

Wavelet trees

Ben Langmead



JOHNS HOPKINS

WHITING SCHOOL
of ENGINEERING

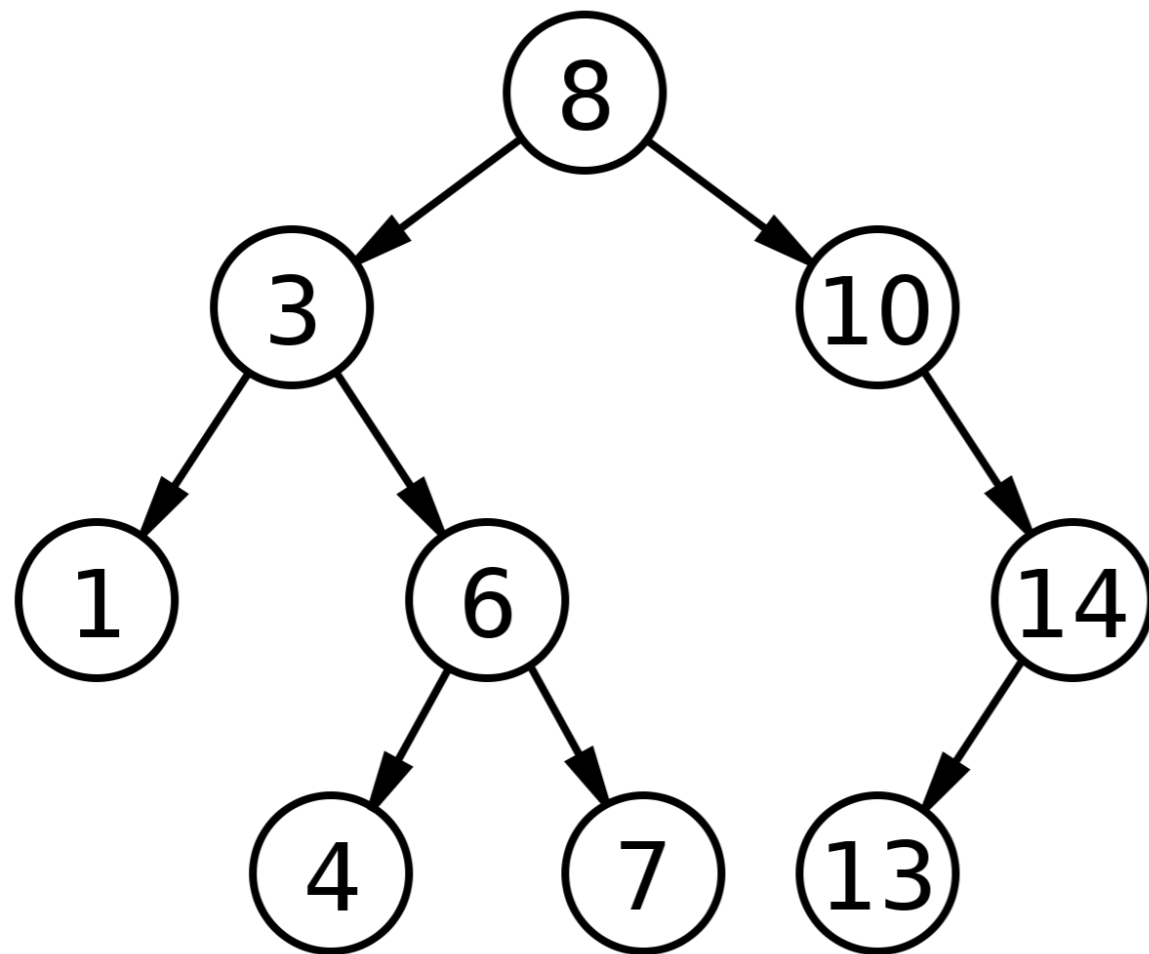
Department of Computer Science



Please sign guestbook (www.langmead-lab.org/teaching-materials) to tell me briefly how you are using the slides. For original Keynote files, email me (ben.langmead@gmail.com).

Trees

We're used to binary trees that repeatedly partition "value space"



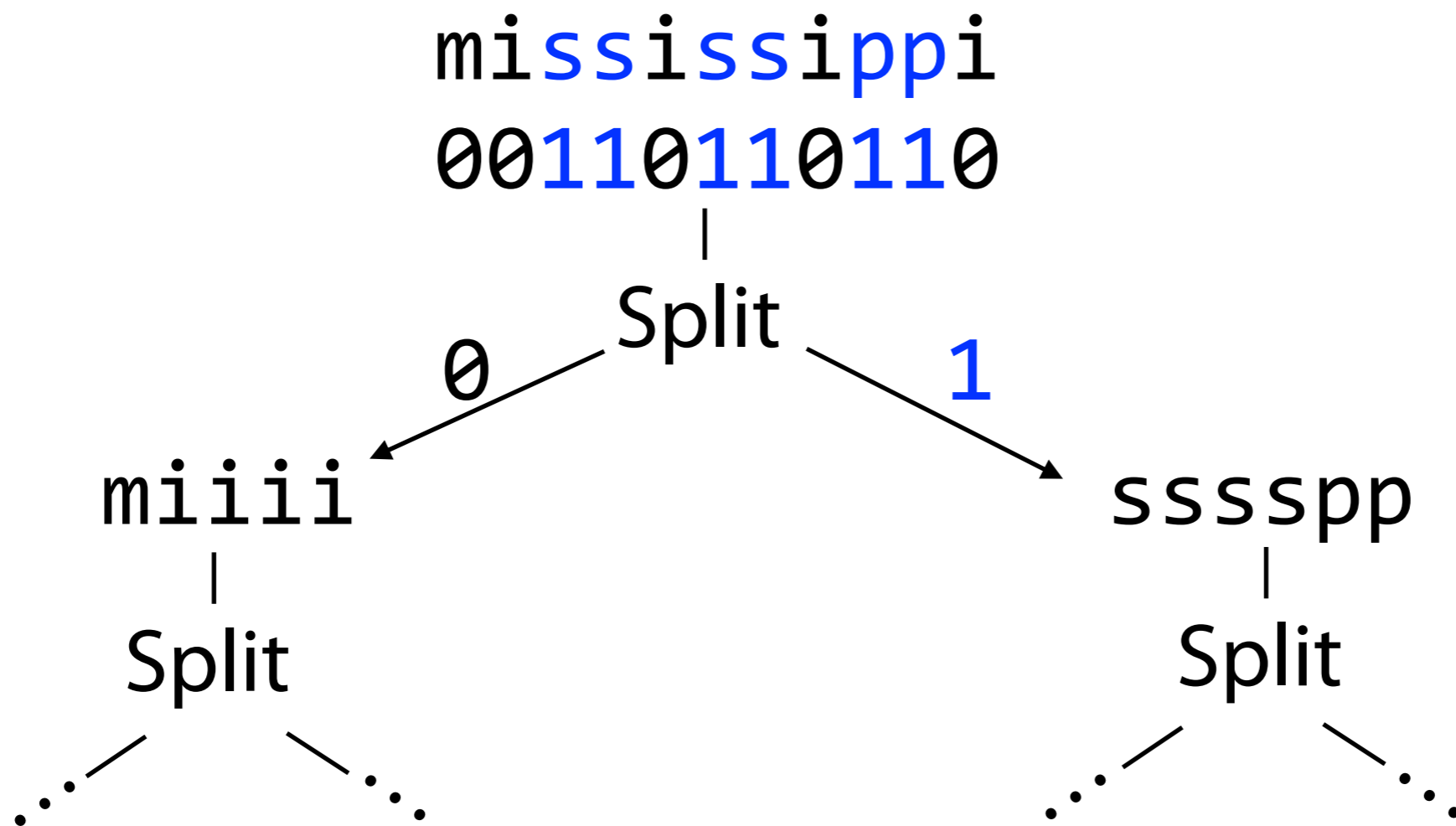
Rank and select are about alphabet space; where are the 0s and 1s? Where are the a's, c's, t's and g's?

Idea: partition *alphabet* space

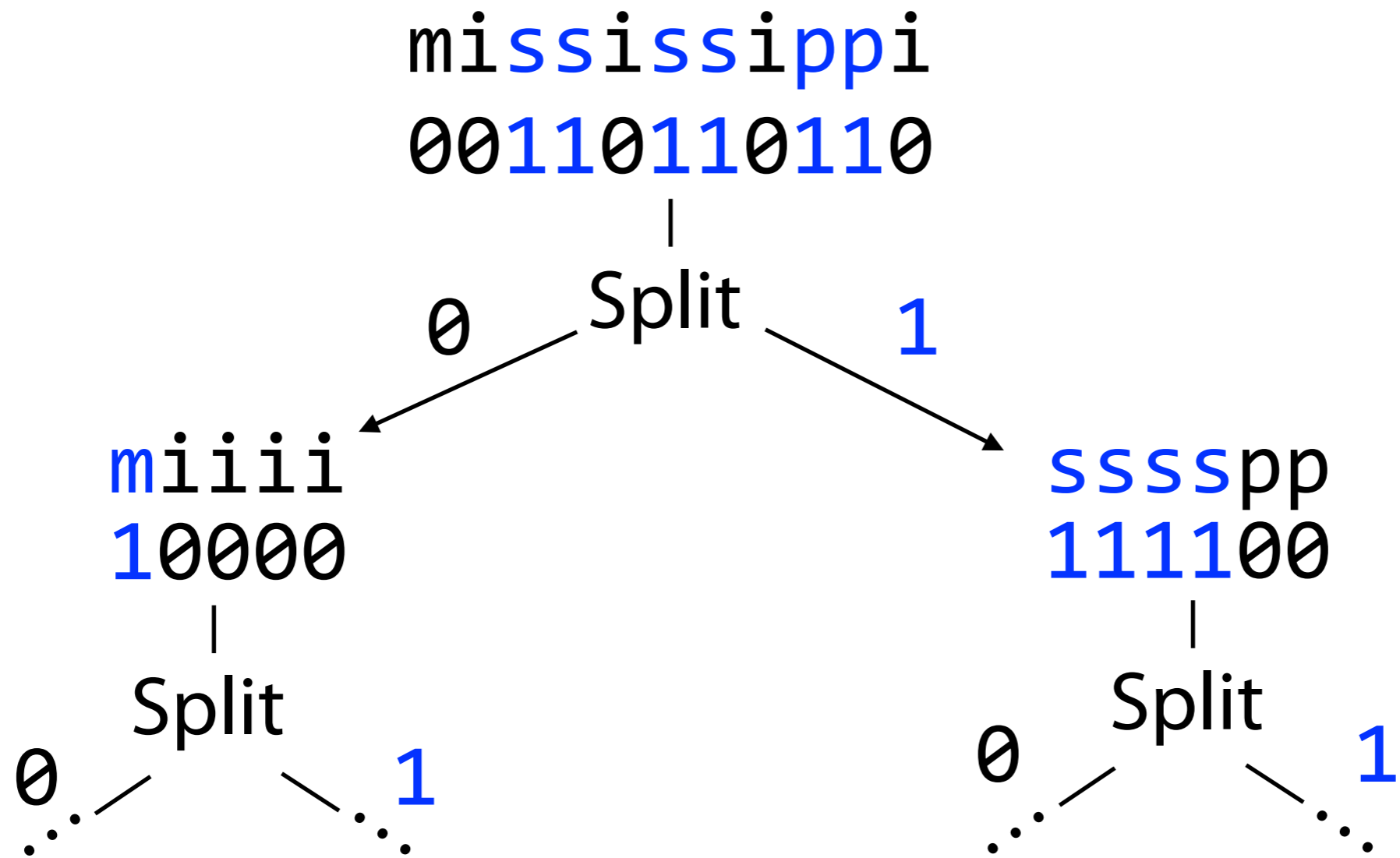
Wavelet trees

mississippi

Partition alphabet into $\{i, m\}$, $\{p, s\}$

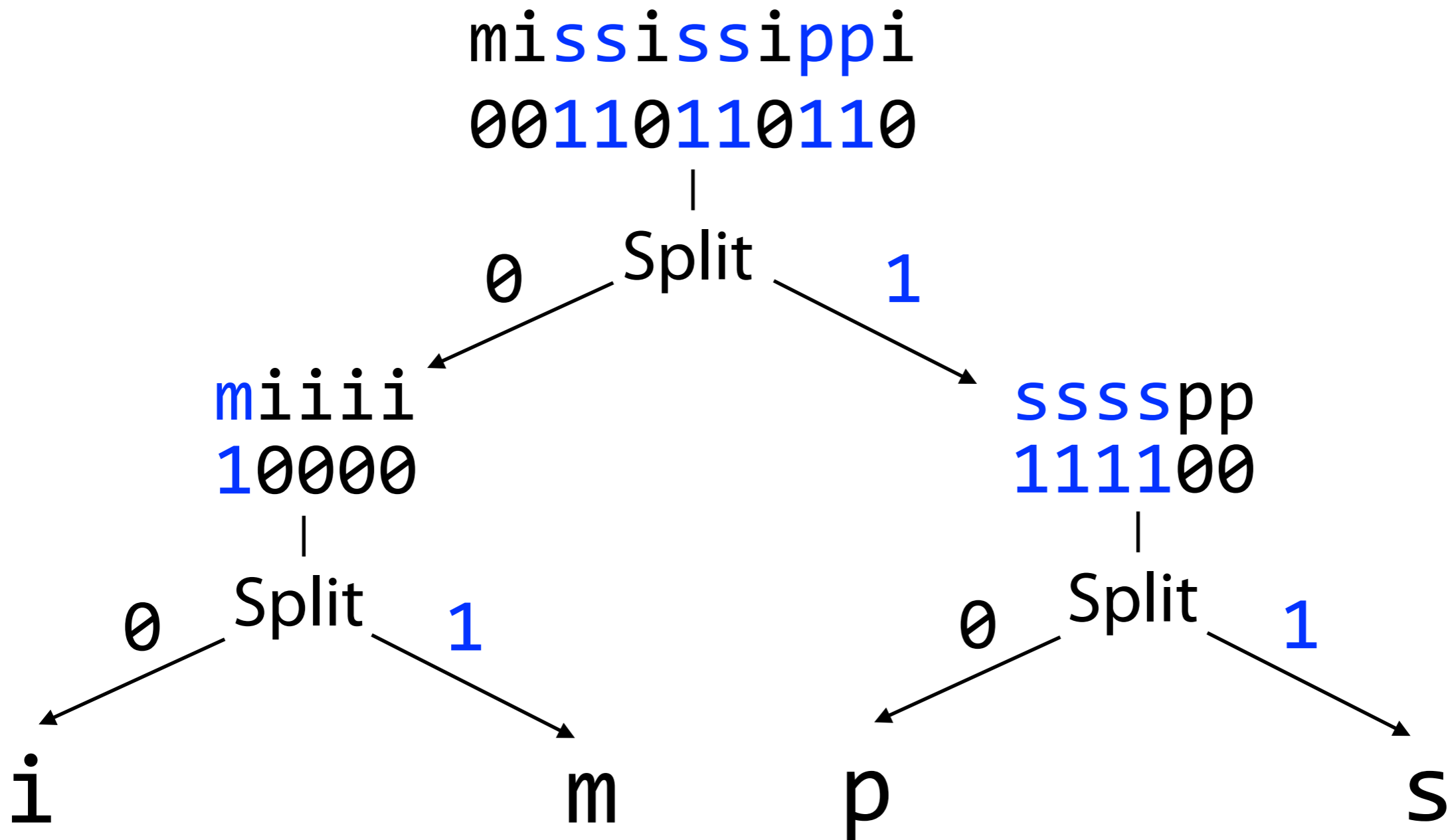


Wavelet trees

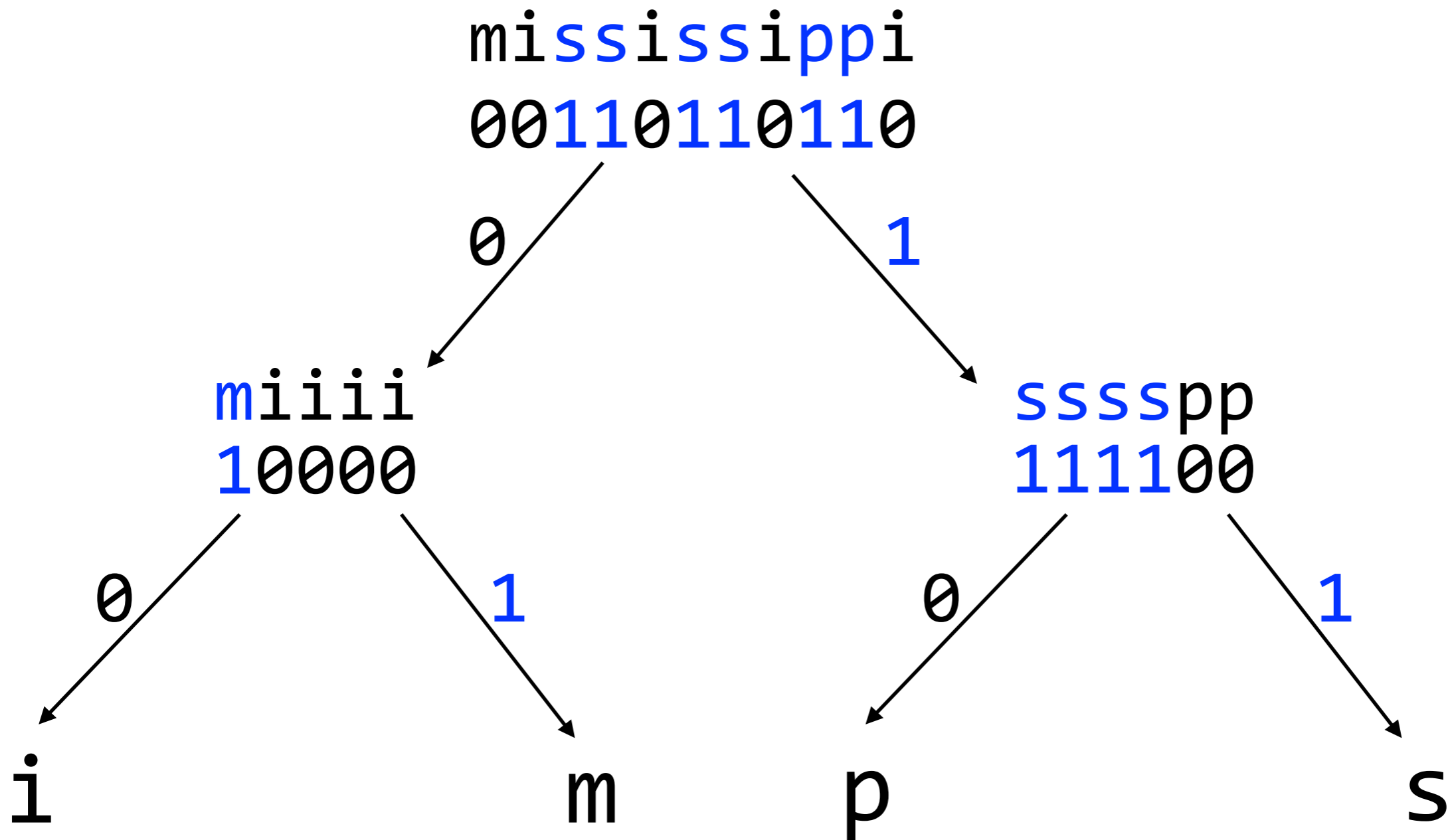


What goes in this next layer?

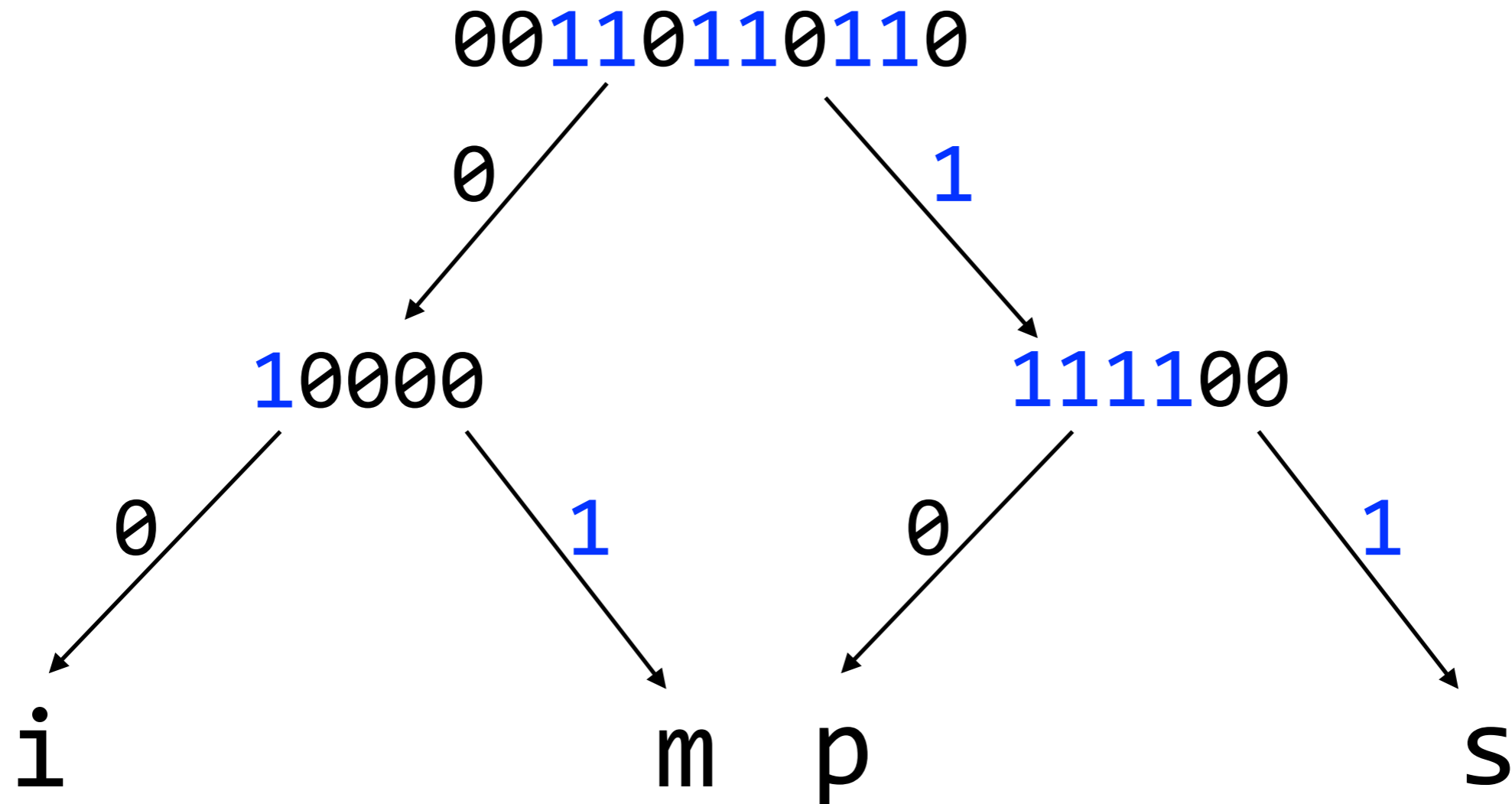
Wavelet trees



Wavelet trees



Wavelet trees



Can we do full-alphabet versions of access, rank and select?

How big is this?

Wavelet trees

RSA queries extend naturally to strings:

$$S . \text{access}(i) = S[i]$$

$$S . \text{rank}_c(i) = \sum_{j=0}^{i-1} \begin{cases} 1 & \text{if } S[j] = c \\ 0 & \text{otherwise} \end{cases}$$

$$S . \text{select}_c(i) = \max \{ j \mid S . \text{rank}_c(j) = i \}$$

Where S is a string with alphabet Σ and $c \in \Sigma$

Wavelet trees

$$S . \text{access}(i) = S[i]$$

$$S . \text{rank}_c(i) = \sum_{j=0}^{i-1} \begin{cases} 1 & \text{if } S[j] = c \\ 0 & \text{otherwise} \end{cases}$$

$$S . \text{select}_c(i) = \max \{ j \mid S . \text{rank}_c(j) = i \}$$

Wavelet trees

$S . \text{access}(4)$



$00110110110 . \text{access}(4) = 0$

$\{i, m\}$ 0

Descend to
left child

1 $\{p, s\}$

10000

111100

$\{i\}$ 0

1 $\{m\}$

0

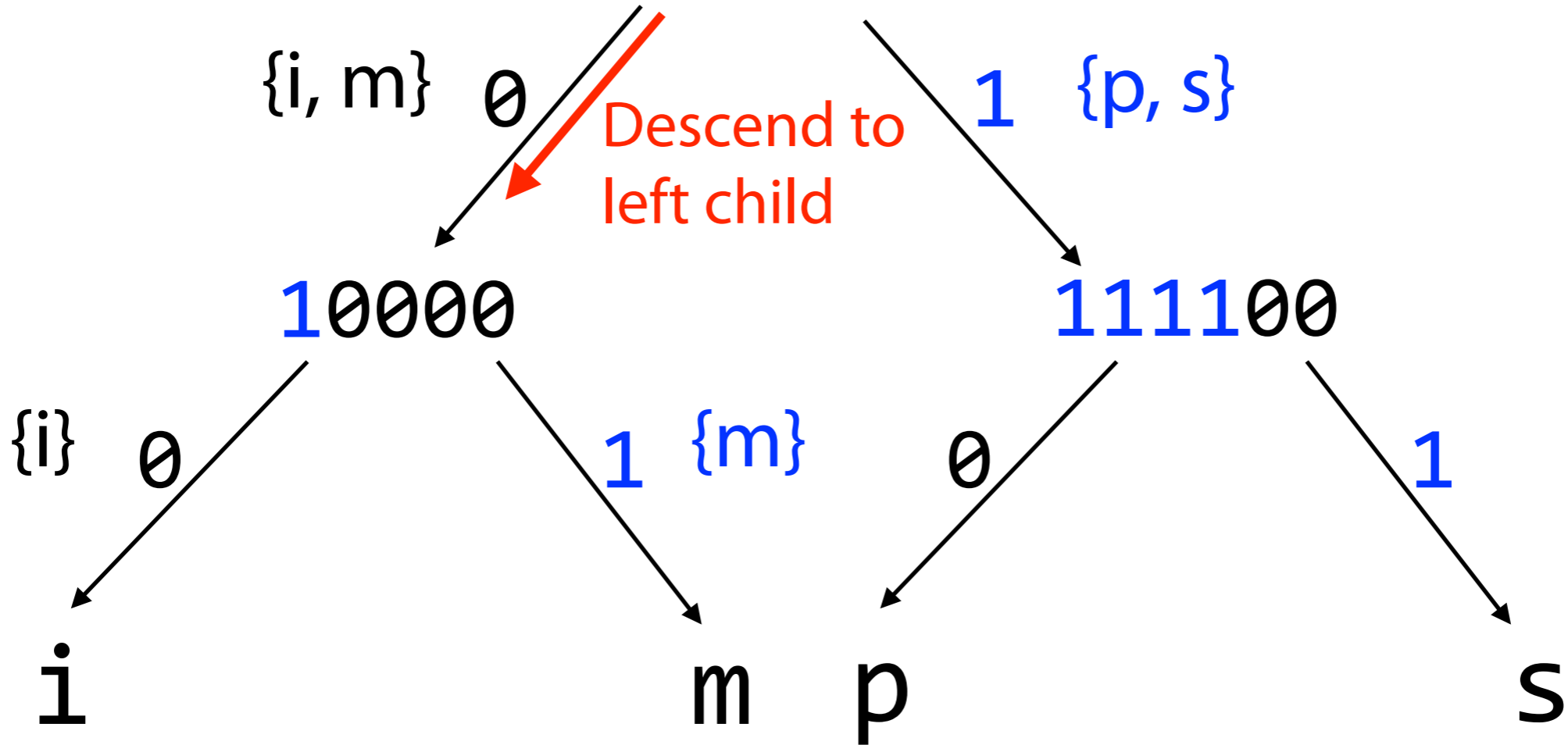
1

i

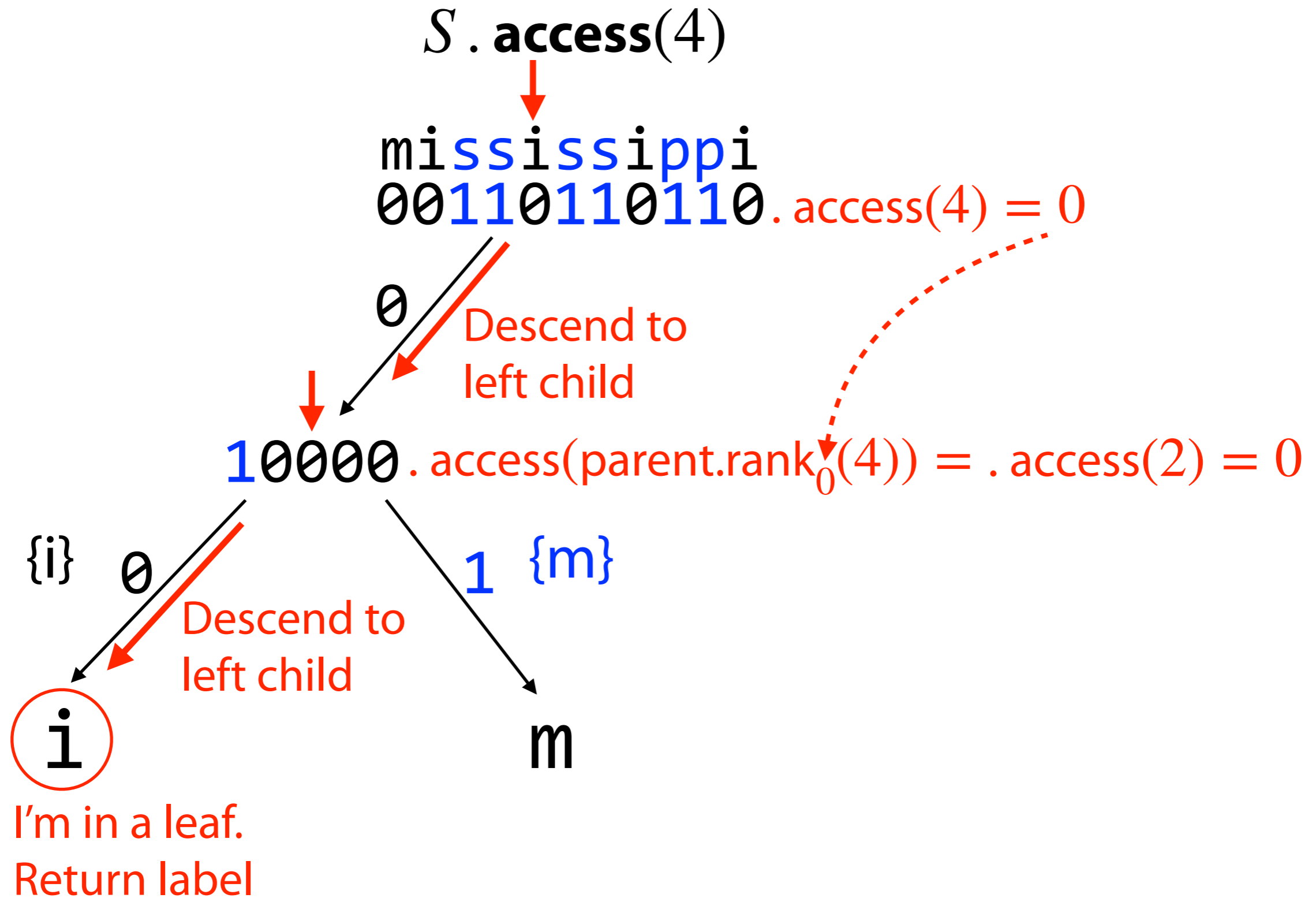
m

p

s



Wavelet trees



Wavelet trees

Wavelet tree access(i):

Given offset i :

$N \leftarrow \text{root}$

while N is not leaf

$B \leftarrow N.\text{bitvector}$

$b \leftarrow B[i]$

$N \leftarrow N.\text{child}(b)$

$i \leftarrow B.\text{rank}_b(i)$

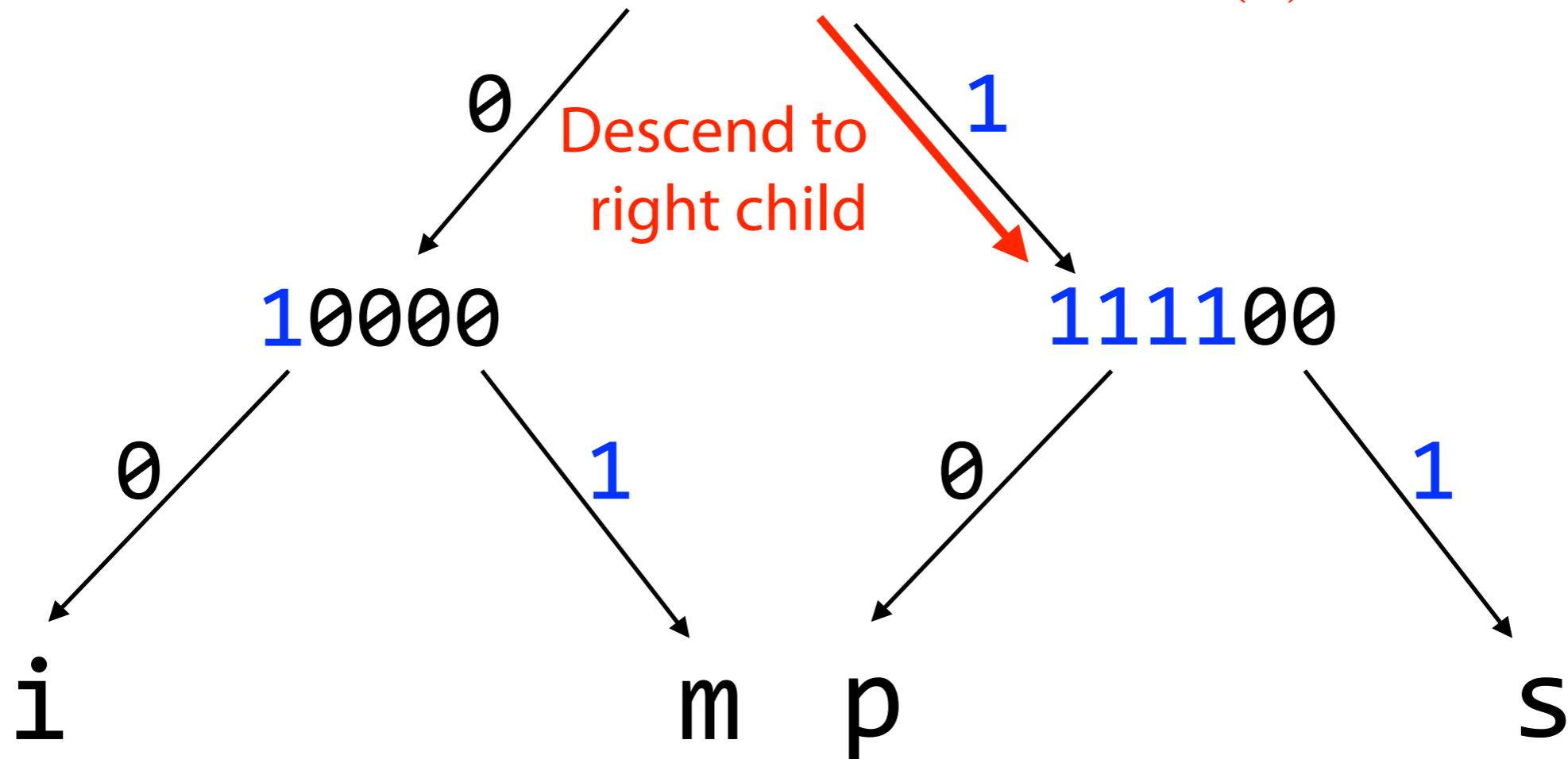
return $N.\text{label}$

Wavelet trees

$S . \text{access}(6)$



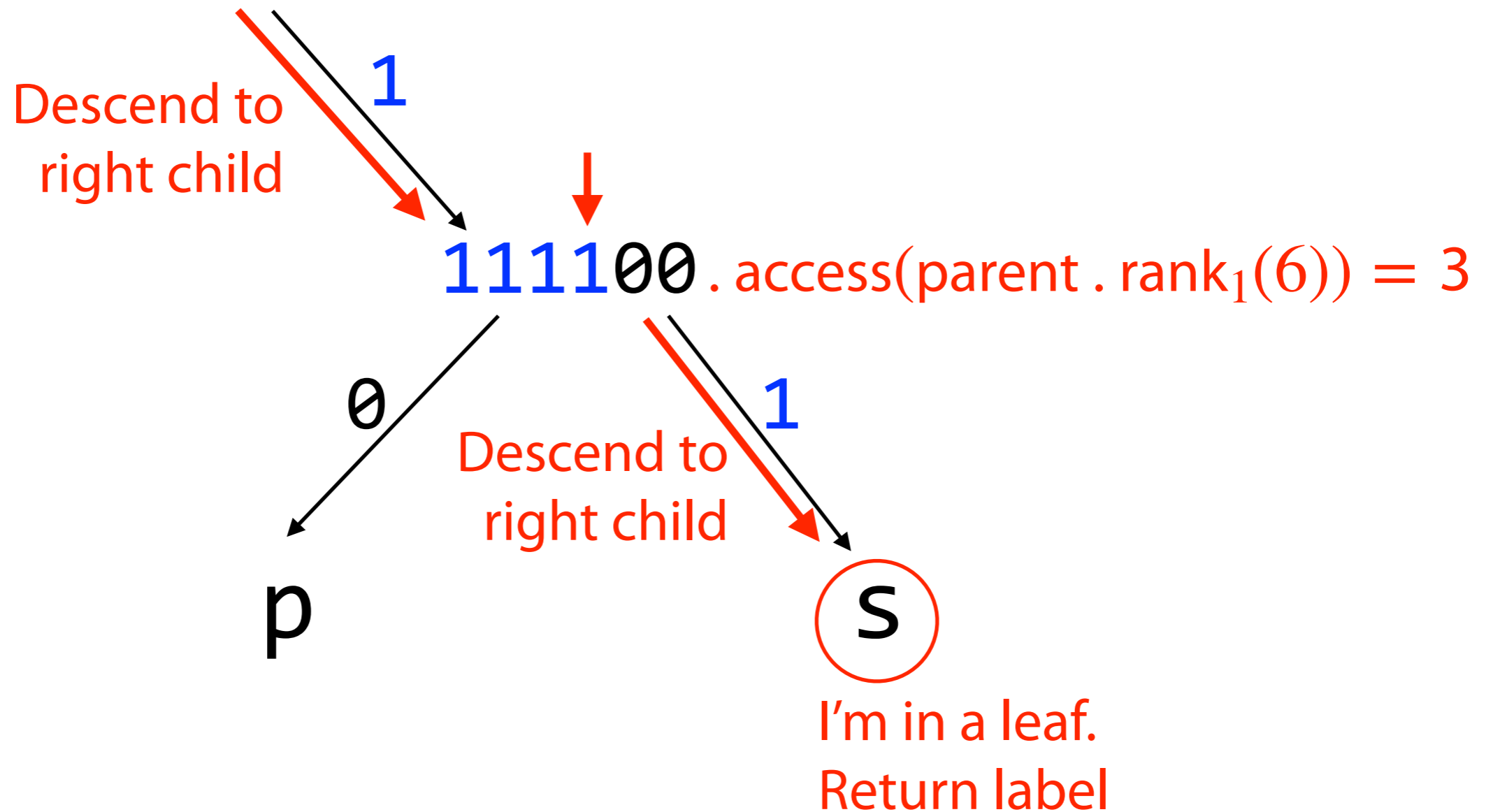
00**110110110**.access(6) = 1



Wavelet trees

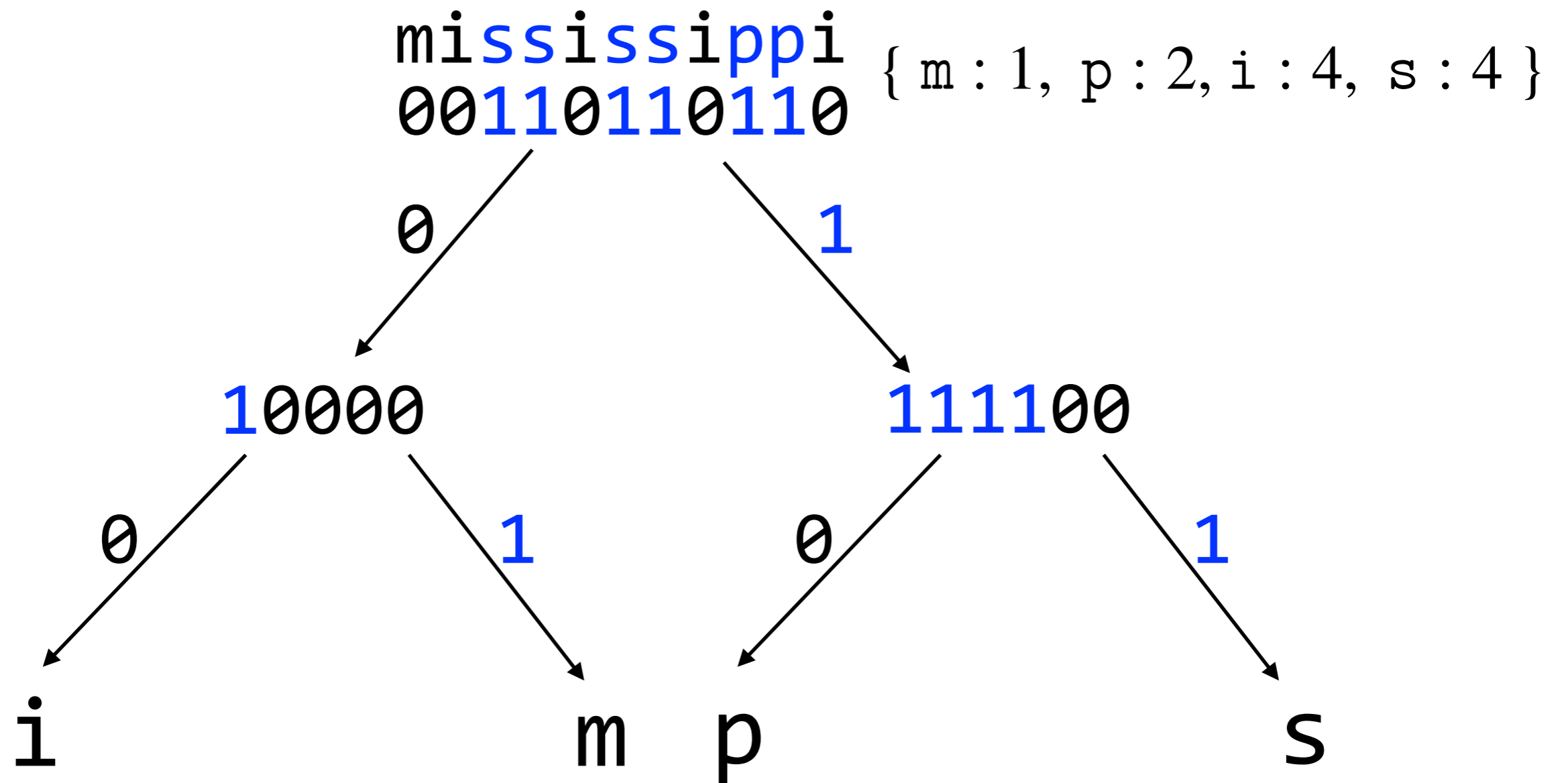
S . **access**(6)

mississippi
00110110110. **access**(6) = 1



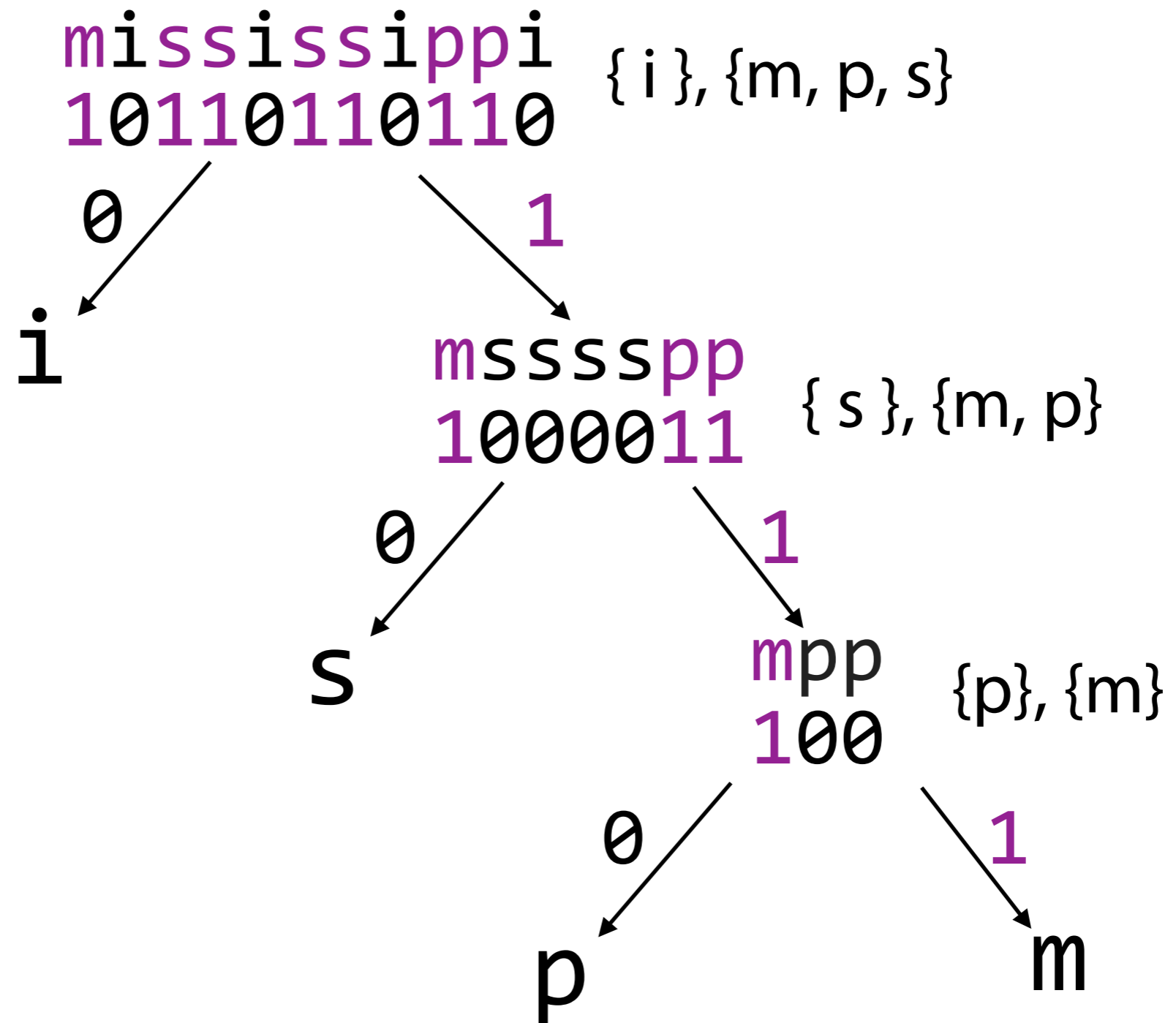
Wavelet trees

Could have picked a different shape for the tree



Wavelet trees

Could have picked a different shape for the tree



Wavelet trees

Tree shape
defines a (prefix)
code

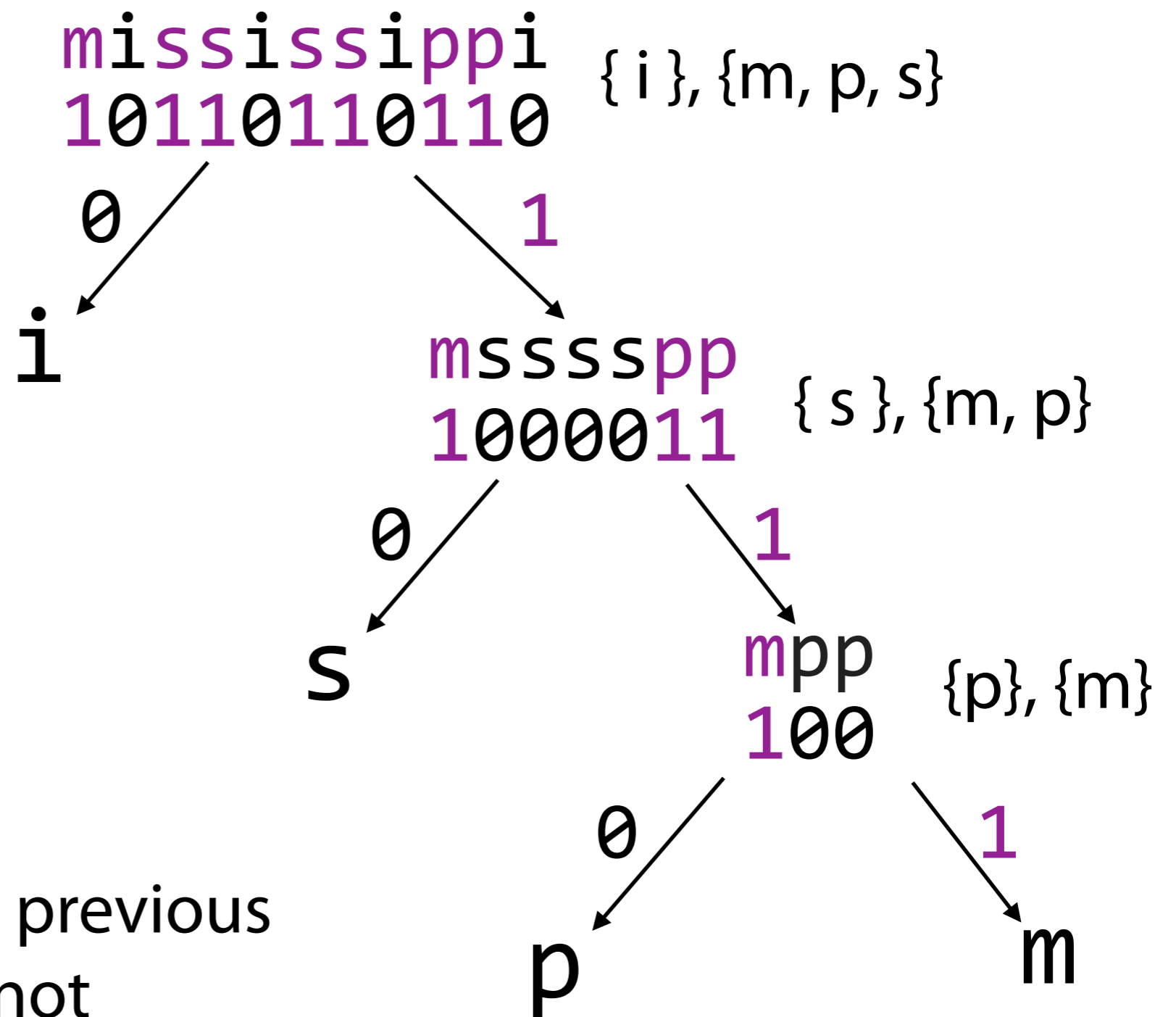
$$C(i) = 0$$

$$C(s) = 10$$

$$C(p) = 110$$

$$C(m) = 111$$

This tree is Huffman; previous
(balanced) tree was not



Wavelet trees

$$S . \text{access}(i) = S[i]$$

$$S . \text{rank}_c(i) = \sum_{j=0}^{i-1} \begin{cases} 1 & \text{if } S[j] = c \\ 0 & \text{otherwise} \end{cases}$$

$$S . \text{select}_c(i) = \max \{ j \mid S . \text{rank}_c(j) = i \}$$

Note that rank can ask about *any* character c at *any* position i

Wavelet trees

$$S.\text{rank}_c(i) = \sum_{j=0}^{i-1} \begin{cases} 1 & \text{if } S[j] = c \\ 0 & \text{otherwise} \end{cases}$$

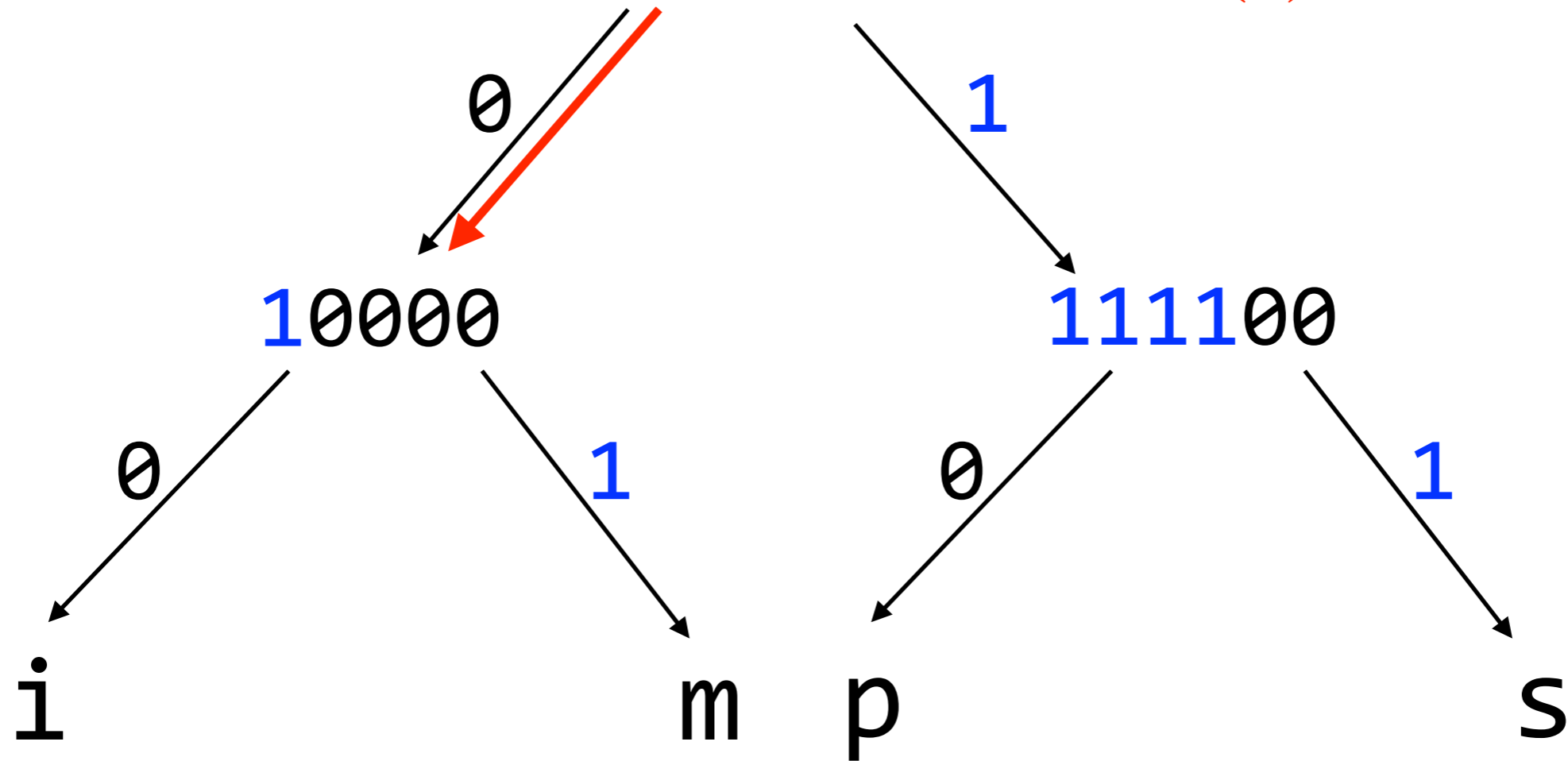
Algorithm will be similar to access...

Wavelet trees

$S . \text{rank}_i(7)$



00**11**0**11**0**11**0 . $\text{access}(7) = 0$



Wavelet trees

$S.\text{rank}_i(7)$



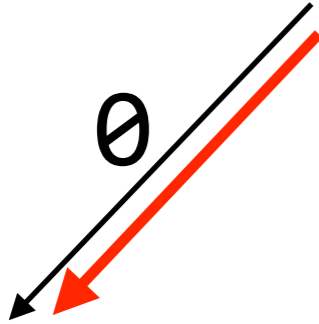
00**11**0**11**0**11**0 . $\text{access}(7) = 0$

0



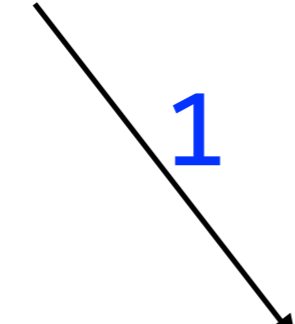
10000 . $\text{access}(\text{parent}.\text{rank}_0(7)) = .\text{access}(3) = 0$

0



i

1



m

I'm in a leaf.
Return label

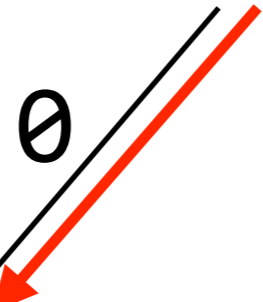
So....where's the answer?

Wavelet trees

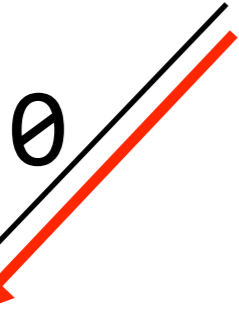
$S.\text{rank}_i(7)$



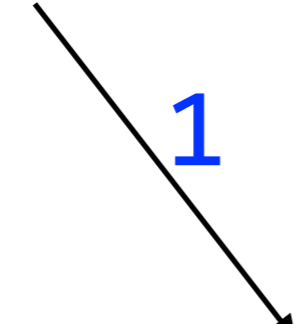
00**11**0**11**0**11**0 . $\text{access}(7) = 0$



10000 . $\text{access}(\text{parent}.\text{rank}_0(7)) = \text{access}(3) = 0$



i



m

I'm in a leaf.
Return label

Answer: $\text{parent}.\text{rank}_0(3) = 2$

Wavelet trees

$$S.\text{rank}_c(i) = \sum_{j=0}^{i-1} \begin{cases} 1 & \text{if } S[j] = c \\ 0 & \text{otherwise} \end{cases}$$

Algorithm will be similar to access...

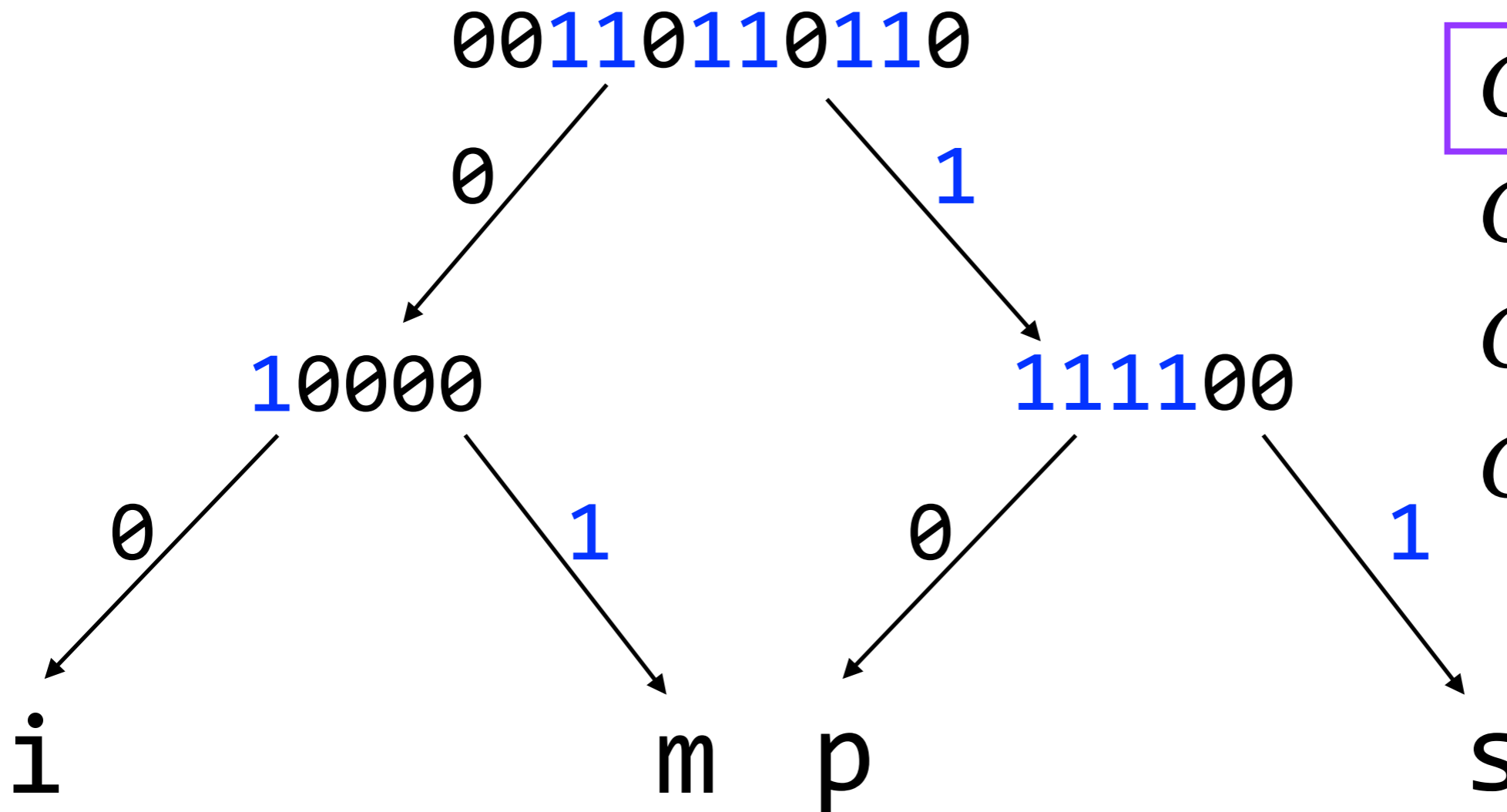
But the path we follow corresponds to c , which isn't necessarily the character at $S[i]$

Wavelet trees

$S.\text{rank}_i(6)$



Note: $S[6] = s \neq i$
mississ**s**ippi



$$C(i) = 00$$

$$C(m) = 01$$

$$C(p) = 10$$

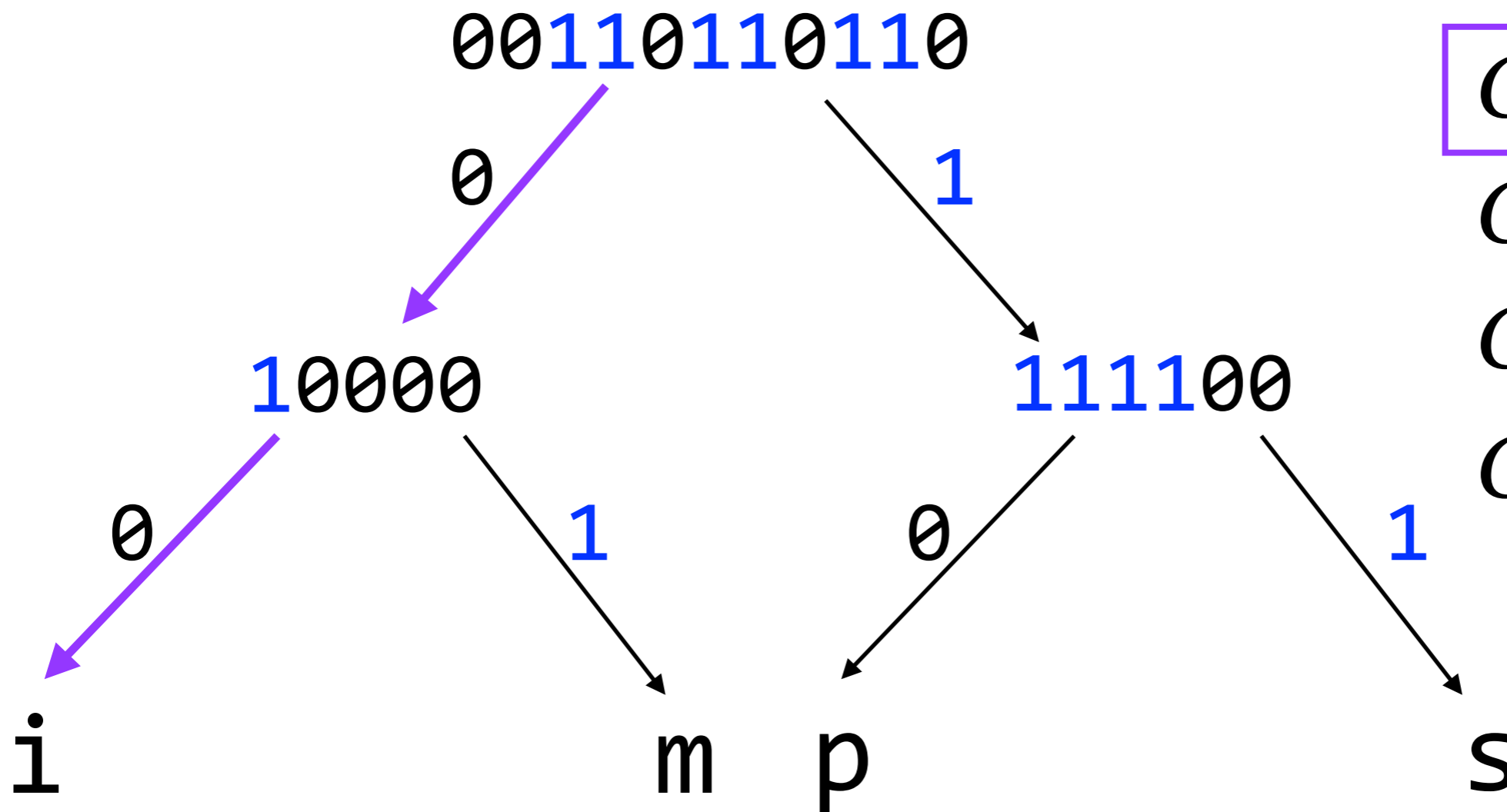
$$C(s) = 11$$

Wavelet trees

$S.\text{rank}_i(6)$



Note: $S[6] = s \neq i$
mississippi



$$C(i) = 00$$

$$C(m) = 01$$

$$C(p) = 10$$

$$C(s) = 11$$

Wavelet trees

$S . \text{rank}_i(6)$

$$C(i) = 00$$

$00110110110 . \text{rank}_0(6) = 3$

\emptyset

$100000 . \text{rank}_0(3) = 2$

Answer: 2

\emptyset

i

Movements through tree & subscripts of rank queries come from code C

Wavelet trees

Wavelet tree $\text{rank}_x(i)$:

Given character x and offset i :

$N \leftarrow \text{root}$

$k \leftarrow 0$

while N is not leaf

$B \leftarrow N.\text{bitvector}$

$b \leftarrow c(x)[k]$

$i \leftarrow B.\text{rank}_b(i)$

$N \leftarrow N.\text{child}(b)$

$k \leftarrow k + 1$

return i

Wavelet trees

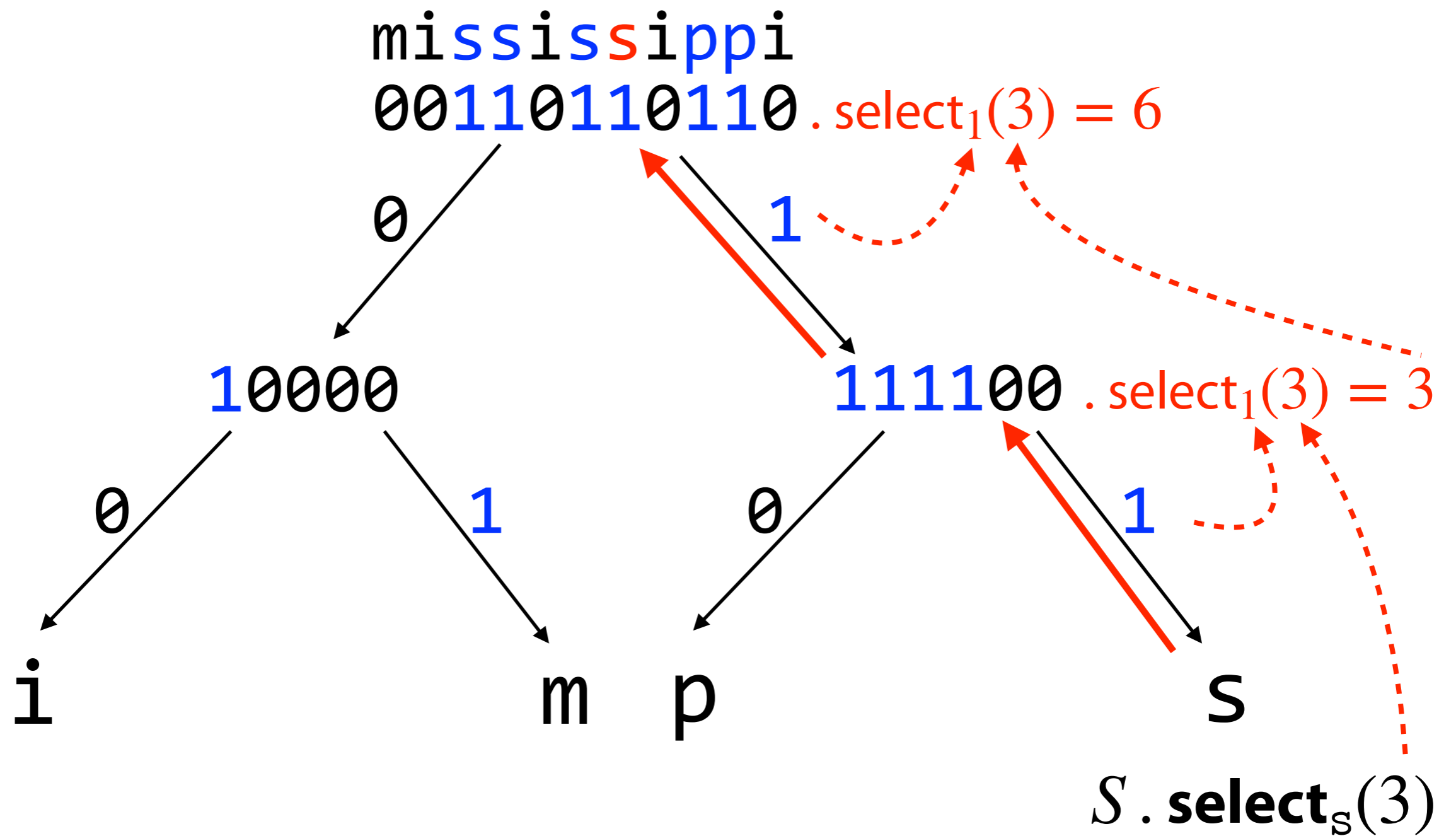
$$S . \text{access}(i) = S[i]$$

$$S . \text{rank}_c(i) = \sum_{j=0}^{i-1} \begin{cases} 1 & \text{if } S[j] = c \\ 0 & \text{otherwise} \end{cases}$$

$$S . \text{select}_c(i) = \max \{ j \mid S . \text{rank}_c(j) = i \}$$

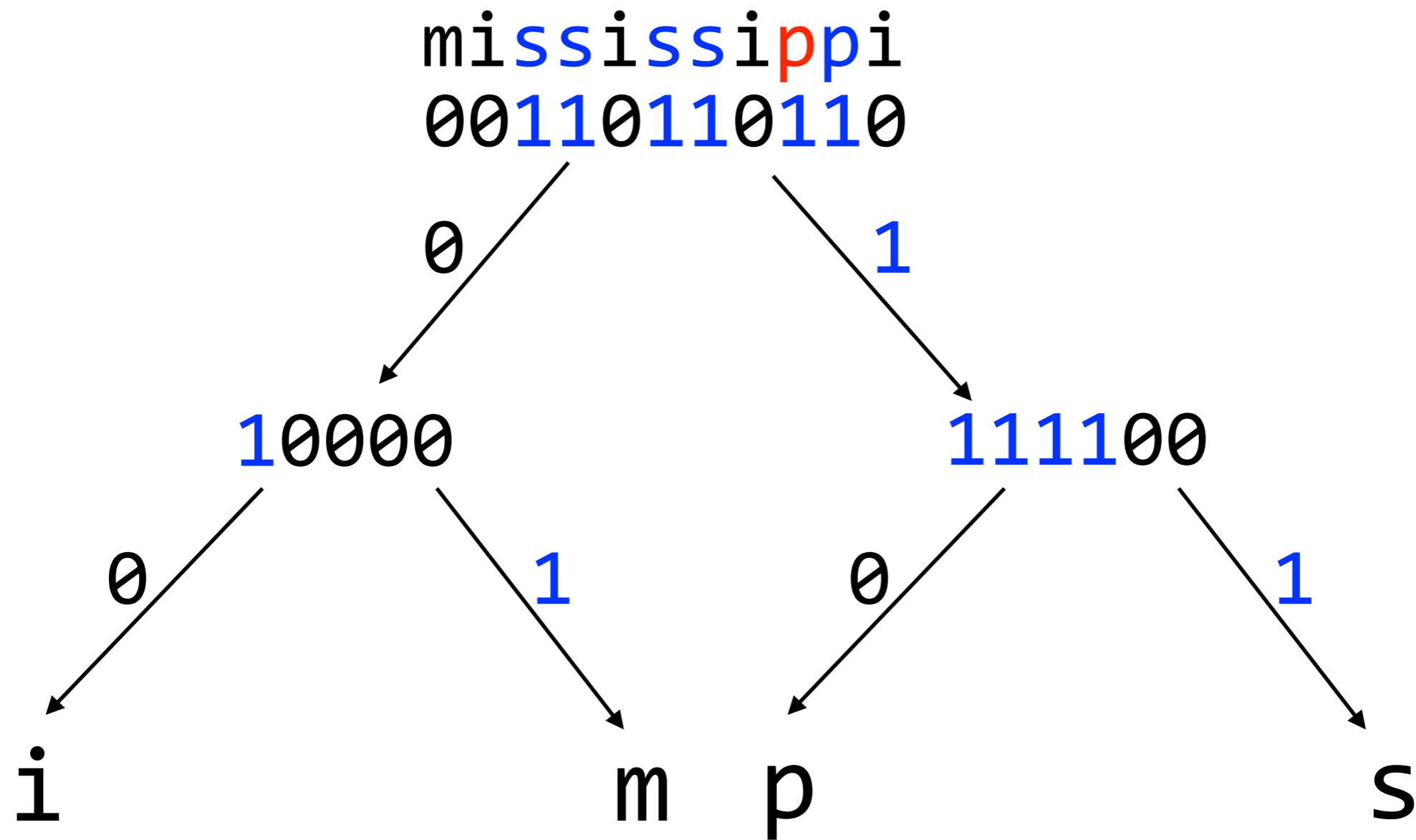
Wavelet trees

Answer: 6



Wavelet trees

Answer: 8



$S . \text{select}_p(0)$

Wavelet trees

Wavelet tree $\text{select}_x(i)$:

Given character x and rank i :

$N \leftarrow \text{leaf}(x)$

$l \leftarrow |c(x)| - 1$

while N is not root

$N \leftarrow N.\text{parent}()$

$B \leftarrow N.\text{bitvector}$

$b \leftarrow c(x)[k]$

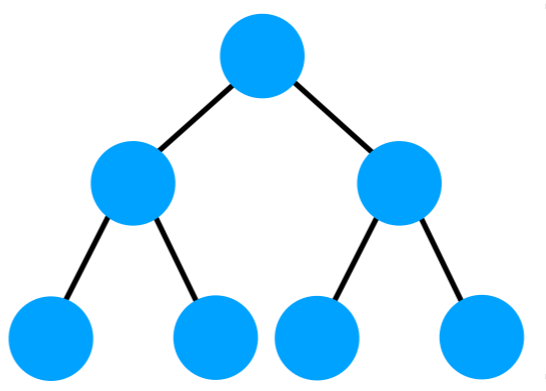
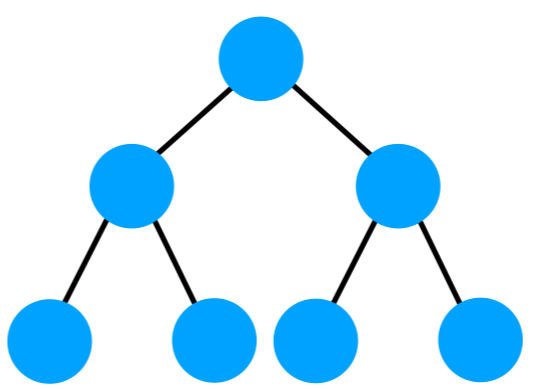
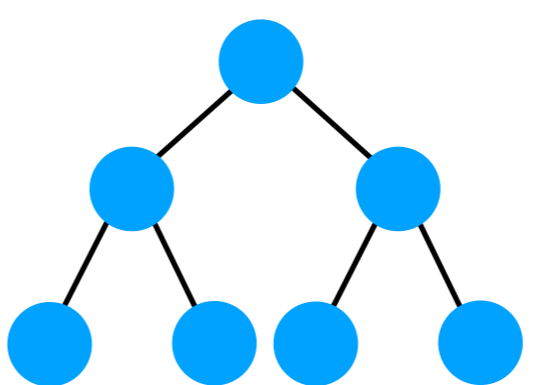
$i \leftarrow B.\text{select}_b(i)$

$k \leftarrow k - 1$

return i

Wavelet trees

Assuming ***balanced*** tree

	Time	Note
$S . \text{access}$	$O(\log \sigma)$	 $\log_2 \sigma$ Jacobson's rank is $O(1)$ time
$S . \text{rank}_c$	$O(\log \sigma)$	 $\log_2 \sigma$ Jacobson's rank is $O(1)$ time
$S . \text{select}_c$	$O(\log \sigma)$	 $\log_2 \sigma$ Clark's select is $O(1)$ time

Wavelet trees

Exercise: do similar analysis for Huffman-shaped tree, with results in terms of H_0

Exercise: space analysis, assuming bitvectors at internal nodes can be combined in a single level-wise bitvector