

PARALLEL  
WAVELET TREE CONSTRUCTION  
SIMPLE, FAST AND LIGHTWEIGHT

Johannes Fischer    *Florian Kurpicz*

# RANK AND SELECT

$\text{rank}_\alpha(i)$  # of  $\alpha$  before position  $i$

$\text{select}_\alpha(j)$  position of  $j$ -th  $\alpha$

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

# RANK AND SELECT

$\text{rank}_\alpha(i)$  # of  $\alpha$  before position  $i$

$\text{select}_\alpha(j)$  position of  $j$ -th  $\alpha$

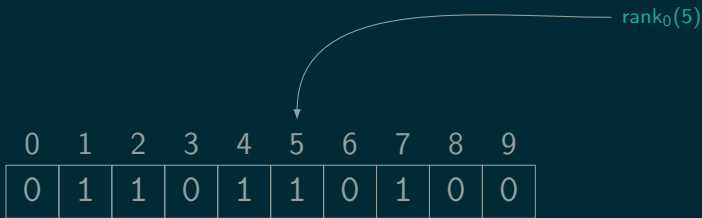
$\text{rank}_0(5)$

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

# RANK AND SELECT

$\text{rank}_\alpha(i)$  # of  $\alpha$  before position  $i$

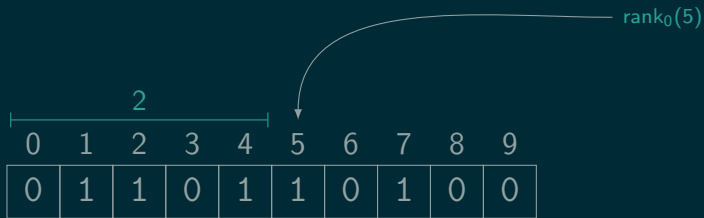
$\text{select}_\alpha(j)$  position of  $j$ -th  $\alpha$



# RANK AND SELECT

$\text{rank}_\alpha(i)$  # of  $\alpha$  before position  $i$

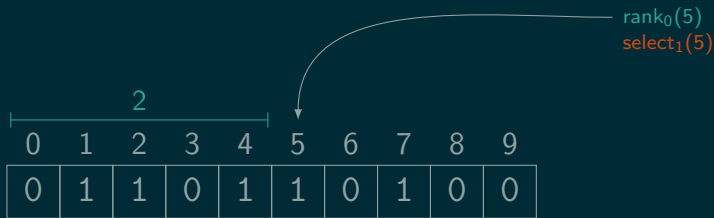
$\text{select}_\alpha(j)$  position of  $j$ -th  $\alpha$



# RANK AND SELECT

$\text{rank}_\alpha(i)$  # of  $\alpha$  before position  $i$

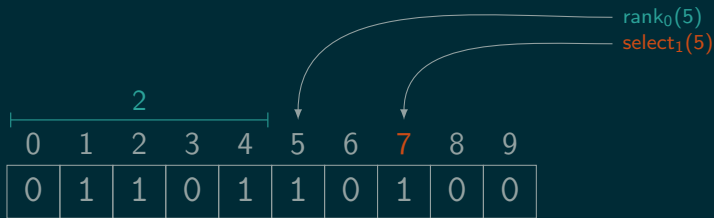
$\text{select}_\alpha(j)$  position of  $j$ -th  $\alpha$



# RANK AND SELECT

$\text{rank}_\alpha(i)$  # of  $\alpha$  before position  $i$

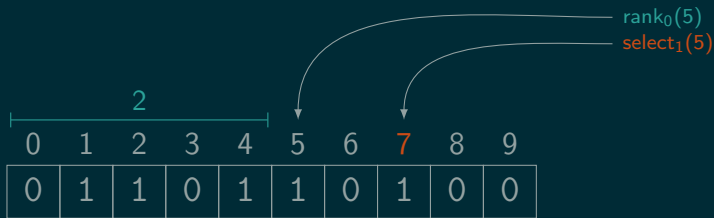
$\text{select}_\alpha(j)$  position of  $j$ -th  $\alpha$



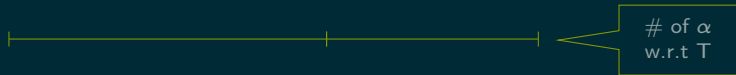
# RANK AND SELECT

$\text{rank}_\alpha(i)$  # of  $\alpha$  before position  $i$

$\text{select}_\alpha(j)$  position of  $j$ -th  $\alpha$



super-block

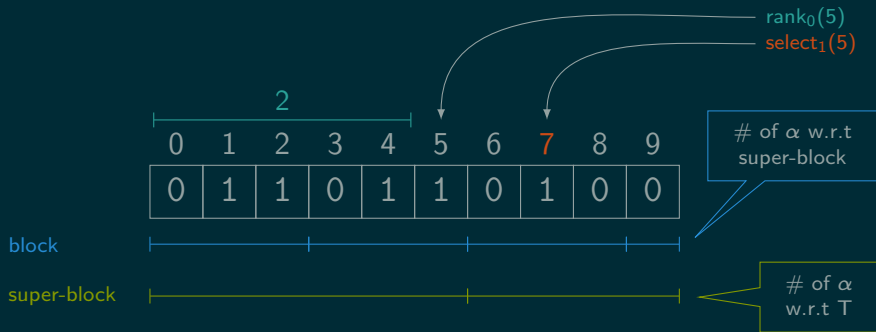




# RANK AND SELECT

$\text{rank}_\alpha(i)$  # of  $\alpha$  before position  $i$

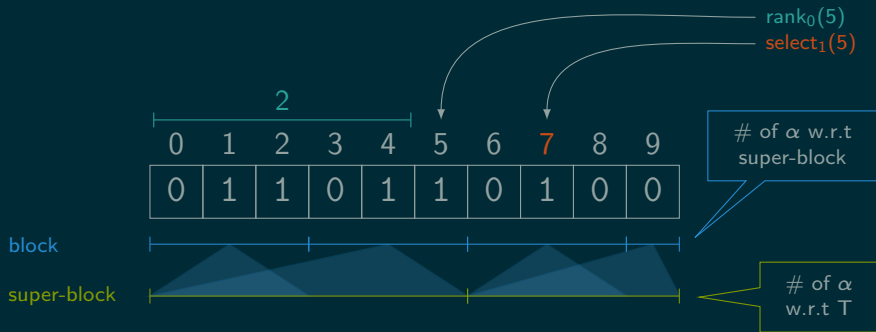
$\text{select}_\alpha(j)$  position of  $j$ -th  $\alpha$



# RANK AND SELECT

$\text{rank}_\alpha(i)$  # of  $\alpha$  before position  $i$

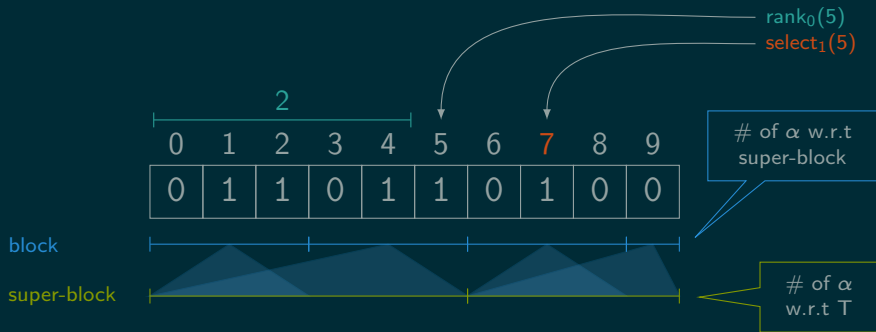
$\text{select}_\alpha(j)$  position of  $j$ -th  $\alpha$



# RANK AND SELECT

$\text{rank}_\alpha(i)$  # of  $\alpha$  before position  $i$

$\text{select}_\alpha(j)$  position of  $j$ -th  $\alpha$



What if  $|\Sigma| > 2$

$$\Sigma = \{0, 1, 2, 3, 4, 5, 6, 7\}$$

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

$$\Sigma = \{0, 1, 2, 3, 4, 5, 6, 7\}$$

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

MSB	0	0	1	1	0	1	1	0	1	0
	0	0	1	1	0	0	0	1	1	1
LSB	0	1	0	1	1	1	0	0	0	1

$$\Sigma = \{0, 1, 2, 3, 4, 5, 6, 7\}$$

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

MSB	0	0	1	1	0	1	1	0	1	0
	0	0	1	1	0	0	0	1	1	1
LSB	0	1	0	1	1	1	0	0	0	1

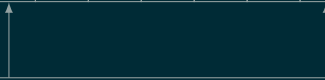
$\text{rank}_6(9)$

$$\Sigma = \{0, 1, 2, 3, 4, 5, 6, 7\}$$

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

MSB	0	0	1	1	0	1	1	0	1	0
	0	0	1	1	0	0	0	1	1	1
LSB	0	1	0	1	1	1	0	0	0	1

$\text{rank}_6(9)$



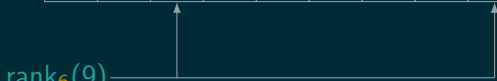
$$\Sigma = \{0, 1, 2, 3, 4, 5, 6, 7\}$$

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

MSB	0	0	1	1	0	1	1	0	1	0
	0	0	1	1	0	0	0	1	1	1
LSB	0	1	0	1	1	1	0	0	0	1

↓ level by level

rank<sub>6</sub>(9)





# WAVELET TREE

POINTER-BASED

→

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0
0	0	1	1	0	0	0	1	1	1
0	1	0	1	1	1	0	0	0	1

[0, 7]

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0

# WAVELET TREE

POINTER-BASED

→

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0
0	0	1	1	0	0	0	1	1	1
0	1	0	1	1	1	0	0	0	1

[0, 7]

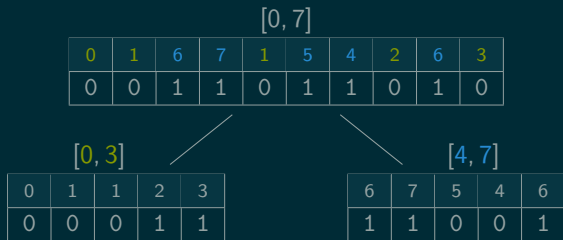
0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0

# WAVELET TREE

POINTER-BASED

→

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0
0	0	1	1	0	0	0	1	1	1
0	1	0	1	1	1	0	0	0	1

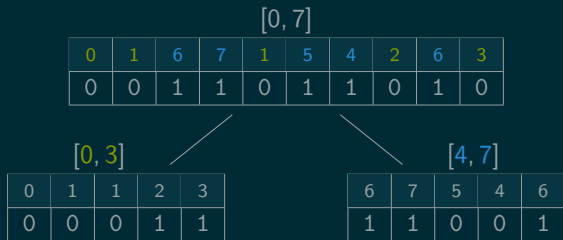


# WAVELET TREE

POINTER-BASED

→

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0
0	0	1	1	0	0	0	1	1	1
0	1	0	1	1	1	0	0	0	1

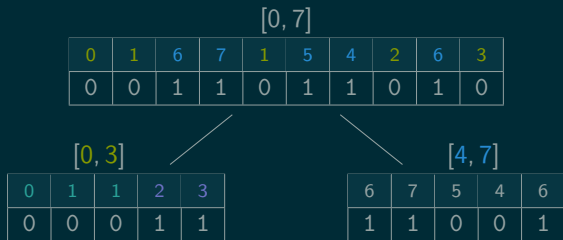


# WAVELET TREE

POINTER-BASED

→

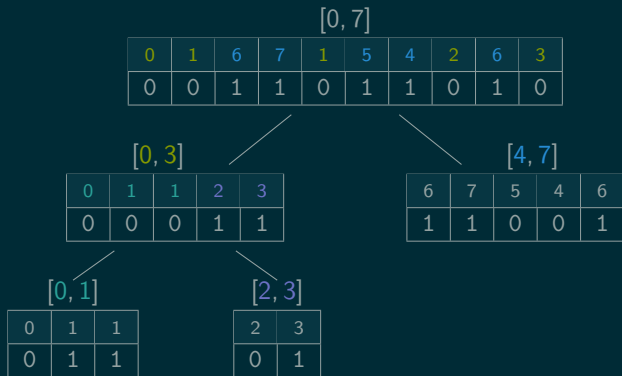
0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0
0	0	1	1	0	0	0	1	1	1
0	1	0	1	1	1	0	0	0	1



# WAVELET TREE

POINTER-BASED

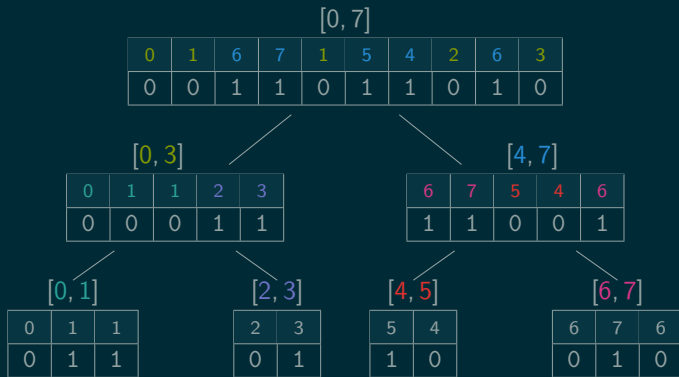
0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0
0	0	1	1	0	0	0	1	1	1
→	0	1	0	1	1	1	0	0	1



# WAVELET TREE

POINTER-BASED

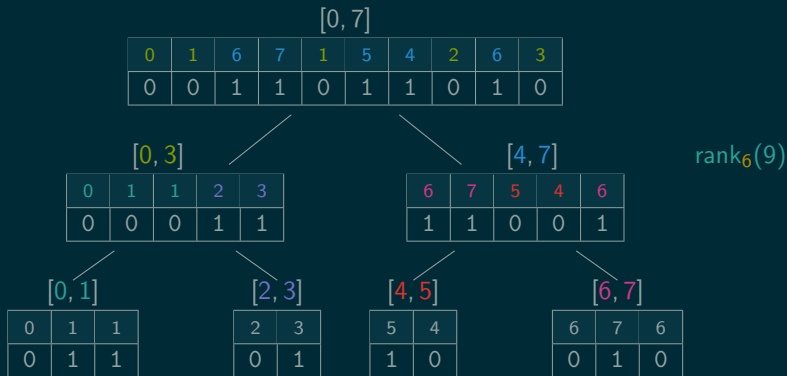
0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0
0	0	1	1	0	0	0	1	1	1
→	0	1	0	1	1	1	0	0	1



# WAVELET TREE

POINTER-BASED

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0
0	0	1	1	0	0	0	1	1	1
0	1	0	1	1	1	0	0	0	1

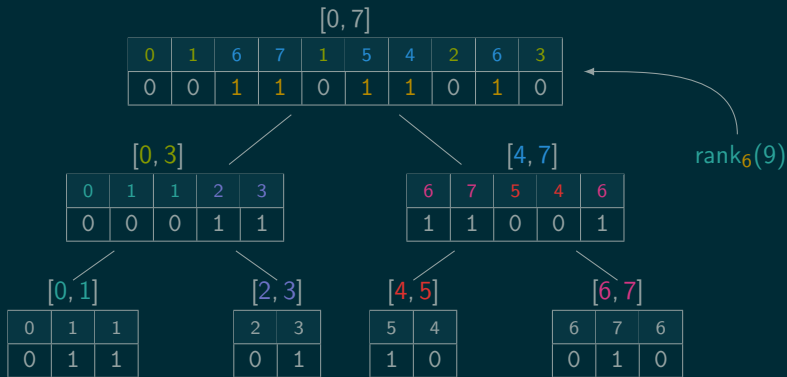




# WAVELET TREE

POINTER-BASED

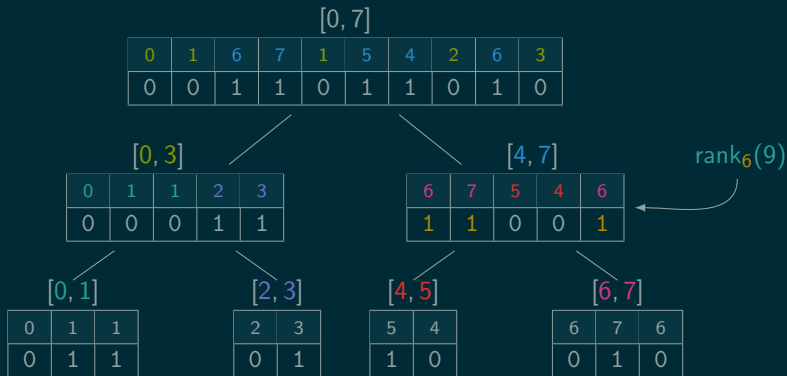
0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0
0	0	1	1	0	0	0	1	1	1
0	1	0	1	1	1	0	0	0	1



# WAVELET TREE

POINTER-BASED

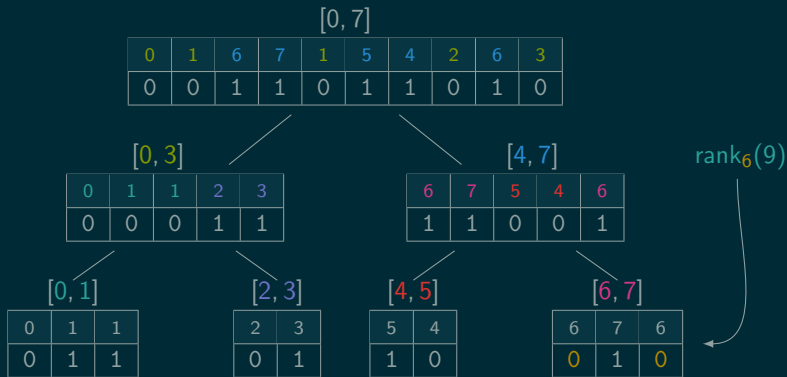
0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0
0	0	1	1	0	0	0	1	1	1
0	1	0	1	1	1	0	0	0	1



# WAVELET TREE

POINTER-BASED

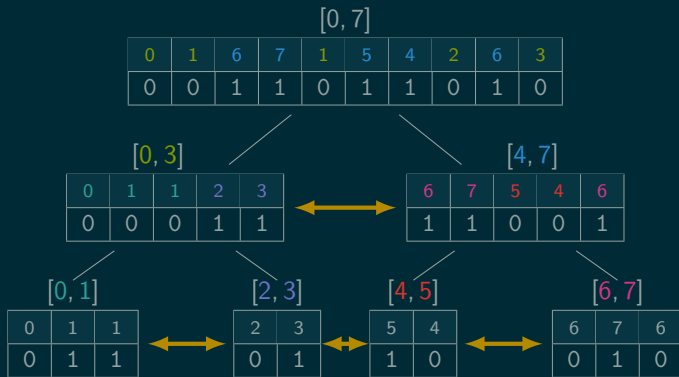
0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0
0	0	1	1	0	0	0	1	1	1
0	1	0	1	1	1	0	0	0	1



# WAVELET TREE

POINTER-BASED

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0
0	0	1	1	0	0	0	1	1	1
0	1	0	1	1	1	0	0	0	1



# WAVELET TREE

LEVEL-WISE

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0
0	0	1	1	0	0	0	1	1	1
0	1	0	1	1	1	0	0	0	1

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0

0	1	1	2	3	6	7	5	4	6
0	0	0	1	1	1	1	0	0	1

0	1	1	2	3	5	4	6	7	6
0	1	1	0	1	1	0	0	1	0

# WAVELET TREE

LEVEL-WISE

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0
0	0	1	1	0	0	0	1	1	1
0	1	0	1	1	1	0	0	0	1

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0

MSB of each character

0	1	1	2	3	6	7	5	4	6
0	0	0	1	1	1	1	0	0	1

0	1	1	2	3	5	4	6	7	6
0	1	1	0	1	1	0	0	1	0

# WAVELET TREE

LEVEL-WISE

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0
0	0	1	1	0	0	0	1	1	1
0	1	0	1	1	1	0	0	0	1

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0

MSB of each character

0	1	1	2	3	6	7	5	4	6
0	0	0	1	1	1	1	0	0	1

2<sup>nd</sup> bit of each character with MSB

0 | 1

0	1	1	2	3	5	4	6	7	6
0	1	1	0	1	1	0	0	1	0

# WAVELET TREE

LEVEL-WISE

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0
0	0	1	1	0	0	0	1	1	1
0	1	0	1	1	1	0	0	0	1

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0

MSB of each character

0	1	1	2	3	6	7	5	4	6
0	0	0	1	1	1	1	0	0	1

2<sup>nd</sup> bit of each character with MSB

0 | 1

0	1	1	2	3	5	4	6	7	6
0	1	1	0	1	1	0	0	1	0

3<sup>rd</sup> bit of each character with MSBs

00 | 01 | 10 | 11



# WAVELET TREE

LEVEL-WISE

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0
0	0	1	1	0	0	0	1	1	1
0	1	0	1	1	1	0	0	0	1

Interval on level  $\ell$  is given by *bit prefix* of length  $\ell$

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0

MSB of each character

0	1	1	2	3	6	7	5	4	6
0	0	0	1	1	1	1	0	0	1

2<sup>nd</sup> bit of each character with MSB

0 | 1

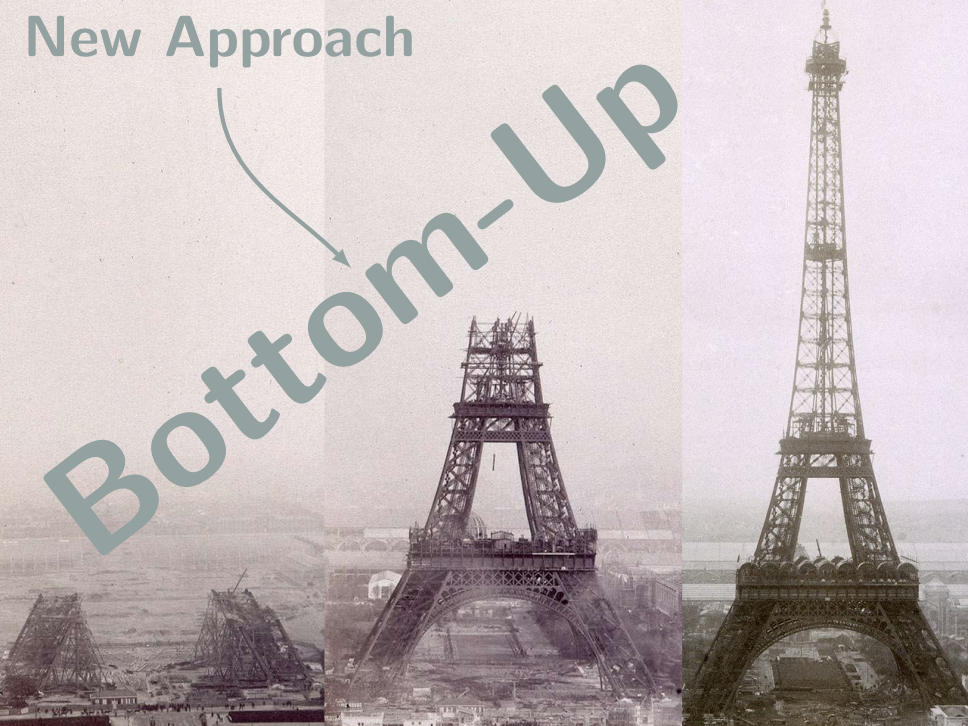
0	1	1	2	3	5	4	6	7	6
0	1	1	0	1	1	0	0	1	0

3<sup>rd</sup> bit of each character with MSBs

00 | 01 | 10 | 11

New Approach

Bottom-Up



# WAVELET TREE

BOTTOM-UP

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0
0	0	1	1	0	0	0	1	1	1
0	1	0	1	1	1	0	0	0	1

Interval on level  $\ell$  is given by *bit prefix* of length  $\ell$

$\ell = 0$

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0

$\ell = 1$

0	1	1	2	3	6	7	5	4	6
0	0	0	1	1	1	1	0	0	1

$\ell = 2$

0	1	1	2	3	5	4	6	7	6
0	1	1	0	1	1	0	0	1	0

# WAVELET TREE

BOTTOM-UP

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0
0	0	1	1	0	0	0	1	1	1
0	1	0	1	1	1	0	0	0	1

Interval on level  $\ell$  is given by *bit prefix* of length  $\ell$

$\ell = 0$

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0

00 $\rightarrow$ [0, 1]	0
01 $\rightarrow$ [2, 3]	0
10 $\rightarrow$ [4, 5]	0
11 $\rightarrow$ [6, 7]	0

$\ell = 1$

0	1	1	2	3	6	7	5	4	6
0	0	0	1	1	1	1	0	0	1

$\ell = 2$

0	1	1	2	3	5	4	6	7	6
0	1	1	0	1	1	0	0	1	0

# WAVELET TREE

BOTTOM-UP

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0
0	0	1	1	0	0	0	1	1	1
0	1	0	1	1	1	0	0	0	1

Interval on level  $\ell$  is given by *bit prefix* of length  $\ell$

$\ell = 0$

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0

$\ell = 1$

0	1	1	2	3	6	7	5	4	6
0	0	0	1	1	1	1	0	0	1

$\ell = 2$

0	1	1	2	3	5	4	6	7	6
0	1	1	0	1	1	0	0	1	0

00 $\rightarrow$ [0, 1]	1
01 $\rightarrow$ [2, 3]	0
10 $\rightarrow$ [4, 5]	0
11 $\rightarrow$ [6, 7]	0

# WAVELET TREE

BOTTOM-UP

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0
0	0	1	1	0	0	0	1	1	1
0	1	0	1	1	1	0	0	0	1

Interval on level  $\ell$  is given by *bit prefix* of length  $\ell$

$\ell = 0$

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0

$\ell = 1$

0	1	1	2	3	6	7	5	4	6
0	0	0	1	1	1	1	0	0	1

$\ell = 2$

0	1	1	2	3	5	4	6	7	6
0	1	1	0	1	1	0	0	1	0

00 $\rightarrow$ [0, 1]	2
01 $\rightarrow$ [2, 3]	0
10 $\rightarrow$ [4, 5]	0
11 $\rightarrow$ [6, 7]	0

# WAVELET TREE

BOTTOM-UP

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0
0	0	1	1	0	0	0	1	1	1
0	1	0	1	1	1	0	0	0	1

Interval on level  $\ell$  is given by *bit prefix* of length  $\ell$

$\ell = 0$

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0

$\ell = 1$

0	1	1	2	3	6	7	5	4	6
0	0	0	1	1	1	1	0	0	1

$\ell = 2$

0	1	1	2	3	5	4	6	7	6
0	1	1	0	1	1	0	0	1	0

00 $\rightarrow$ [0, 1]	2
01 $\rightarrow$ [2, 3]	0
10 $\rightarrow$ [4, 5]	0
11 $\rightarrow$ [6, 7]	1

# WAVELET TREE

BOTTOM-UP

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0
0	0	1	1	0	0	0	1	1	1
0	1	0	1	1	1	0	0	0	1

Interval on level  $\ell$  is given by *bit prefix* of length  $\ell$

$\ell = 0$

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0

$\ell = 1$

0	1	1	2	3	6	7	5	4	6
0	0	0	1	1	1	1	0	0	1

$\ell = 2$

0	1	1	2	3	5	4	6	7	6
0	1	1	0	1	1	0	0	1	0

00 $\rightarrow$ [0, 1]	2
01 $\rightarrow$ [2, 3]	0
10 $\rightarrow$ [4, 5]	0
11 $\rightarrow$ [6, 7]	2



# WAVELET TREE

BOTTOM-UP

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0
0	0	1	1	0	0	0	1	1	1
0	1	0	1	1	1	0	0	0	1

Interval on level  $\ell$  is given by *bit prefix* of length  $\ell$

$\ell = 0$

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0

00 $\rightarrow$ [0, 1]	3
01 $\rightarrow$ [2, 3]	2
10 $\rightarrow$ [4, 5]	2
11 $\rightarrow$ [6, 7]	3

$\ell = 1$

0	1	1	2	3	6	7	5	4	6
0	0	0	1	1	1	1	0	0	1

$\ell = 2$

0	1	1	2	3	5	4	6	7	6
0	1	1	0	1	1	0	0	1	0

# WAVELET TREE

BOTTOM-UP

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0
0	0	1	1	0	0	0	1	1	1
0	1	0	1	1	1	0	0	0	1

Interval on level  $\ell$  is given by *bit prefix* of length  $\ell$

$\ell = 0$

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0

$\ell = 1$

0	1	1	2	3	6	7	5	4	6
0	0	0	1	1	1	1	0	0	1

$\ell = 2$

0	1	1	2	3	5	4	6	7	6
0	1	1	0	1	1	0	0	1	0

00 $\rightarrow$ [0, 1]	3
01 $\rightarrow$ [2, 3]	2
10 $\rightarrow$ [4, 5]	2
11 $\rightarrow$ [6, 7]	3

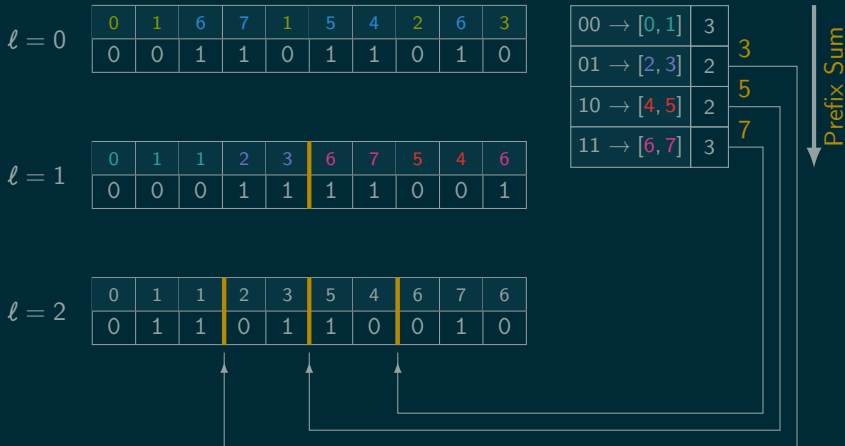
Prefix Sum

# WAVELET TREE

BOTTOM-UP

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0
0	0	1	1	0	0	0	1	1	1
0	1	0	1	1	1	0	0	0	1

Interval on level  $\ell$  is given by *bit prefix* of length  $\ell$



# WAVELET TREE

BOTTOM-UP

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0
0	0	1	1	0	0	0	1	1	1
0	1	0	1	1	1	0	0	0	1

Interval on level  $\ell$  is given by *bit prefix* of length  $\ell$

$\ell = 0$

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0

00 $\rightarrow$ [0, 1]	3
01 $\rightarrow$ [2, 3]	2
10 $\rightarrow$ [4, 5]	2
11 $\rightarrow$ [6, 7]	3

$\ell = 1$

0	1	1	2	3	6	7	5	4	6
0	0	0	1	1	1	1	0	0	1

$\ell = 2$

0	1	1	2	3	5	4	6	7	6
0	1	1	0	1	1	0	0	1	0

# WAVELET TREE

BOTTOM-UP

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0
0	0	1	1	0	0	0	1	1	1
0	1	0	1	1	1	0	0	0	1

Interval on level  $\ell$  is given by *bit prefix* of length  $\ell$

$\ell = 0$

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0

$\ell = 1$

0	1	1	2	3	6	7	5	4	6
0	0	0	1	1	1	1	0	0	1

$\ell = 2$

0	1	1	2	3	5	4	6	7	6
0	1	1	0	1	1	0	0	1	0

00 $\rightarrow$ [0, 1]	3
01 $\rightarrow$ [2, 3]	2
10 $\rightarrow$ [4, 5]	2
11 $\rightarrow$ [6, 7]	3

Compute the borders for the **previous** level without rescanning T.

# WAVELET TREE

BOTTOM-UP

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0
0	0	1	1	0	0	0	1	1	1
0	1	0	1	1	1	0	0	0	1

Interval on level  $\ell$  is given by *bit prefix* of length  $\ell$

$\ell = 0$

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0

$\ell = 1$

0	1	1	2	3	6	7	5	4	6
0	0	0	1	1	1	1	0	0	1

$\ell = 2$

0	1	1	2	3	5	4	6	7	6
0	1	1	0	1	1	0	0	1	0

00 $\rightarrow$ [0, 1]	3
01 $\rightarrow$ [2, 3]	5
10 $\rightarrow$ [4, 5]	2
11 $\rightarrow$ [6, 7]	5

Compute the borders for the **previous** level without rescanning T.

# RECAP

## BOTTOM-UP

1. Compute histogram for level  $\ell = \lceil \lg \sigma \rceil$  (plus bit vector of first level)
2. Compute borders for level  $\ell$
3. Fill bit vector for level  $\ell$  accordingly
4. Compute histogram for level  $\ell - 1$  using the current histogram
5. If  $\ell \geq 1$  set  $\ell = \ell - 1$  and repeat at step 2

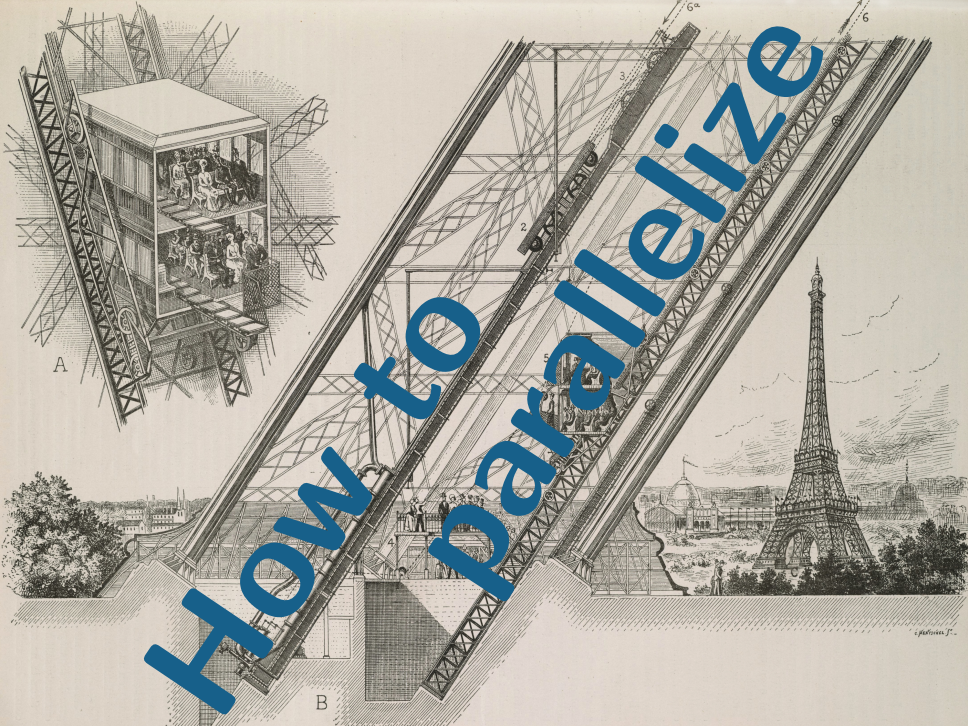
# RECAP

## BOTTOM-UP

1. Compute histogram for level  $\ell = \lceil \lg \sigma \rceil$  (plus bit vector of first level)
2. Compute borders for level  $\ell$
3. Fill bit vector for level  $\ell$  accordingly
4. Compute histogram for level  $\ell - 1$  using the current histogram
5. If  $\ell \geq 1$  set  $\ell = \ell - 1$  and repeat at step 2



How to parallelize



A

B

J. G. W. J.

# PREFIX COUNTING

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0
0	0	1	1	0	0	0	1	1	1
0	1	0	1	1	1	0	0	0	1

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0

0	1	1	2	3	6	7	5	4	6
0	0	0	1	1	1	1	0	0	1

0	1	1	2	3	5	4	6	7	6
0	1	1	0	1	1	0	0	1	0

# PREFIX COUNTING

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0
0	0	1	1	0	0	0	1	1	1
0	1	0	1	1	1	0	0	0	1

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0

---

0	1	1	2	3	6	7	5	4	6
0	0	0	1	1	1	1	0	0	1

---

0	1	1	2	3	5	4	6	7	6
0	1	1	0	1	1	0	0	1	0

# PREFIX COUNTING

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0
0	0	1	1	0	0	0	1	1	1
0	1	0	1	1	1	0	0	0	1

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0

---

0	1	1	2	3	6	7	5	4	6
0	0	0	1	1	1	1	0	0	1

---

0	1	1	2	3	5	4	6	7	6
0	1	1	0	1	1	0	0	1	0

$0 \rightarrow [0, 3]$	3
$1 \rightarrow [4, 7]$	2

$00 \rightarrow [0, 1]$	3
$01 \rightarrow [2, 3]$	2
$10 \rightarrow [4, 5]$	2
$11 \rightarrow [6, 7]$	3

# PREFIX COUNTING

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0
0	0	1	1	0	0	0	1	1	1
0	1	0	1	1	1	0	0	0	1

**Problem:** Utilizes at most  $\lceil \lg \sigma \rceil$  cores

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0

---

0	1	1	2	3	6	7	5	4	6
0	0	0	1	1	1	1	0	0	1

---

0	1	1	2	3	5	4	6	7	6
0	1	1	0	1	1	0	0	1	0

0 → [0, 3]	3
1 → [4, 7]	2

00 → [0, 1]	3
01 → [2, 3]	2
10 → [4, 5]	2
11 → [6, 7]	3

# PREFIX SORTING

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0
0	0	1	1	0	0	0	1	1	1
0	1	0	1	1	1	0	0	0	1

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0

0	1	1	2	3	6	7	5	4	6
0	0	0	1	1	1	1	0	0	1

0	1	1	2	3	5	4	6	7	6
0	1	1	0	1	1	0	0	1	0

# PREFIX SORTING

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0
0	0	1	1	0	0	0	1	1	1
0	1	0	1	1	1	0	0	0	1

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0

0	1	1	2	3	6	7	5	4	6
0	0	0	1	1	1	1	0	0	1

0	1	1	2	3	5	4	6	7	6
0	1	1	0	1	1	0	0	1	0

# PREFIX SORTING

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0
0	0	1	1	0	0	0	1	1	1
0	1	0	1	1	1	0	0	0	1

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0

0	1	1	2	3	6	7	5	4	6
0	0	0	1	1	1	1	0	0	1

0	1	1	2	3	5	4	6	7	6
0	1	1	0	1	1	0	0	1	0

Starting with the last level:

- ▶ Compute borders
- ▶ Sort text (keeping original)
- ▶ Fill bit vector



# DOMAIN DECOMPOSITION

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0
0	0	1	1	0	0	0	1	1	1
0	1	0	1	1	1	0	0	0	1

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0

0	1	1	2	3	6	7	5	4	6
0	0	0	1	1	1	1	0	0	1

0	1	1	2	3	5	4	6	7	6
0	1	1	0	1	1	0	0	1	0

# DOMAIN DECOMPOSITION

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0
0	0	1	1	0	0	0	1	1	1
0	1	0	1	1	1	0	0	0	1

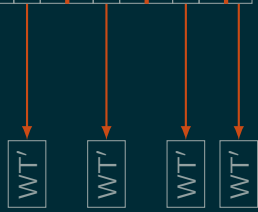
0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0

0	1	1	2	3	6	7	5	4	6
0	0	0	1	1	1	1	0	0	1

0	1	1	2	3	5	4	6	7	6
0	1	1	0	1	1	0	0	1	0

# DOMAIN DECOMPOSITION

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0
0	0	1	1	0	0	0	1	1	1
0	1	0	1	1	1	0	0	0	1



0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0

0	1	1	2	3	6	7	5	4	6
0	0	0	1	1	1	1	0	0	1

0	1	1	2	3	5	4	6	7	6
0	1	1	0	1	1	0	0	1	0

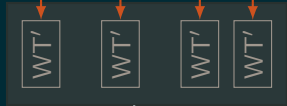
# DOMAIN DECOMPOSITION

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0
0	0	1	1	0	0	0	1	1	1
0	1	0	1	1	1	0	0	0	1

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0

0	1	1	2	3	6	7	5	4	6
0	0	0	1	1	1	1	0	0	1

0	1	1	2	3	5	4	6	7	6
0	1	1	0	1	1	0	0	1	0



Merge



# RECAP

## PARALLELIZATION

- ▶ Prefix Counting is simple but not scaling

# RECAP

## PARALLELIZATION

- ▶ Prefix Counting is simple but not scaling
- ▶ Prefix Sorting requires more space but scales

# RECAP

## PARALLELIZATION

- ▶ Prefix Counting is simple but not scaling
- ▶ Prefix Sorting requires more space but scales
- ▶ Domain Decomposition scales and is flexible (but needs more space)

# WAVELET MATRIX





# WAVELET MATRIX

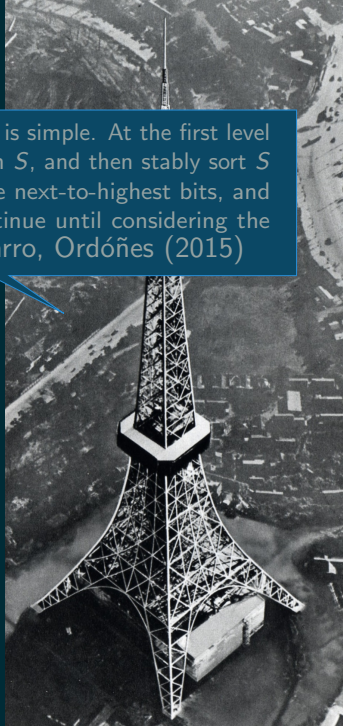
SIMILAR BUT DIFFERENT



# WAVELET MATRIX

SIMILAR BUT DIFFERENT

**Construction:** The construction of the wavelet matrix is simple. At the first level we keep in bitmap  $\tilde{B}_0$  the highest bit of the symbols in  $S$ , and then stably sort  $S$  by those highest bits. Now we keep in bitmap  $\tilde{B}_1$  the next-to-highest bits, and stably sort  $S$  by those next-to-highest bits. We continue until considering the lowest bit. This takes  $\mathcal{O}(n \lg \sigma)$  time. Claude, Navarro, Ordóñez (2015)



# WAVELET MATRIX

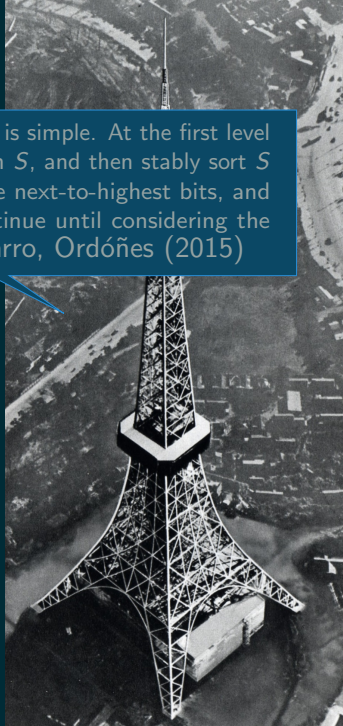
SIMILAR BUT DIFFERENT

**Construction:** The construction of the wavelet matrix is simple. At the first level we keep in bitmap  $\tilde{B}_0$  the highest bit of the symbols in  $S$ , and then stably sort  $S$  by those highest bits. Now we keep in bitmap  $\tilde{B}_1$  the next-to-highest bits, and stably sort  $S$  by those next-to-highest bits. We continue until considering the lowest bit. This takes  $\mathcal{O}(n \lg \sigma)$  time. Claude, Navarro, Ordóñez (2015)

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0

0	1	1	2	3	6	7	5	4	6
0	0	0	1	1	1	1	0	0	1

0	1	1	5	4	2	3	6	7	6
0	1	1	1	0	0	1	0	1	0



# WAVELET MATRIX

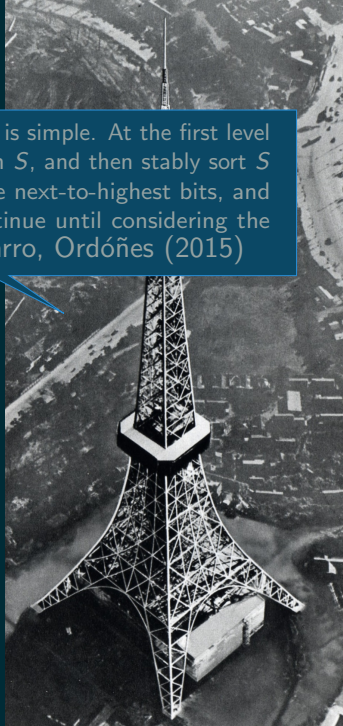
SIMILAR BUT DIFFERENT

**Construction:** The construction of the wavelet matrix is simple. At the first level we keep in bitmap  $\tilde{B}_0$  the highest bit of the symbols in  $S$ , and then stably sort  $S$  by those highest bits. Now we keep in bitmap  $\tilde{B}_1$  the next-to-highest bits, and stably sort  $S$  by those next-to-highest bits. We continue until considering the lowest bit. This takes  $\mathcal{O}(n \lg \sigma)$  time. Claude, Navarro, Ordóñez (2015)

0	1	6	7	1	5	4	2	6	3
0	0	1	1	0	1	1	0	1	0

0	1	1	2	3	6	7	5	4	6
0	0	0	1	1	1	1	0	0	1

0	1	1	5	4	2	3	6	7	6
0	1	1	1	0	0	1	0	1	0



# WAVELET MATRIX

CONSIDERED BIT PREFIXES

- ▶ Bit-reversal permutation

# WAVELET MATRIX

CONSIDERED BIT PREFIXES

- ▶ Bit-reversal permutation

0	1
---	---

# WAVELET MATRIX

CONSIDERED BIT PREFIXES

- ▶ Bit-reversal permutation

0	1		
00	10	01	11

# WAVELET MATRIX

CONSIDERED BIT PREFIXES

- ▶ Bit-reversal permutation

0	1						
00	10	01	11				
000	100	010	110	001	101	011	111



# WAVELET MATRIX

CONSIDERED BIT PREFIXES

- ▶ Bit-reversal permutation

0	1														
00	10	01	11												
000	100	010	110	001	101	011	111								
0000	1000	0100	1100	0010	1010	0110	1110	0001	1001	0101	1101	0011	1011	0111	1111

# WAVELET MATRIX

CONSIDERED BIT PREFIXES

- ▶ Bit-reversal permutation
- ▶ Compute borders in this order

0	1														
00	10	01	11												
000	100	010	110	001	101	011	111								
0000	1000	0100	1100	0010	1010	0110	1110	0001	1001	0101	1101	0011	1011	0111	1111

# WAVELET MATRIX

CONSIDERED BIT PREFIXES

- ▶ Bit-reversal permutation
- ▶ Compute borders in this order

0	1														
00	10	01	11												
000	100	010	110	001	101	011	111								
0000	1000	0100	1100	0010	1010	0110	1110	0001	1001	0101	1101	0011	1011	0111	1111

# WAVELET MATRIX

CONSIDERED BIT PREFIXES

- ▶ Bit-reversal permutation
- ▶ Compute borders in this order

0	1														
00	10	01	11												
000	100	010	110	001	101	011	111								
0000	1000	0100	1100	0010	1010	0110	1110	0001	1001	0101	1101	0011	1011	0111	1111

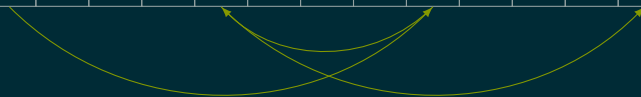
The diagram illustrates a bit-reversal permutation. Three yellow curved arrows point from the left side of the table to the right side, indicating the mapping of bit-reversed prefixes to their original order. The arrows connect 0000 to 1111, 0010 to 1010, and 0001 to 1001.

# WAVELET MATRIX

CONSIDERED BIT PREFIXES

- ▶ Bit-reversal permutation
- ▶ Compute borders in this order

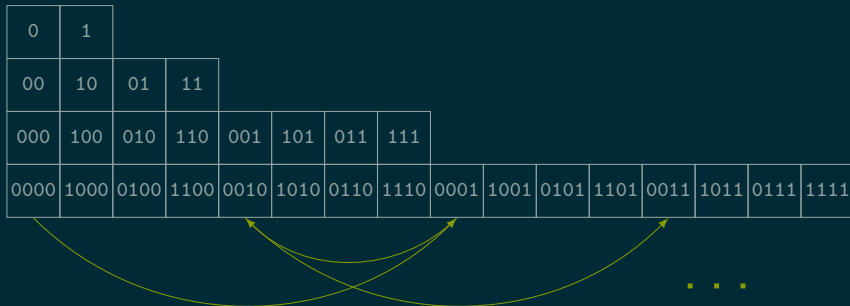
0	1														
00	10	01	11												
000	100	010	110	001	101	011	111								
0000	1000	0100	1100	0010	1010	0110	1110	0001	1001	0101	1101	0011	1011	0111	1111



# WAVELET MATRIX

CONSIDERED BIT PREFIXES

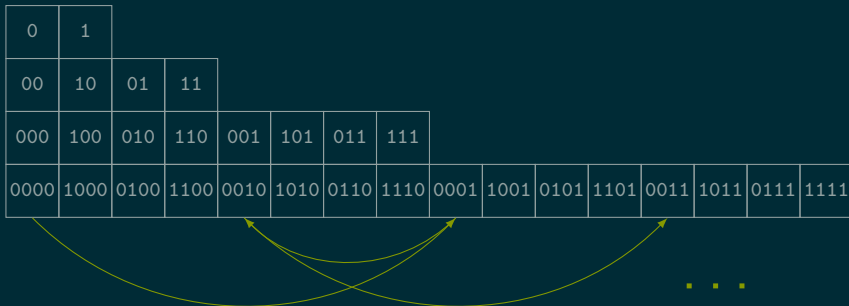
- ▶ Bit-reversal permutation
- ▶ Compute borders in this order



# WAVELET MATRIX

CONSIDERED BIT PREFIXES

- ▶ Bit-reversal permutation
- ▶ Compute borders in this order



- ▶ All algorithms presented can also compute the Wavelet Matrix

# EXPERIMENTS

## SETTING

- ▶ Workstation with 32 cores and 256 GB RAM
- ▶ 11.8 GB Proteins with  $\sigma = 27$
- ▶ 4 GB DNA with  $\sigma = 16$

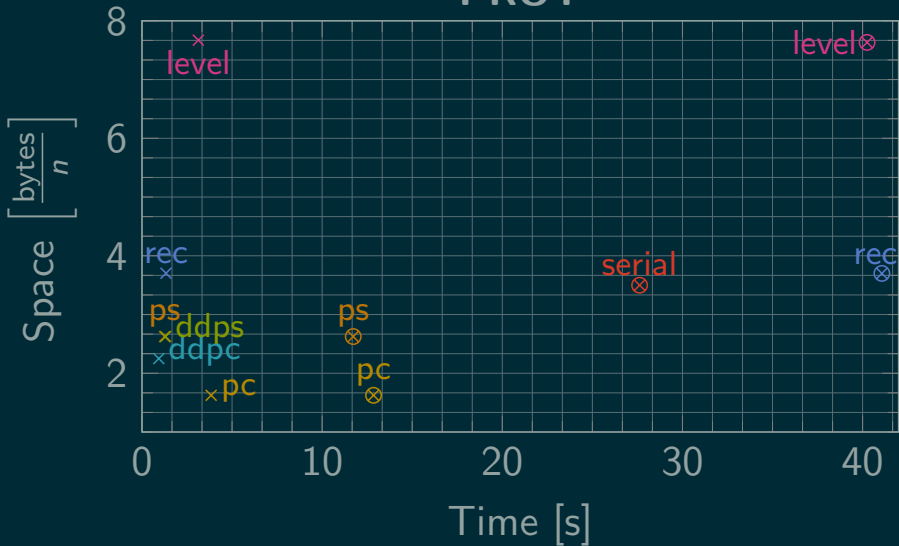


# EXPERIMENTS

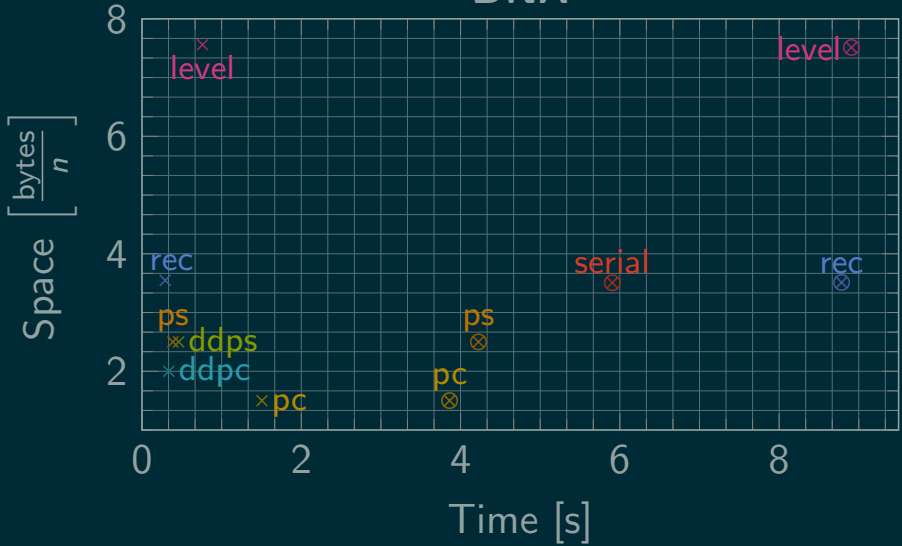
## SETTING

- ▶ Workstation with 32 cores and 256 GB RAM
- ▶ 11.8 GB Proteins with  $\sigma = 27$
- ▶ 4 GB DNA with  $\sigma = 16$
- ▶ Test it yourself: [www.github.com/kurpicz/pwm](http://www.github.com/kurpicz/pwm)

# PROT



# DNA



# CONCLUSION & OUTLOOK

- ▶ New algorithms for the Wavelet Tree
- ▶ First practical algorithms for the Wavelet Matrix
- ▶ Fast, simple and lightweight

# CONCLUSION & OUTLOOK

- ▶ New algorithms for the Wavelet Tree
- ▶ First practical algorithms for the Wavelet Matrix
- ▶ Fast, simple and lightweight
- ▶ Better scaling possible?

# CONCLUSION & OUTLOOK

- ▶ New algorithms for the Wavelet Tree
- ▶ First practical algorithms for the Wavelet Matrix
- ▶ Fast, simple and lightweight
- ▶ Better scaling possible?

Thank You!