4 & 2 & 9 & 3 & 0 & 5 \\ 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 & 14 & 15 & 16 & 17 & 18 & 19 & 20 \end{matrix}$



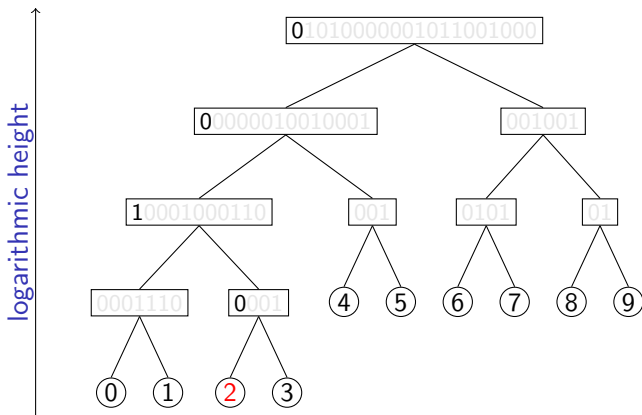
# Wavelet Trees [Grossi, Gupta, Vitter; SODA 2003]

Wavelet tree of  $S = \begin{matrix} 2 & 6 & 0 & 7 & 0 & 0 & 2 & 1 & 4 & 1 & 8 & 1 & 6 & 7 & 4 & 2 & 9 & 3 & 0 & 5 \\ 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 & 14 & 15 & 16 & 17 & 18 & 19 & 20 \end{matrix}$



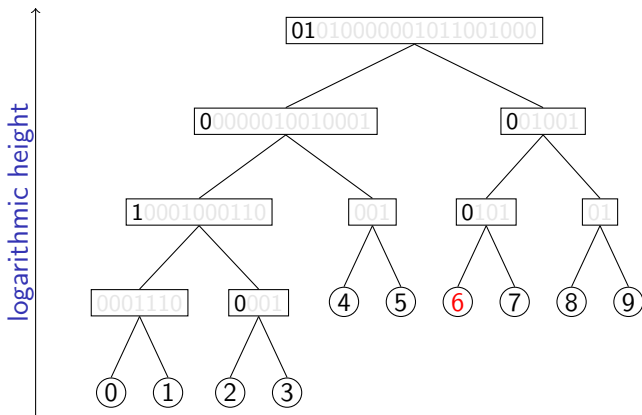
# Wavelet Trees [Grossi, Gupta, Vitter; SODA 2003]

Wavelet tree of  $S = \overset{1}{2} \overset{2}{6} \overset{3}{0} \overset{4}{7} \overset{5}{0} \overset{6}{0} \overset{7}{2} \overset{8}{1} \overset{9}{4} \overset{10}{1} \overset{11}{8} \overset{12}{1} \overset{13}{6} \overset{14}{7} \overset{15}{4} \overset{16}{2} \overset{17}{9} \overset{18}{3} \overset{19}{0} \overset{20}{5}$



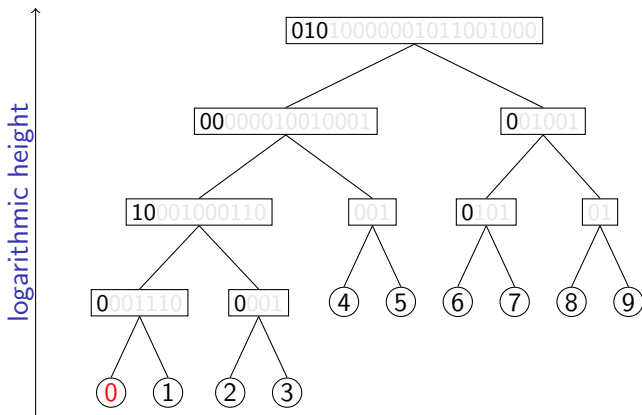
# Wavelet Trees [Grossi, Gupta, Vitter; SODA 2003]

Wavelet tree of  $S = 2 \ 6 \ 0 \ 7 \ 0 \ 0 \ 2 \ 1 \ 4 \ 1 \ 8 \ 1 \ 6 \ 7 \ 4 \ 2 \ 9 \ 3 \ 0 \ 5$   
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20



# Wavelet Trees [Grossi, Gupta, Vitter; SODA 2003]

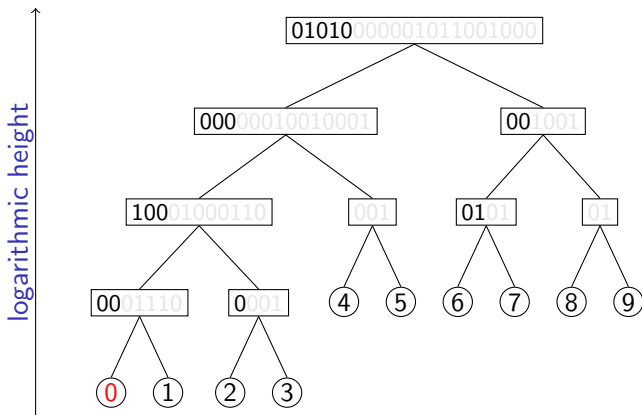
Wavelet tree of  $S = 2\ 6\ 0\ 7\ 0\ 0\ 2\ 1\ 4\ 1\ 8\ 1\ 6\ 7\ 4\ 2\ 9\ 3\ 0\ 5$   
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20





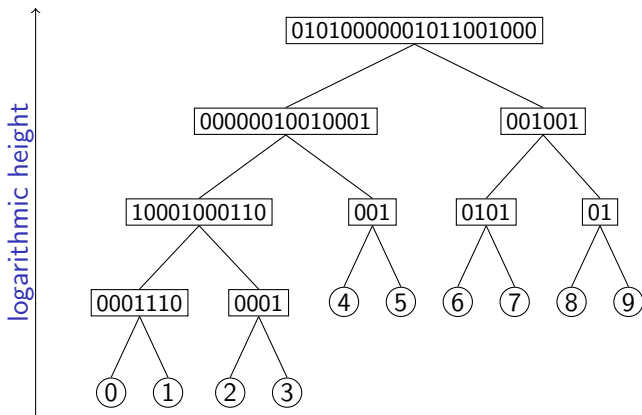
# Wavelet Trees [Grossi, Gupta, Vitter; SODA 2003]

Wavelet tree of  $S = 2\ 6\ 0\ 7\ 0\ 2\ 1\ 4\ 1\ 8\ 1\ 6\ 7\ 4\ 2\ 9\ 3\ 0\ 5$   
<sub>1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20</sub>



# Wavelet Trees [Grossi, Gupta, Vitter; SODA 2003]

Wavelet tree of  $S = \begin{matrix} 2 & 6 & 0 & 7 & 0 & 0 & 2 & 1 & 4 & 1 & 8 & 1 & 6 & 7 & 4 & 2 & 9 & 3 & 0 & 5 \\ 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 & 14 & 15 & 16 & 17 & 18 & 19 & 20 \end{matrix}$



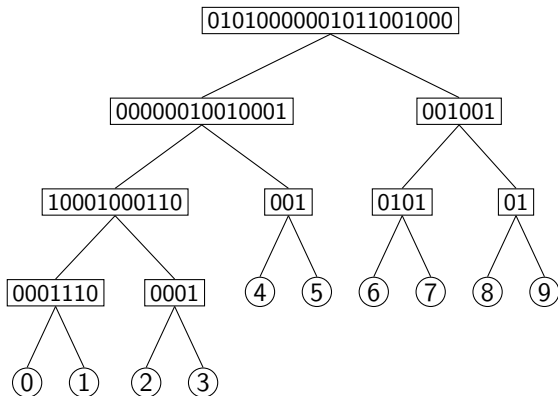
Size  $\mathcal{O}(n)$

Construction  $\mathcal{O}(n\sqrt{\log n})$



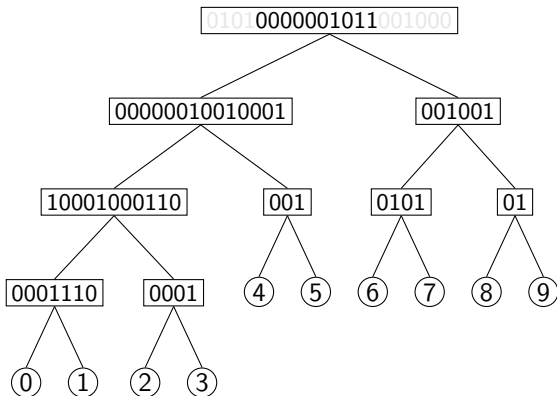
# Range Queries with Wavelet Trees

Wavelet tree of  $S = \underset{1}{2} \underset{2}{6} \underset{3}{0} \underset{4}{7} \underset{5}{0} \underset{6}{0} \underset{7}{2} \underset{8}{1} \underset{9}{4} \underset{10}{1} \underset{11}{8} \underset{12}{1} \underset{13}{6} \underset{14}{7} \underset{15}{4} \underset{16}{2} \underset{17}{9} \underset{18}{3} \underset{19}{0} \underset{20}{5}$



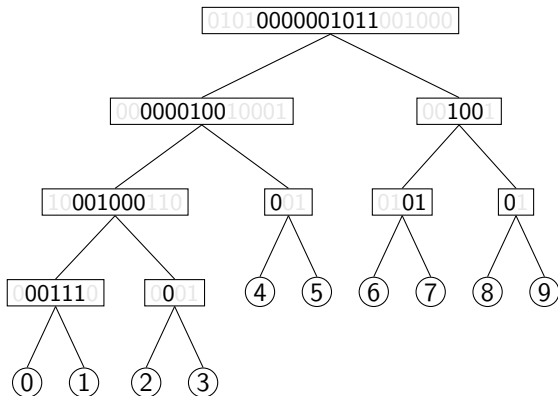
# Range Queries with Wavelet Trees

Wavelet tree of  $S = 2\ 6\ 0\ 7\ 0\ 0\ 2\ 1\ 4\ 1\ 8\ 1\ 6\ 7\ 4\ 2\ 9\ 3\ 0\ 5$   
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20



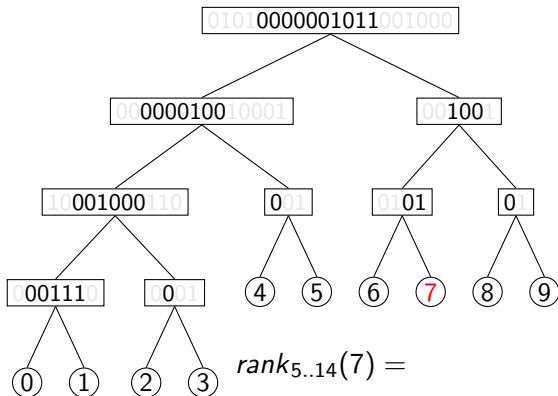
# Range Queries with Wavelet Trees

Wavelet tree of  $S = 2\ 6\ 0\ 7\ 0\ 0\ 2\ 1\ 4\ 1\ 8\ 1\ 6\ 7\ 4\ 2\ 9\ 3\ 0\ 5$   
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20



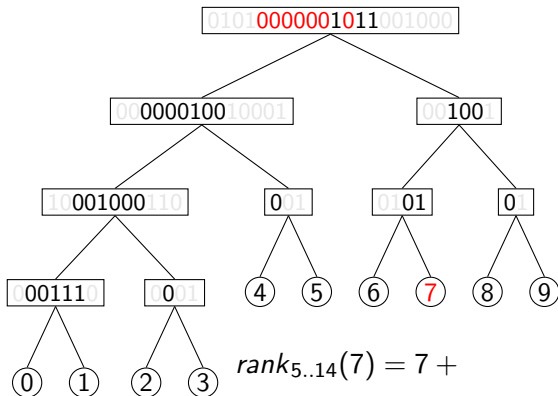
# Range Queries with Wavelet Trees

Wavelet tree of  $S = 2\ 6\ 0\ 7\ 0\ 0\ 2\ 1\ 4\ 1\ 8\ 1\ 6\ 7\ 4\ 2\ 9\ 3\ 0\ 5$   
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20



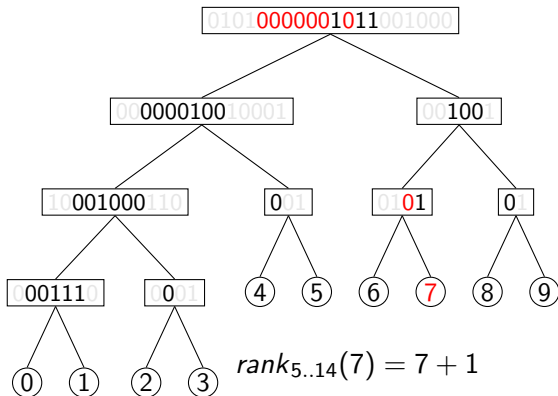
# Range Queries with Wavelet Trees

Wavelet tree of  $S = 2\ 6\ 0\ 7\ 0\ 0\ 2\ 1\ 4\ 1\ 8\ 1\ 6\ 7\ 4\ 2\ 9\ 3\ 0\ 5$   
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20



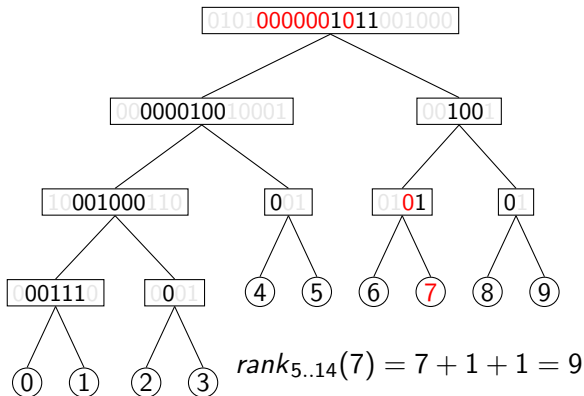
# Range Queries with Wavelet Trees

Wavelet tree of  $S = 2\ 6\ 0\ 7\ \mathbf{0\ 0\ 2\ 1\ 4\ 1\ 8\ 1\ 6\ 7}\ 4\ 2\ 9\ 3\ 0\ 5$   
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20



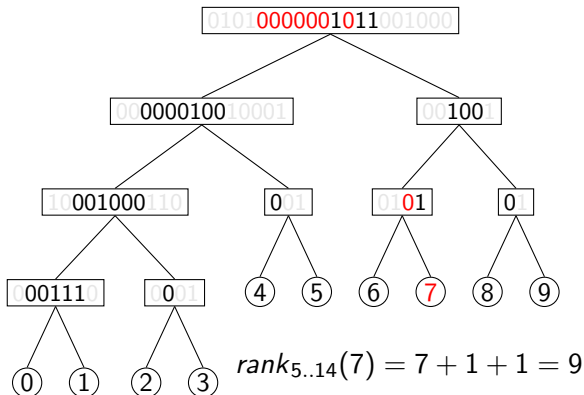
# Range Queries with Wavelet Trees

Wavelet tree of  $S = 2\ 6\ 0\ 7\ 0\ 0\ 2\ 1\ 4\ 1\ 8\ 1\ 6\ 7\ 4\ 2\ 9\ 3\ 0\ 5$   
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20



# Range Queries with Wavelet Trees

Wavelet tree of  $S = 2\ 6\ 0\ 7\ 0\ 0\ 2\ 1\ 4\ 1\ 8\ 1\ 6\ 7\ 4\ 2\ 9\ 3\ 0\ 5$   
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20



Rank  $\mathcal{O}(\log n)$  ( $\mathcal{O}(\frac{\log n}{\log \log n})$ ): higher arity)

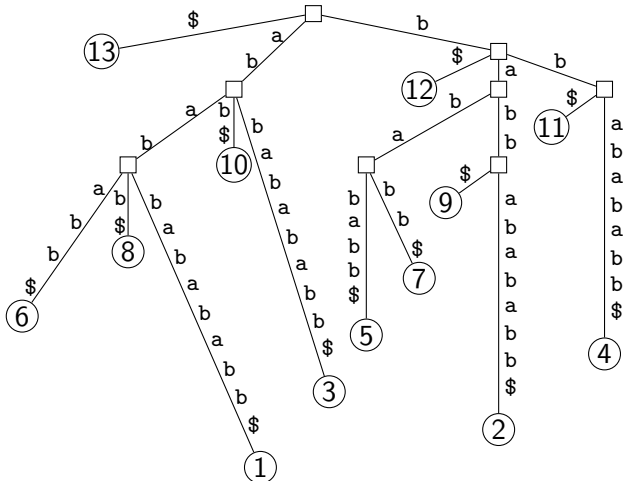
Select  $\mathcal{O}(\log n)$  ( $\mathcal{O}(\frac{\log n}{\log \log n})$ ): higher arity + extra tools)





# Suffix Trees

Suffix tree of  $T = a b a b b a b a b a b b$   
1 2 3 4 5 6 7 8 9 10 11 12

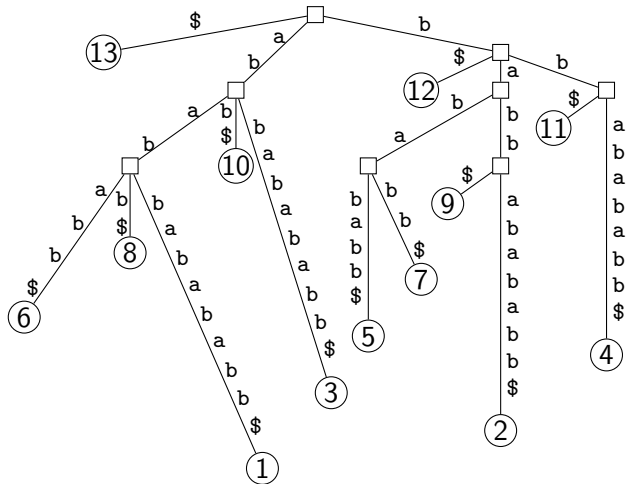


+ subwords of  $T$  and suffixes  
of  $T\$$  in lexicographic order

# Suffix Trees

Suffix tree of  $T = a\ b\ a\ b\ b\ a\ b\ a\ b\ a\ b\ b\ b$

1 2 3 4 5 6 7 8 9 10 11 12

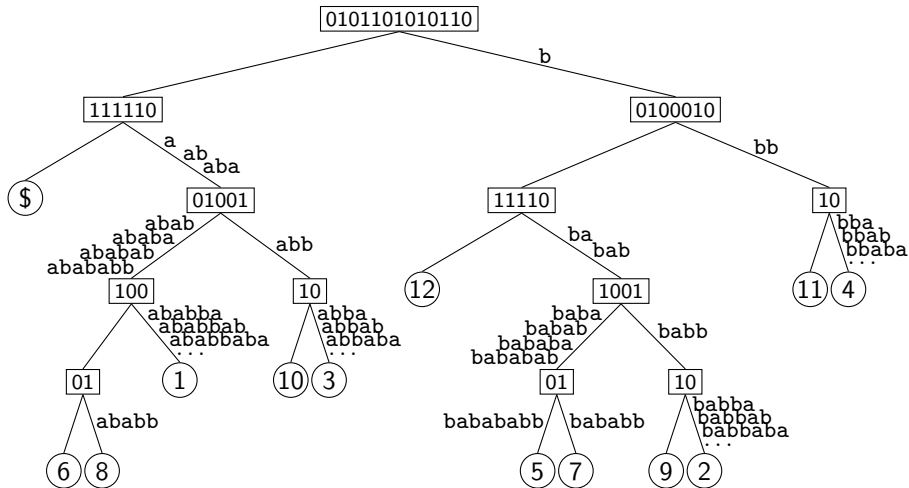


+ subwords of  $T$  and suffixes  
of  $T\$$  in lexicographic order

- high depth  
- non-uniform arity

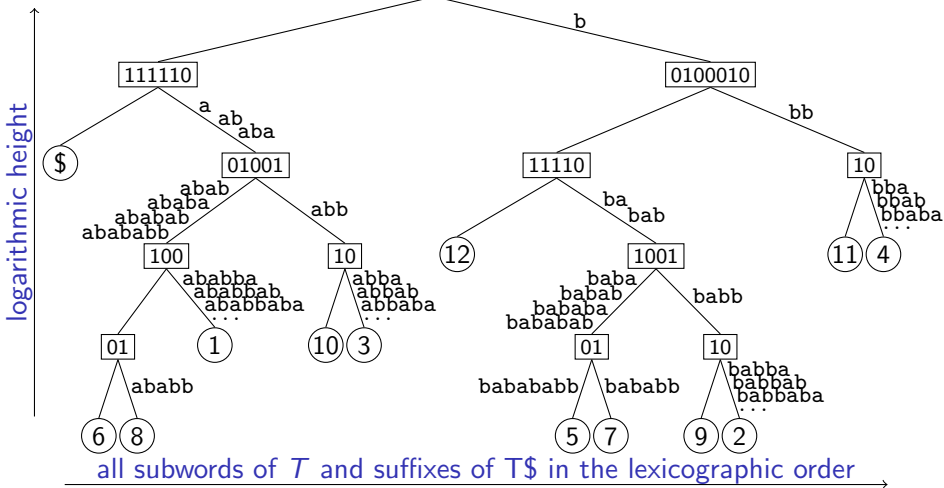
# Wavelet Suffix Trees

Wavelet suffix tree of  $T = \text{a b a b b a b a b a b b}$   
1 2 3 4 5 6 7 8 9 10 11 12



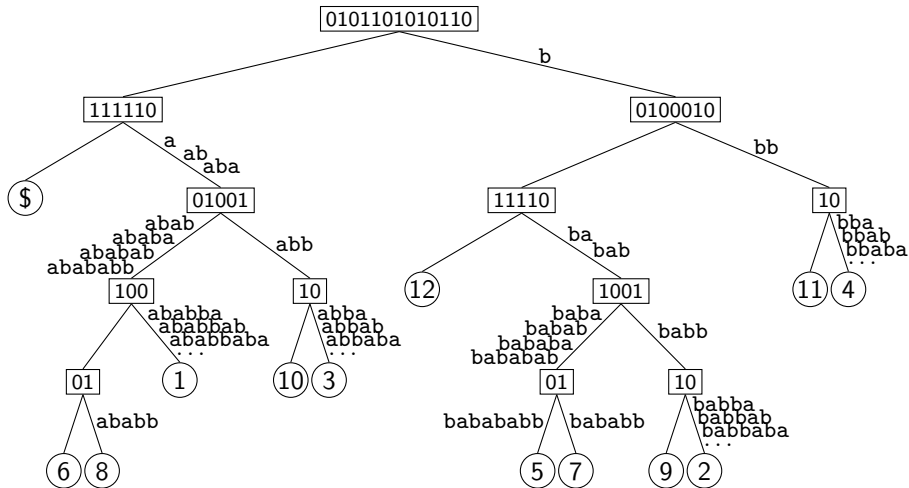
# Wavelet Suffix Trees

Wavelet suffix tree of  $T = \underset{1}{a} \underset{2}{b} \underset{3}{a} \underset{4}{b} \underset{5}{b} \underset{6}{a} \underset{7}{b} \underset{8}{a} \underset{9}{b} \underset{10}{a} \underset{11}{b} \underset{12}{b}$



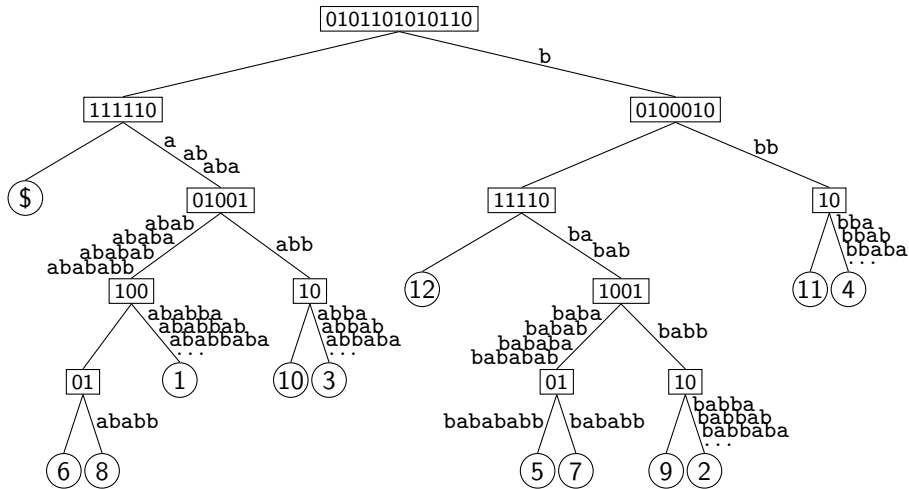
# Subword Suffix Queries with Wavelet Suffix Trees

Wavelet suffix tree of  $T = \underset{1}{a} \underset{2}{b} \underset{3}{a} \underset{4}{b} \underset{5}{b} \underset{6}{a} \underset{7}{b} \underset{8}{a} \underset{9}{b} \underset{10}{a} \underset{11}{b} \underset{12}{b}$



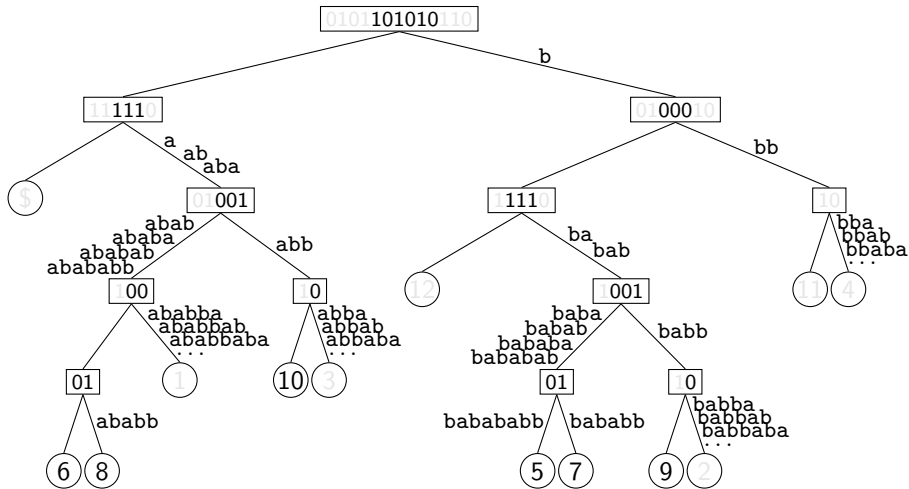
# Subword Suffix Queries with Wavelet Suffix Trees

Wavelet suffix tree of  $T = a b a b \mathbf{b a b a b a} b b$   
1 2 3 4 5 6 7 8 9 10 11 12



# Subword Suffix Queries with Wavelet Suffix Trees

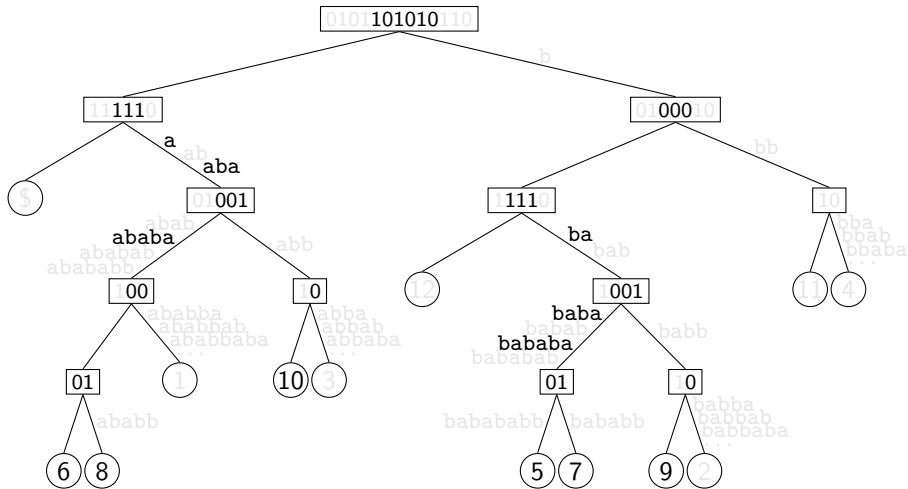
Wavelet suffix tree of  $T = \text{a b a b } \text{b a b a b a b b}$   
1 2 3 4 5 6 7 8 9 10 11 12





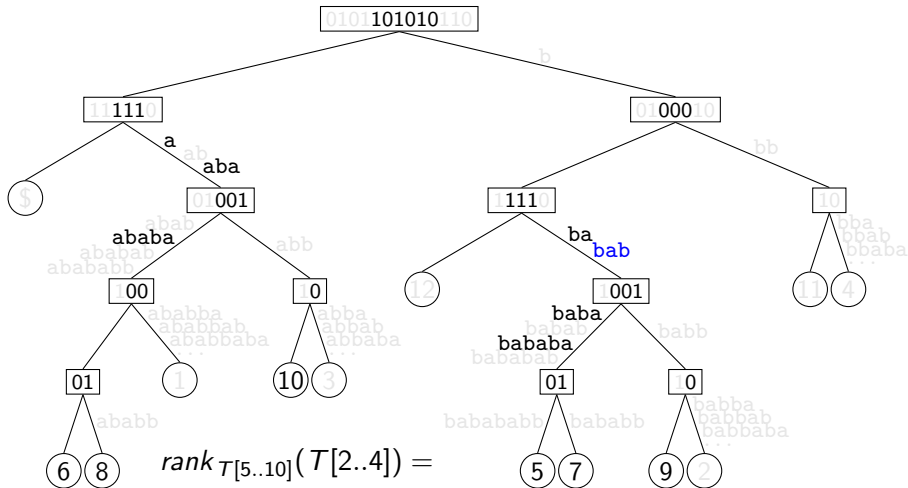
# Subword Suffix Queries with Wavelet Suffix Trees

Wavelet suffix tree of  $T = \text{a b a b b a b a b a b b}$   
1 2 3 4 5 6 7 8 9 10 11 12



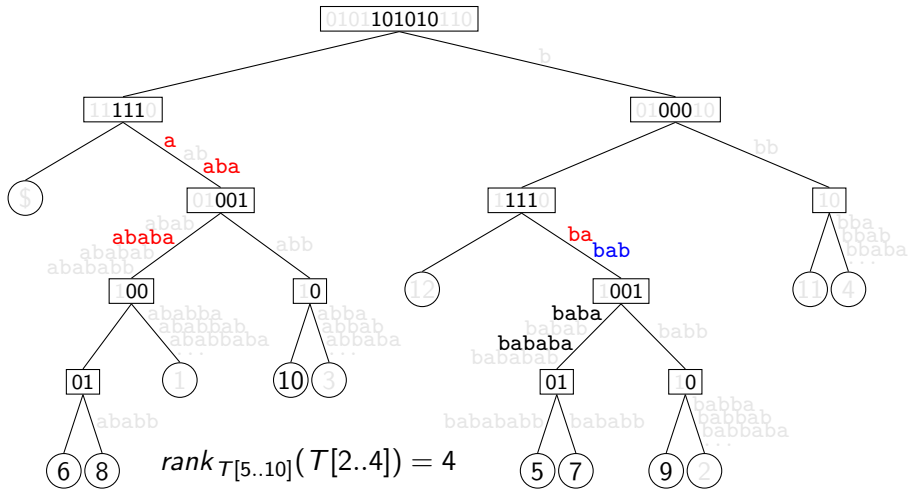
# Subword Suffix Queries with Wavelet Suffix Trees

Wavelet suffix tree of  $T = \text{a b a b } \mathbf{b a b a b a} \text{ b b}$   
1 2 3 4 5 6 7 8 9 10 11 12



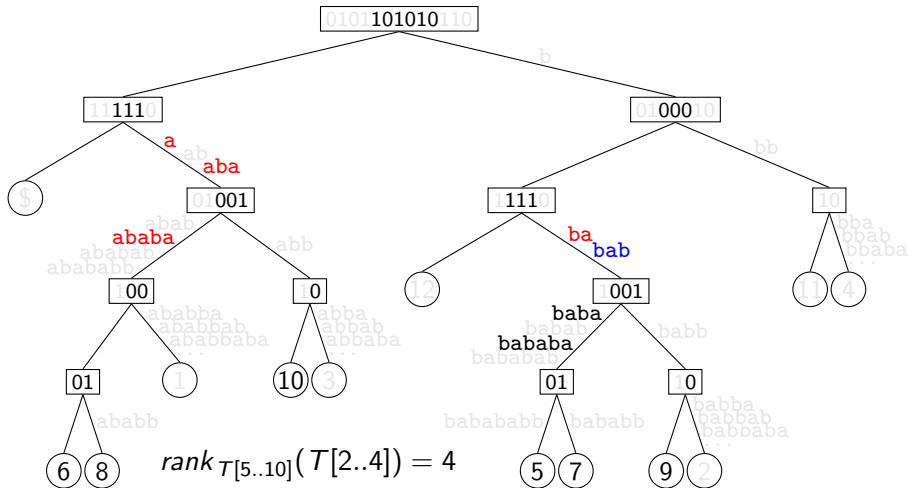
# Subword Suffix Queries with Wavelet Suffix Trees

Wavelet suffix tree of  $T = \underset{1}{a} \underset{2}{b} \underset{3}{a} \underset{4}{b} \underset{5}{b} \underset{6}{a} \underset{7}{b} \underset{8}{a} \underset{9}{b} \underset{10}{a} \underset{11}{b} \underset{12}{b}$



# Subword Suffix Queries with Wavelet Suffix Trees

Wavelet suffix tree of  $T = \underset{1}{a} \underset{2}{b} \underset{3}{a} \underset{4}{b} \underset{5}{b} \underset{6}{a} \underset{7}{b} \underset{8}{a} \underset{9}{b} \underset{10}{a} \underset{11}{b} \underset{12}{b}$



**Operations** count suffixes in subtrees, generate suffixes on an edge.

**Tools** INTERNAL PATTERN MATCHING and bitmasks.

Our results:

- $\mathcal{O}(n\sqrt{\log n})$  construction of wavelet trees,
- simultaneously obtained state-of-the-art construction and query time range selection,
- developed wavelet suffix trees to answer substring suffix rank & selection,
- applied wavelet suffix trees for substring compression with Burrows-Wheeler transform and run-length encoding.

# Conclusions & Open Problems

Our results:

- $\mathcal{O}(n\sqrt{\log n})$  construction of wavelet trees,
- simultaneously obtained state-of-the-art construction and query time range selection,
- developed wavelet suffix trees to answer substring suffix rank & selection,
- applied wavelet suffix trees for substring compression with Burrows-Wheeler transform and run-length encoding.

Open problems:

- Can the  $\mathcal{O}(n\sqrt{\log n})$  construction time be improved?
  - Would affect counting inversions.

# Conclusions & Open Problems

Our results:

- $\mathcal{O}(n\sqrt{\log n})$  construction of wavelet trees,
- simultaneously obtained state-of-the-art construction and query time range selection,
- developed wavelet suffix trees to answer substring suffix rank & selection,
- applied wavelet suffix trees for substring compression with Burrows-Wheeler transform and run-length encoding.

Open problems:

- Can the  $\mathcal{O}(n\sqrt{\log n})$  construction time be improved?
  - Would affect counting inversions.
- Are substring suffix queries inherently harder than analogous range queries?
  - Currently  $\mathcal{O}(\log n)$  vs  $\mathcal{O}\left(\frac{\log n}{\log \log n}\right)$ .

Thank you for your attention!