# Introduction to Wavelet Trees

Giovanni Manzini

# rank and select operations

Given a string $S[1,n]$ over and alphabet $A$ we define:

$\text{rank}_S(c,i)$ = # occs of symbol $c$ in $S[1,i]$

$\text{select}_S(c,i)$ = position of the $i$-th $c$ in $S[1,n]$

Example:

$S$ = abracadabra

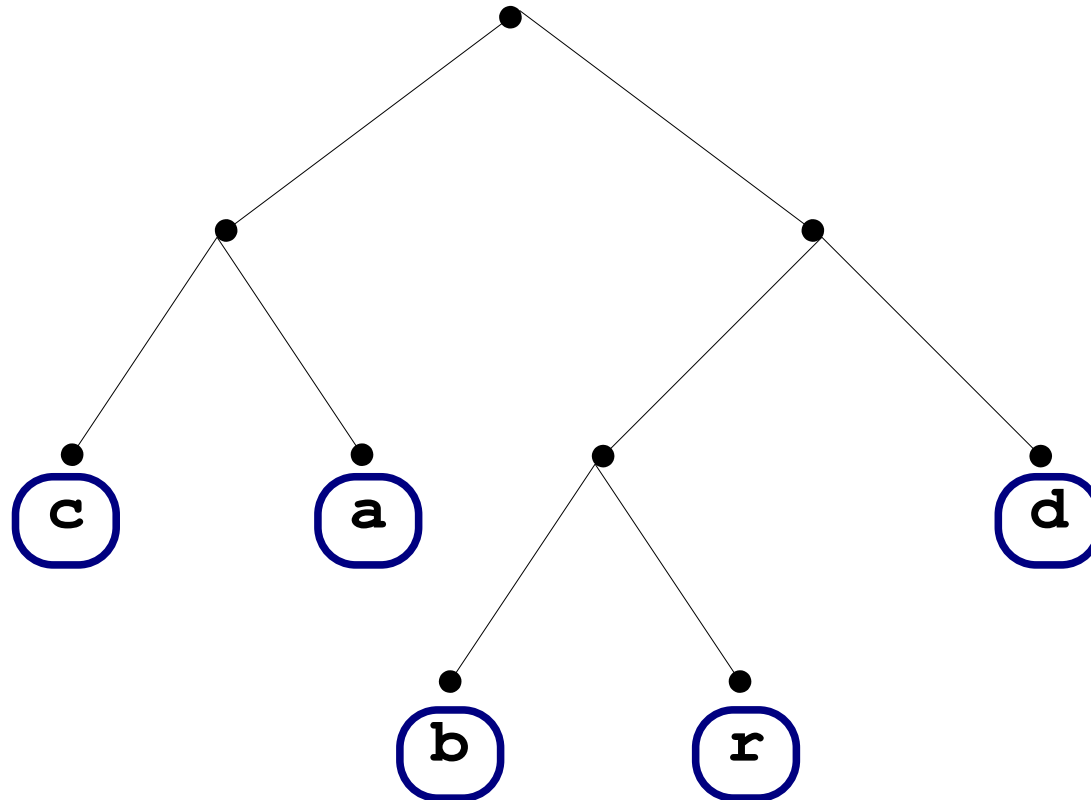$\text{rank}_S(a,3)=1$,    $\text{select}_S(a,3)=6$

# The Wavelet Tree data structure

Wavelet Trees have been introduced to represent compactly a string supporting rank/select operations.
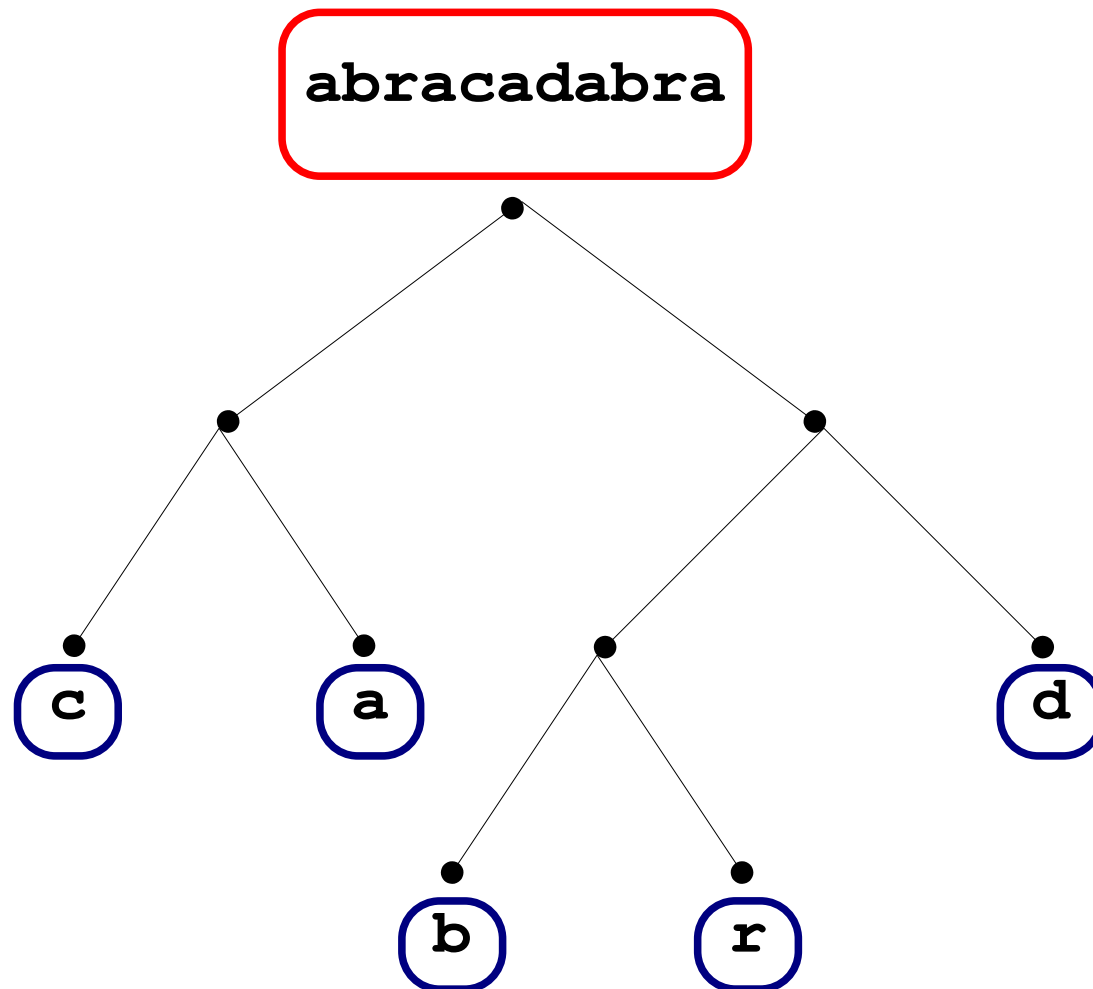
They work assuming we can do rank/select on binary strings,

To build a Wavelet Tree we start with a complete binary tree with a leaf for each alphabet symbol.
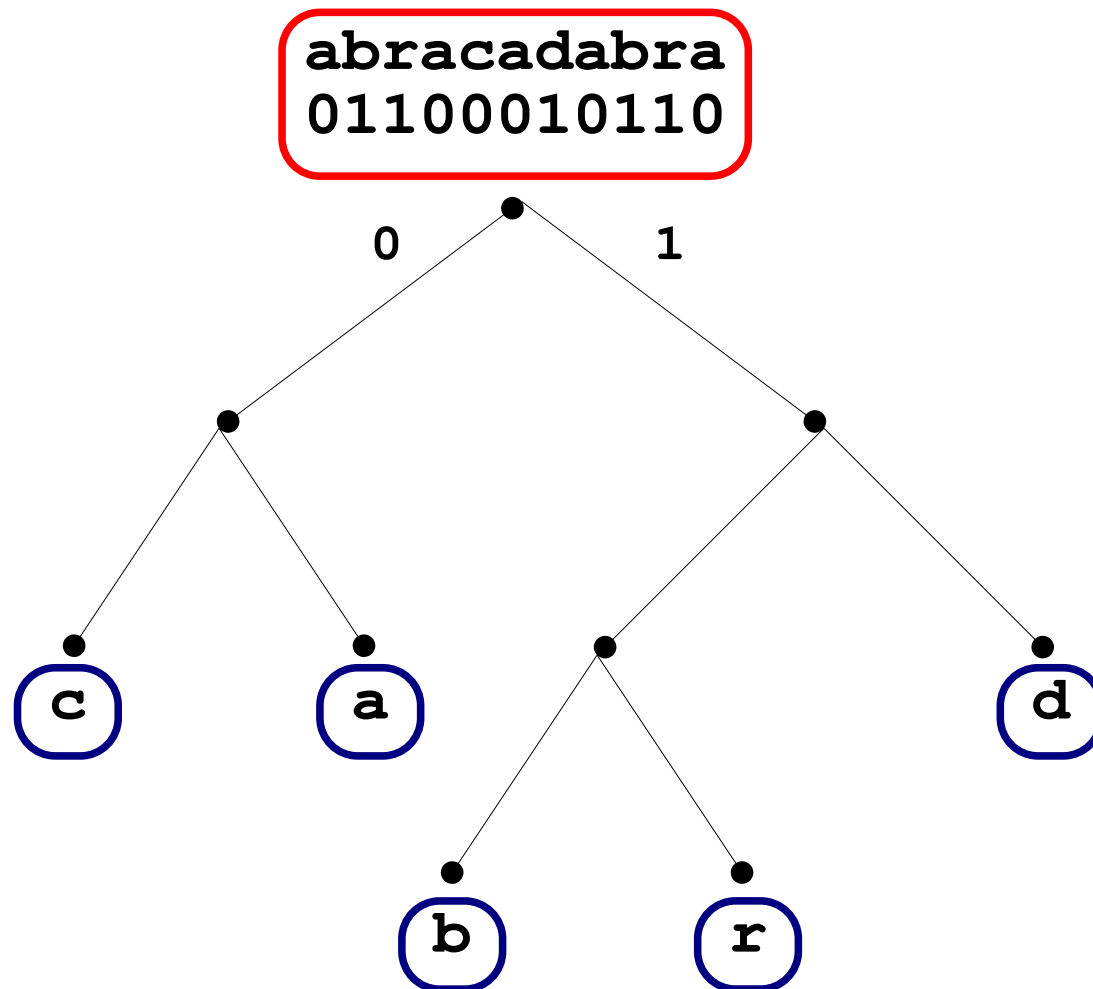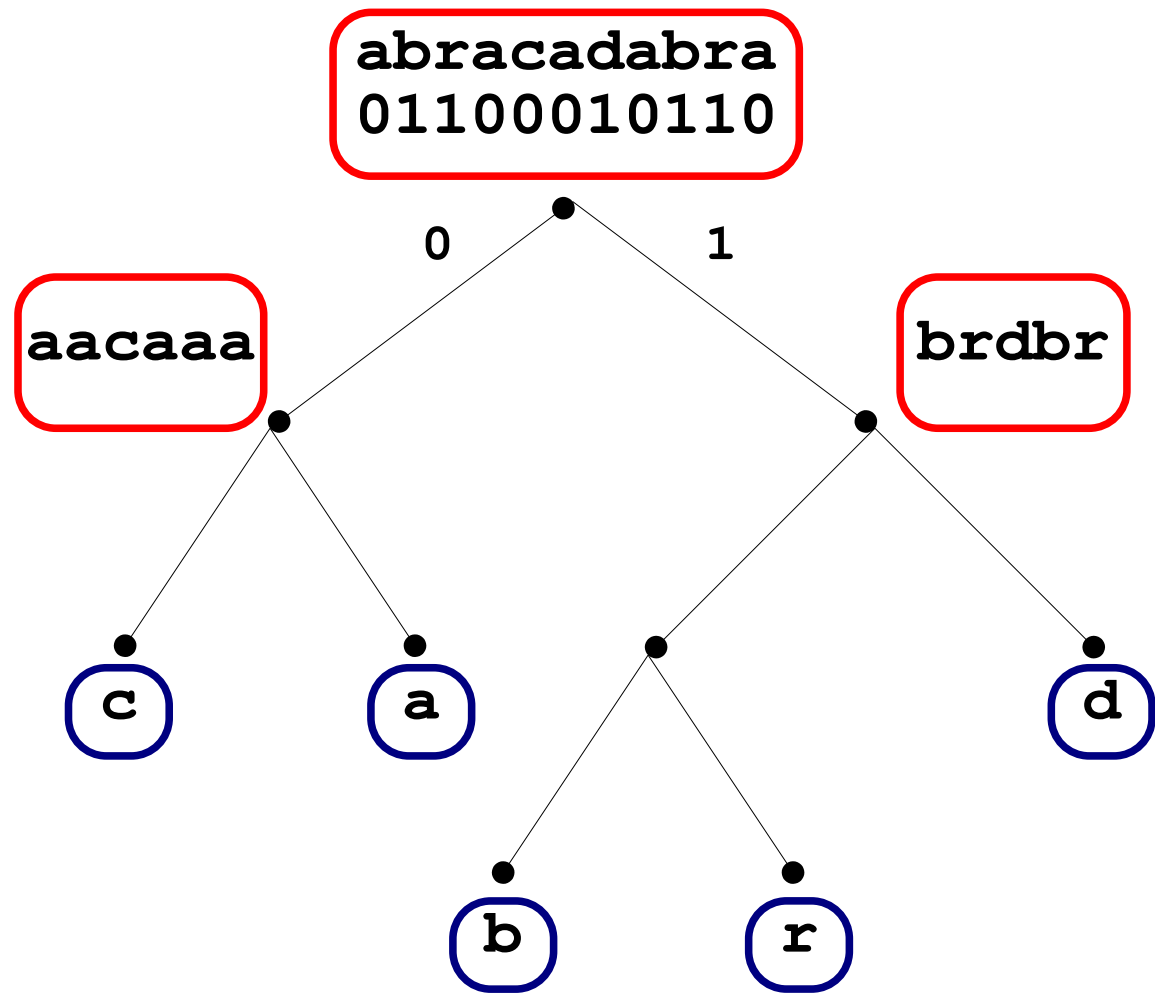
Example.   A = {a, b, c, d, r}

Given the string **abracadabra** we build the corresponding Wavelet Tree as follows:

Given the string **abracadabra** we build the corresponding Wavelet Tree as follows:

Given the string **abracadabra** we build the corresponding Wavelet Tree as follows:

Given the string **abracadabra** we build the corresponding Wavelet Tree as follows:

Given the string **abracadabra** we build the corresponding Wavelet Tree as follows:
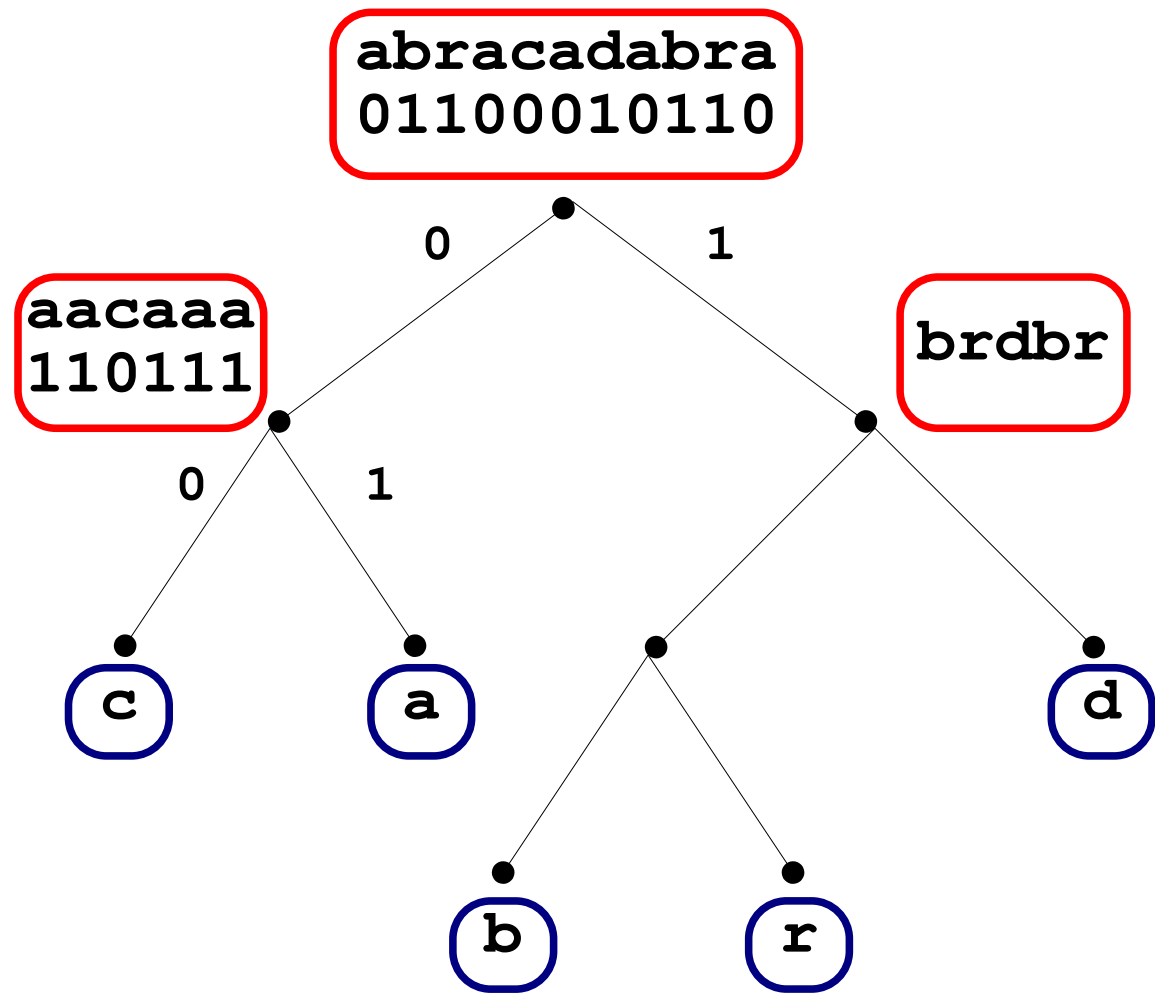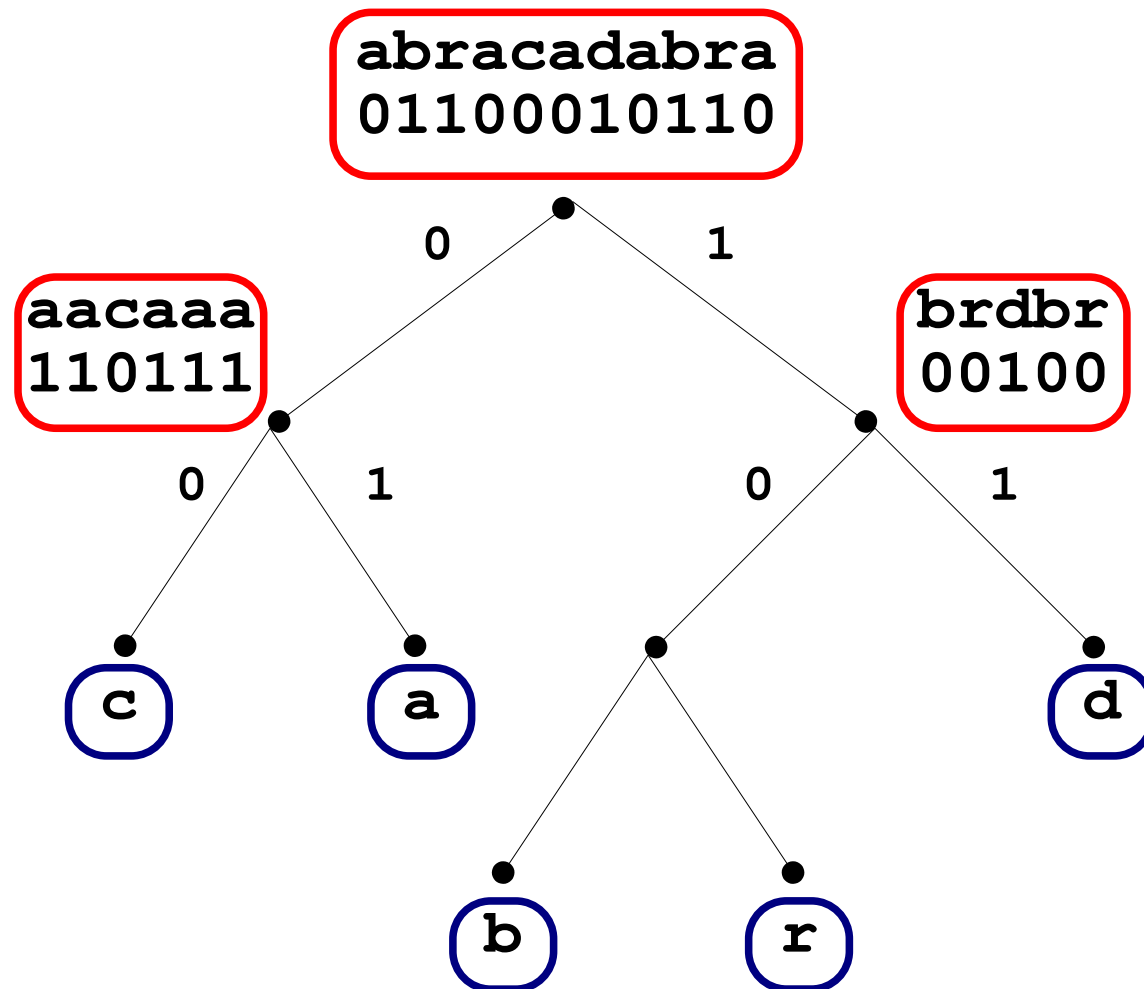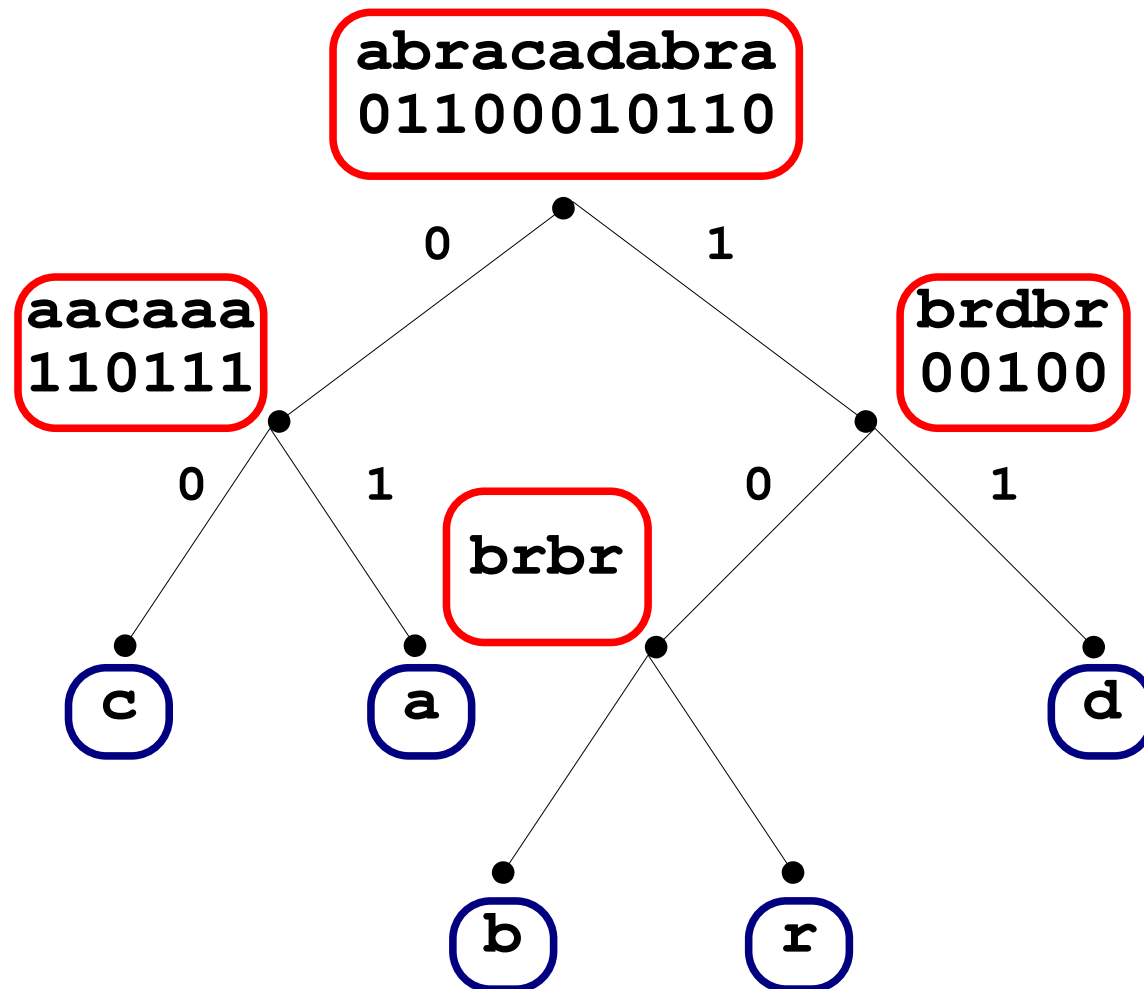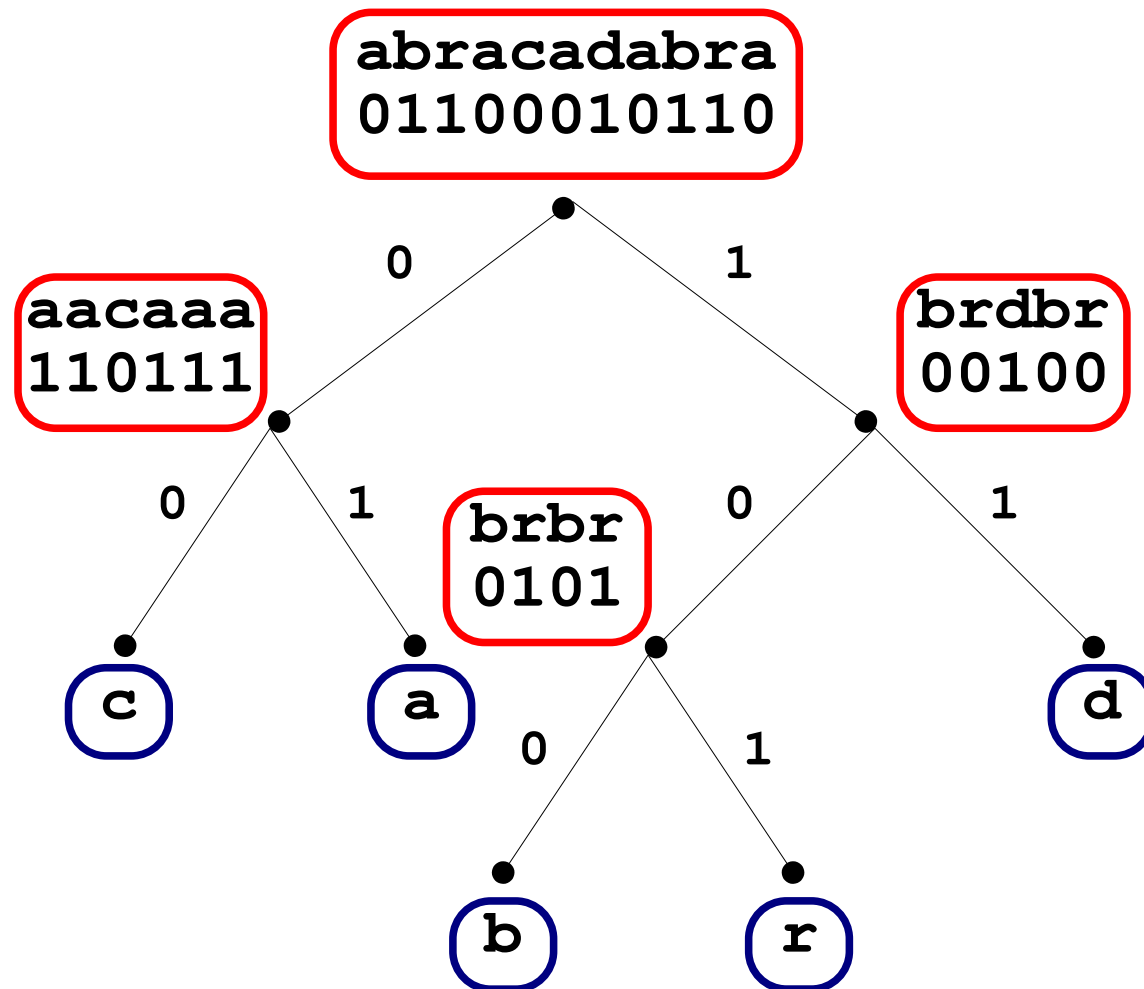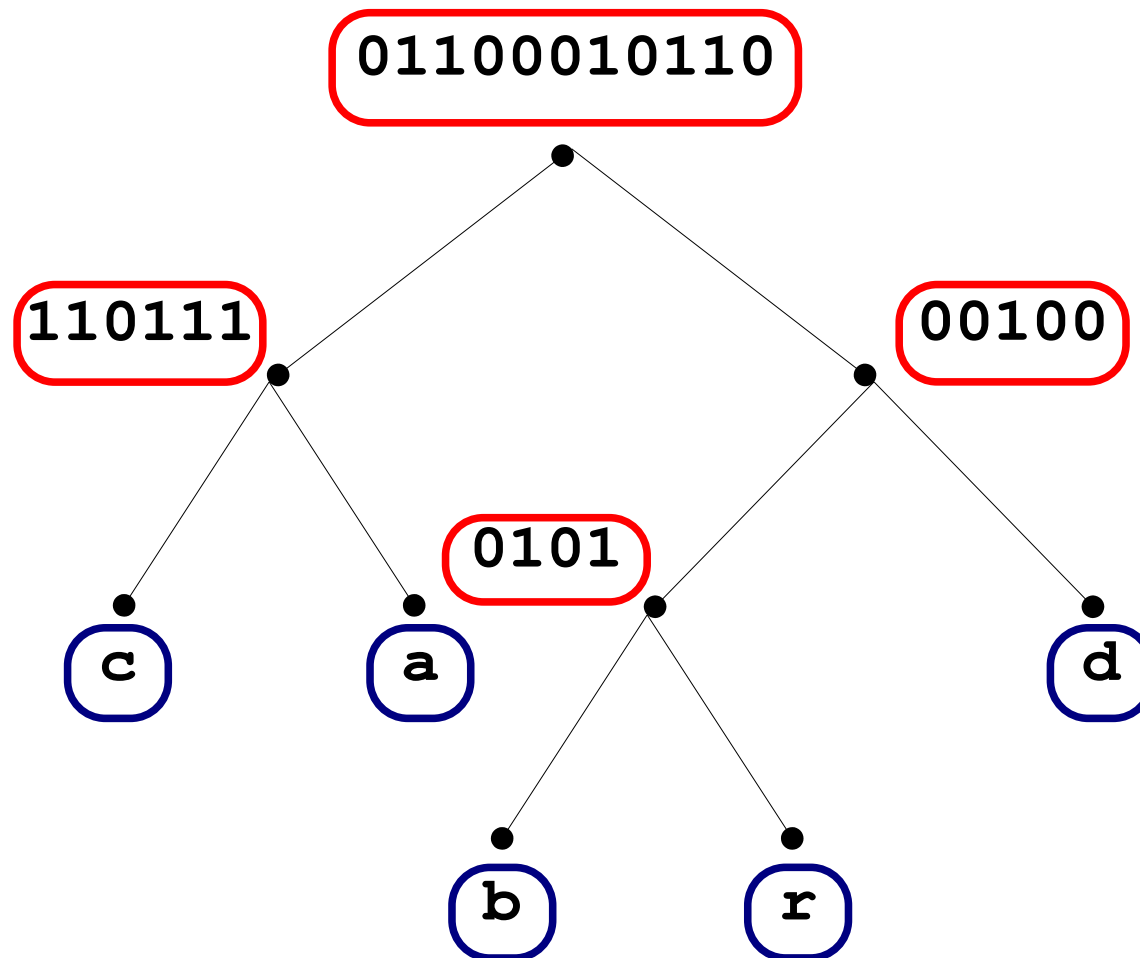
Given the string **abracadabra** we build the corresponding Wavelet Tree as follows:

# Given the string **abracadabra** we build the corresponding Wavelet Tree as follows:

rank/select queries over the original string can be answered via rank/select queries over the binary strings associated to internal nodes.

# Suppose we want to compute rank(b,9)

# Suppose we want to compute rank(b,9)



$rk_1(9) = rk(b,9) + rk(r,9) + rk(d,9) = 4$

**abracadabra**
**01100010110**

0      1

**aacaaa**
**110111**

**brdbr**
**00100**

0      1

0      1

**brbr**
**0101**

c      a

0      1

d

b      r

# Suppose we want to compute rank(b,9)



rk$_1$(9)=rk(b,9)+rk(r,9)+rk(d,9)=4

abracadabra
01100010110

aacaaa
110111

brdbr
00100

0          1

0      1              0      1

brbr
0101

c              a          0      1          d

b          r

# Suppose we want to compute rank(b,9)

**abracadabra**
**01100010110**

$rk_1(9)=rk(b,9)+rk(r,9)+rk(d,9)=4$

0          1

**aacaaa**
**110111**

**brdbr**
**00100**

$rk_0(4)=rk(b,4)+rk(r,4)=3$

0       1                    0          1

**brbr**
**0101**

c          a

0       1

d

b          r

# Suppose we want to compute rank(b,9)



$rk_1(9) = rk(b,9) + rk(r,9) + rk(d,9) = 4$

**abracadabra**
**01100010110**

**aacaaa**
**110111**

**brdbr**
**00100**

$rk_0(4) = rk(b,4) + rk(r,4) = 3$

**brbr**
**0101**

c    a    d

b    r

# Suppose we want to compute rank(b,9)



**abracadabra**
**01100010110**

$rk_1(9)=rk(b,9)+rk(r,9)+rk(d,9)=4$

**aacaaa**
**110111**

**brdbr**
**00100**

$rk_0(4)=rk(b,4)+rk(r,4)=3$

**brbr**
**0101**

$rk_0(3)=rk(b,3)=2$

0

1

0

1

0

1

0

1

c

a

d

b

r

# Suppose we want to compute rank(b,9)

# Select queries go upward. To compute select(a,3)



abracadabra
01100010110

0   1

aacaaa
110111

brdbr
00100

select₁(3) → 4

0   1

0   1

brbr
0101

c   a

0   1

d

b   r

# Select queries go upward. To compute select(a,3)



select$_0$(4) → 6

**abracadabra**
**01100010110**

select$_1$(3) → 4

**aacaaa**
**110111**

**brdbr**
**00100**

0          1

0          1

**brbr**
**0101**

0          1

c          a

0          1

d

b          r

# Select queries go upward. To compute select(a,3)

$select_0(4) \to 6$

**abracadabra**
**01100010110**

0          1

**aacaaa**
**110111**

**brdbr**
**00100**

$select_1(3) \to 4$

0          1

0          1

**brbr**
**0101**

c          a

0          1

d

Answer: select(a,3)=6

b          r

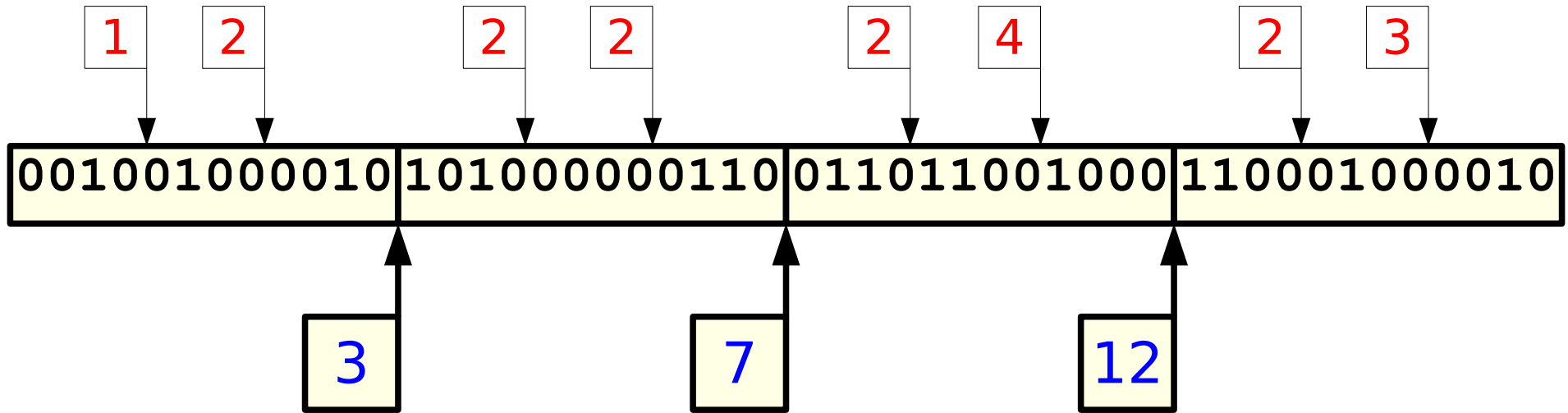# Rank$_1$ queries on binary strings

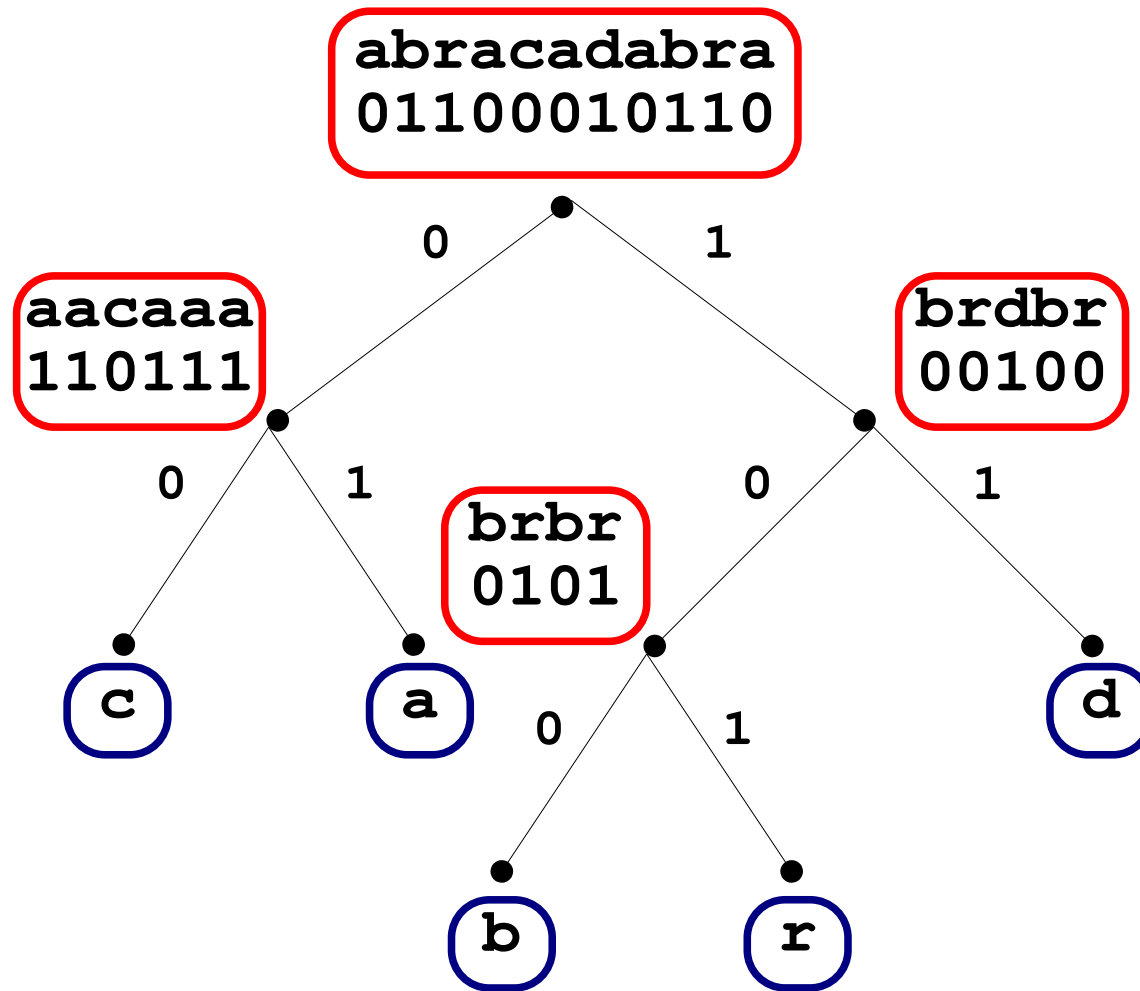Simplest idea: split the binary string into blocks and store a partial sum for each block:

| 00100100000010 | 10100000000110 | 011011001000 | 110000000011 |
|---|---|---|---|

| 3 | 7 | 12 |

We can refine the idea using blocks and mini-blocks:

| 1 | 2 | | 2 | 2 | | 2 | 4 | | 2 | 3 |

001001000010 101000000110 011011001000 110001000010

| 3 | | 7 | | 12 |

Similar techniques for select queries!

# Wavelet Trees as scrambled prefix codes



**abracadabra**
**01100010110**

0 / \ 1

**aacaaa**
**110111**

**brdbr**
**00100**

0 / \ 1

0 / \ 1

**brbr**
**0101**

c    a

0 / \ 1

d

b    r

We are implicitly using the encoding:

$c \to 00$  $a \to 01$  $b \to 100$  $r \to 101$  $d \to 11$

# Wavelet Trees as scrambled prefix codes:

## We are implicitly using the encoding:

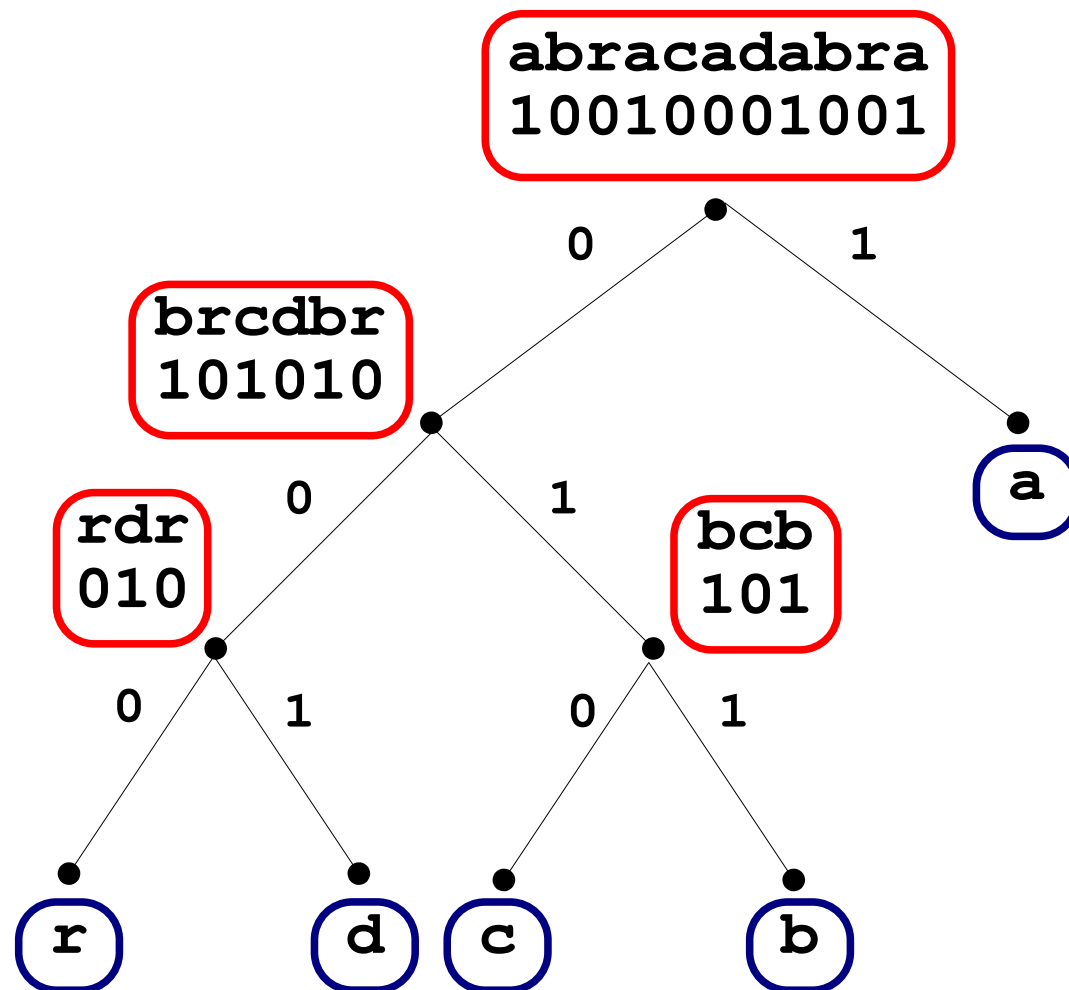$a \to 01$  $c \to 00$  $b \to 100$  $r \to 101$  $d \to 11$

## Traditional prefix codes:

$abracadabra \to 01\ 100\ 101\ 01\ 00\ 01\ 11 \ldots$

## Wavelets:

Changing the shape of the the Wavelet Tree, all properties mentioned above still hold with a different prefix code.



abracadabra
10010001001

0        1

brcdbr
101010

0        1

rdr
010

bcb
101

a

0    1        0    1

r        d  c        b

r→000
d→001
c→010
b→011
a→1

We can choose any prefix-free binary encoding of the characters and build the corresponding Wavelet Tree.

If we choose Huffman codes we have compression for free and direct access to arbitrary positions in the text.

# Summing up

Binary strings operations **+** Wavelet Trees **+**

**=** Compression and efficient rank support **+**

BWT & backward search **=** Compressed full text index

# Other Wavelet Trees virtues

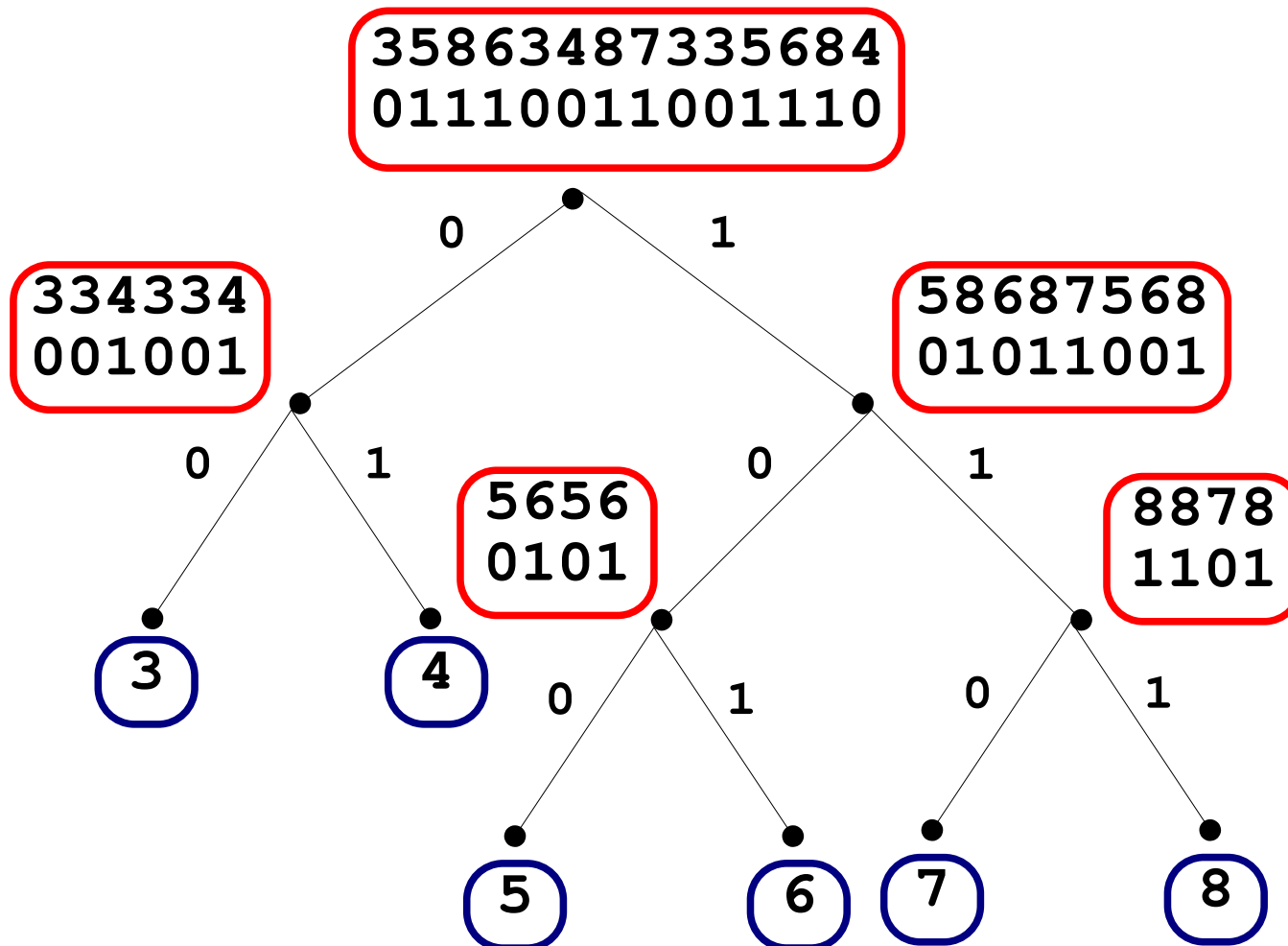Simple solutions for range queries on an integer sequence S[1,n]

range_quantile(S,i,j,k): return the k-th smallest value in S[i···j]

range_next(S,i,j,x): return the smallest value greater than x in S[i···j]
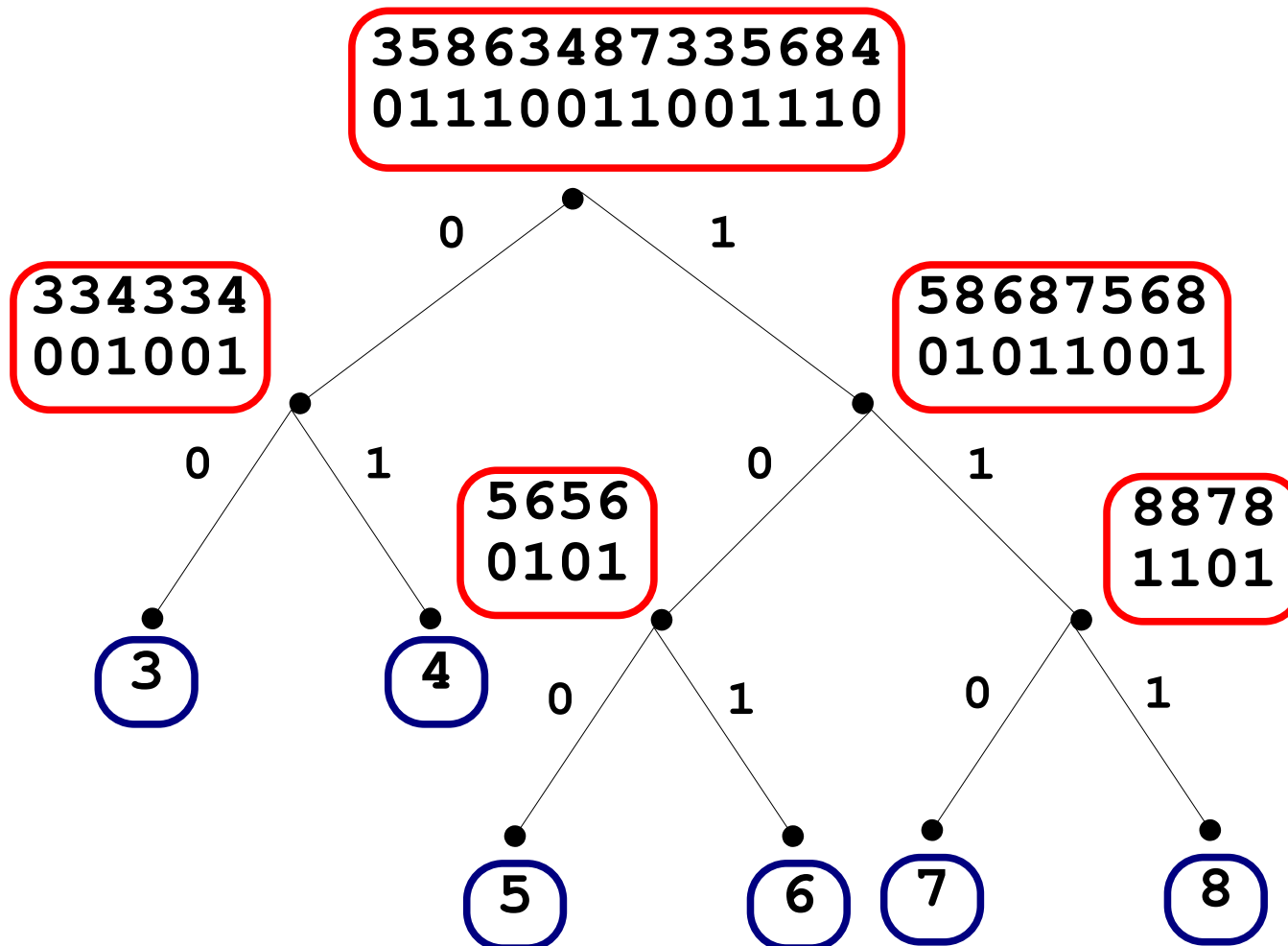
range_count(S,i,j): return # of distinct values in S[i···j]

# Example.

S = [3, 5, 8, 6, 3, 4, 8, 7, 3, 3, 5, 6, 8, 4]

range_quantile(S,2,7,3): 3-rd smallest value in S[2···7]

S = [3,5,**8,6,3,4,8,7**,3,3,5,6,8,4]

358634873356840
01110011001110

334334
001001

58687568
01011001

5656
0101

8878
1101

0   1
0   1
0   1

3   4   5   6   7   8
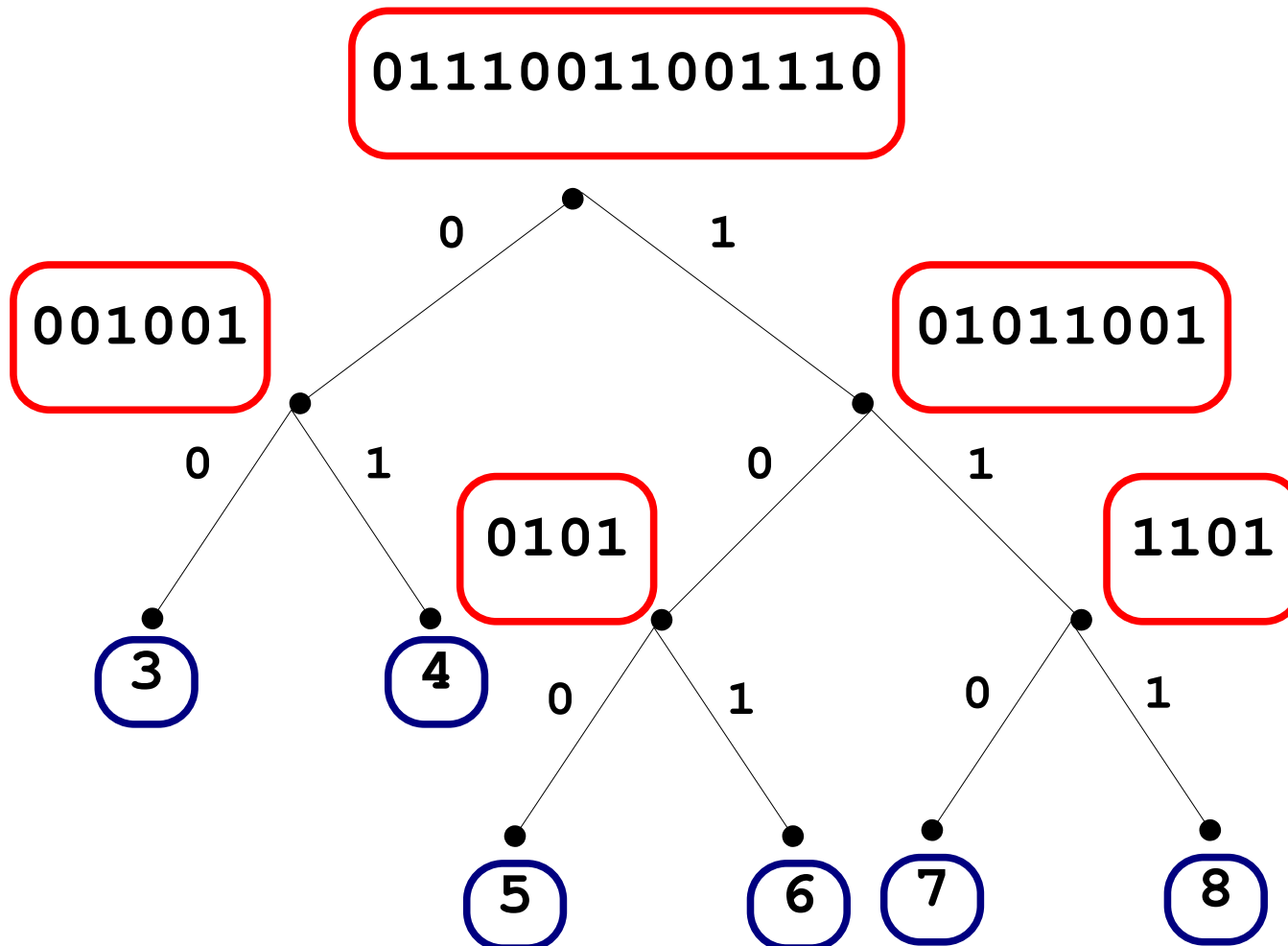
range_quantile(S,2,7,3): 3-rd smallest value in S[2···7]

S = [3,5,8,6,3,4,8,7,3,3,5,6,8,4]

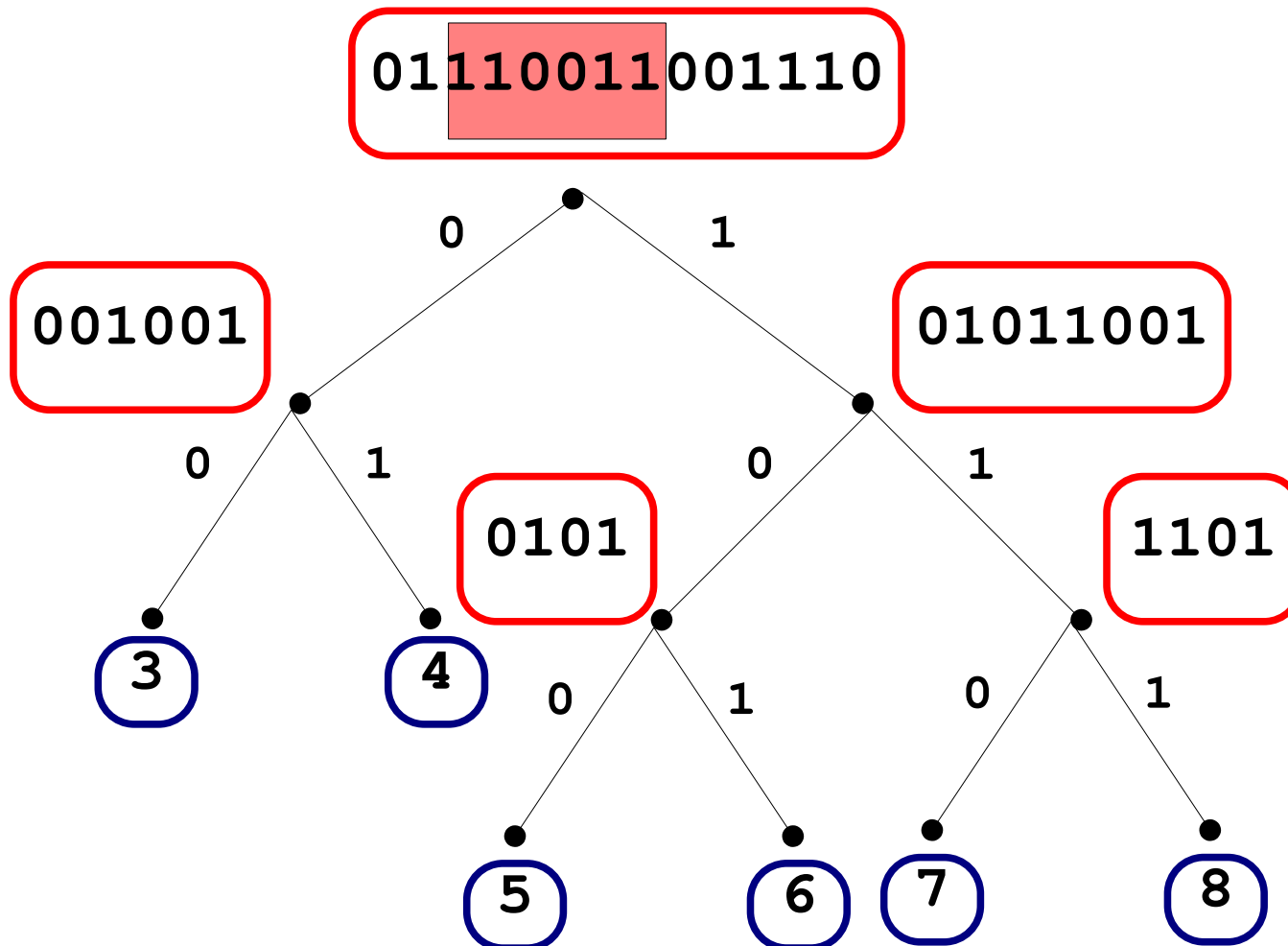range_quantile(S,2,7,3): 3-rd smallest value in S[2⋯7]

S = [3,5,**8,6,3,4,8,7**,3,3,5,6,8,4]

range_quantile(S,2,7,3): 3-rd smallest value in S[2···7]

S = [3,5,**8,6,3,4,8,7**,3,3,5,6,8,4]
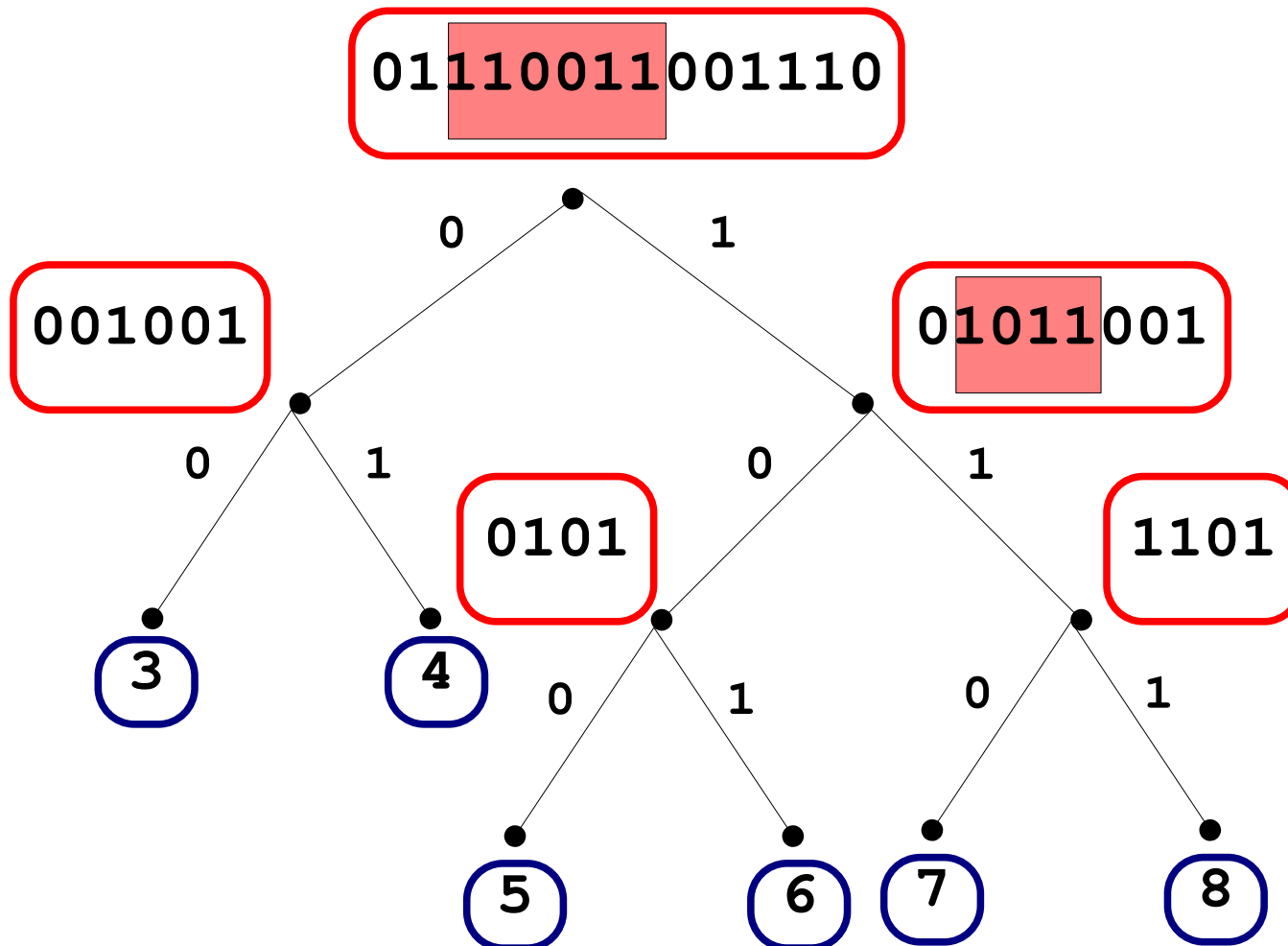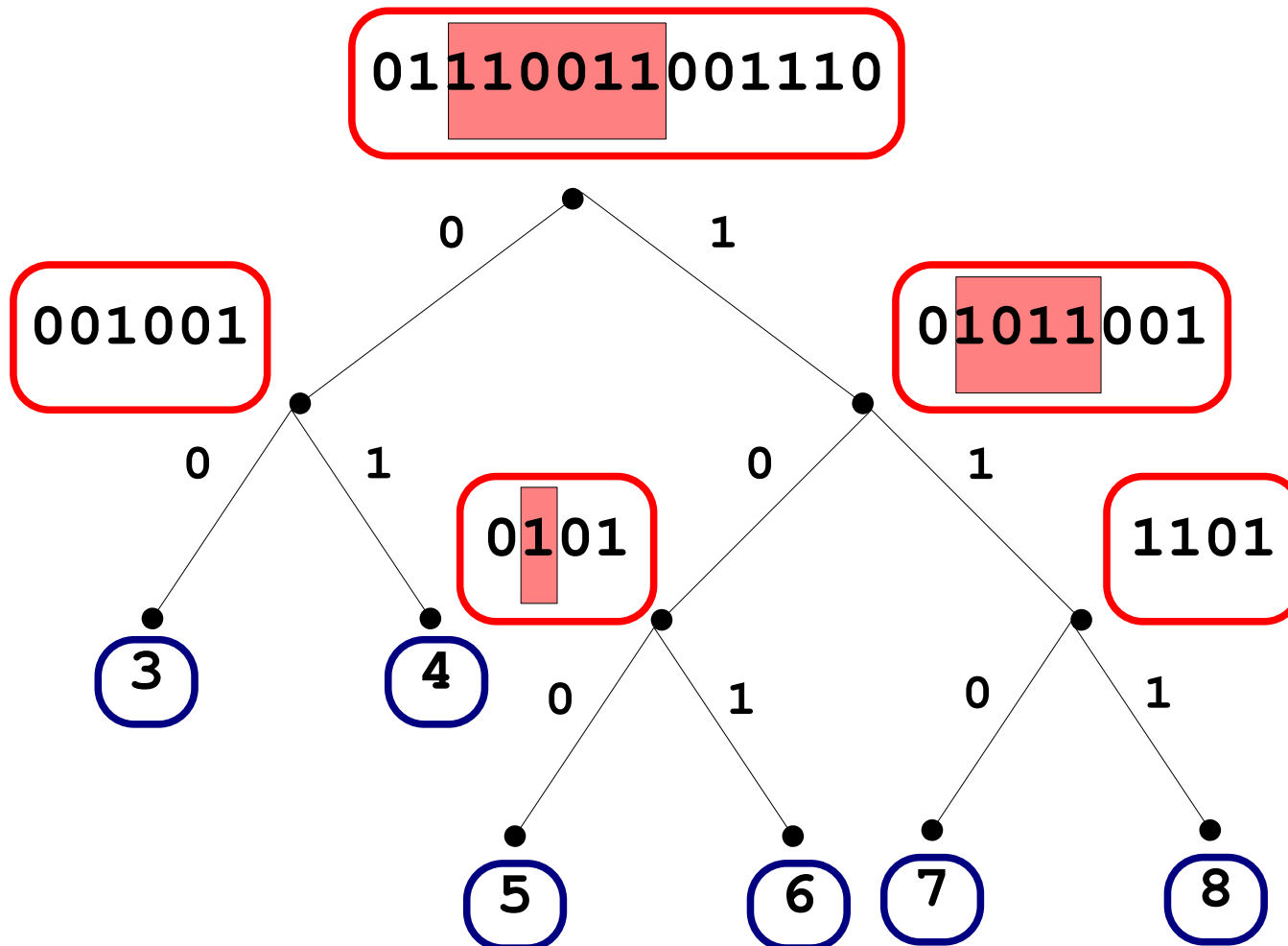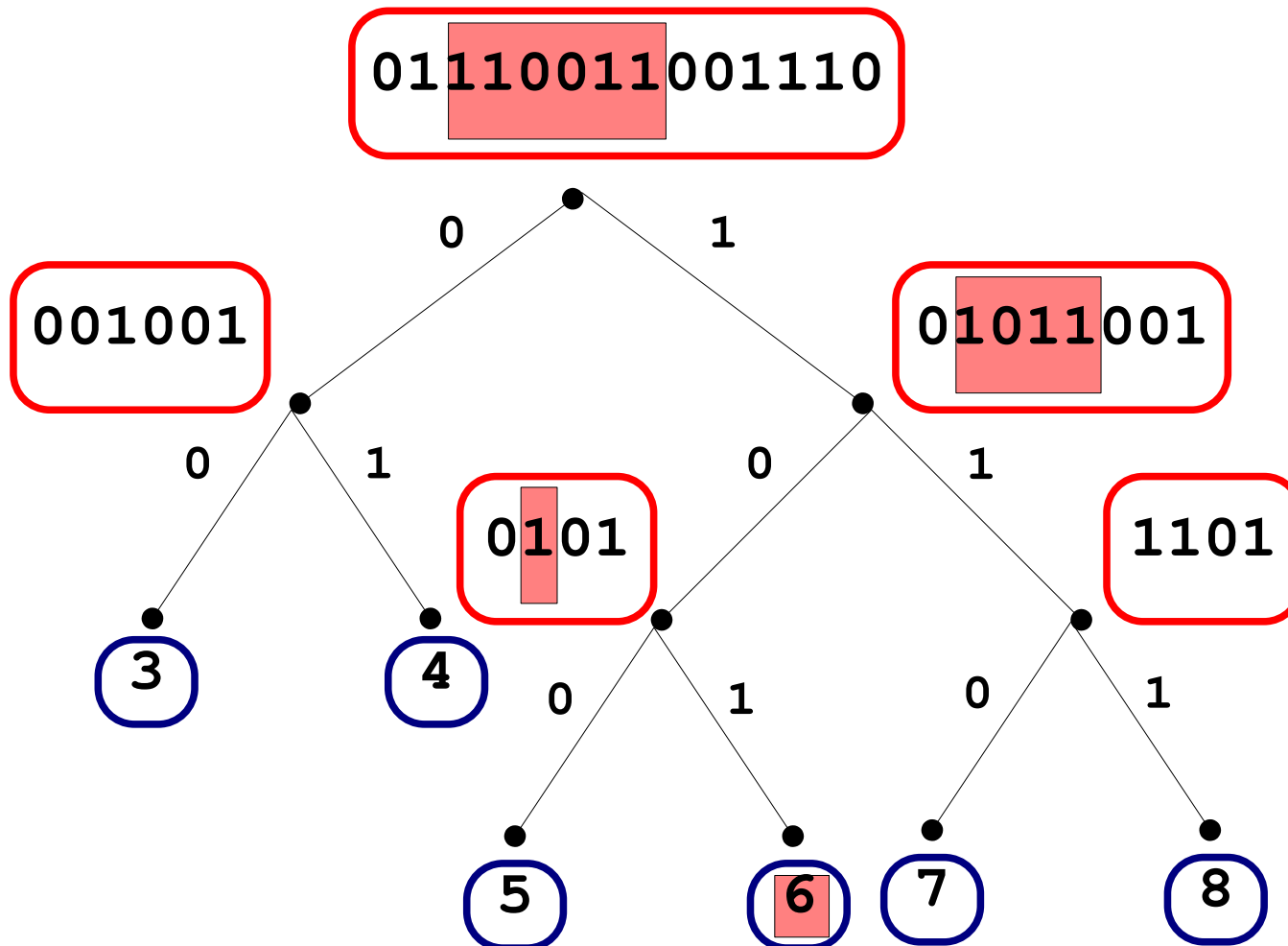
range_quantile(S,2,7,3): 3-rd smallest value in S[2⋯7]

S = [3,5,**8,6,3,4,8,7**,3,3,5,6,8,4]

range_quantile(S,2,7,3): 3-rd smallest value in S[2···7]

S = [3,5,8,6,3,4,8,7,3,3,5,6,8,4]

# Other applications

- Inverted indices (representation of posting lists)

- Computations Geometry (bidimensional range queries)

- Representation of Graphs, Permutations, ...

- Bioinformatics (maximal repeats)

# Implementation

- Conceptually simple, but tricky in a few details

- Several implementations cited in the literature, some available on the web

- Ready to use, open source, modular library:

  https://github.com/simongog/sdsl-lite

# References

- R. Grossi, A. Gupta, J. Vitter High-order entropy-compressed text indexes. Proc. SODA 2003 (paper introducing WTs: somewhat hard to read)

- G. Navarro, V. Mäkinen. Compressed Full-text indexes. ACM Comp. Surv. 2007 (main reference for WT and compressed indices)

- G. Navarro. Wavelet Trees for All, Journal Discrete Algorithms, 2014 (recent review paper on WTs)

# Thank you!