

# Wavelet trees

Ben Langmead



JOHNS HOPKINS

WHITING SCHOOL  
*of* ENGINEERING

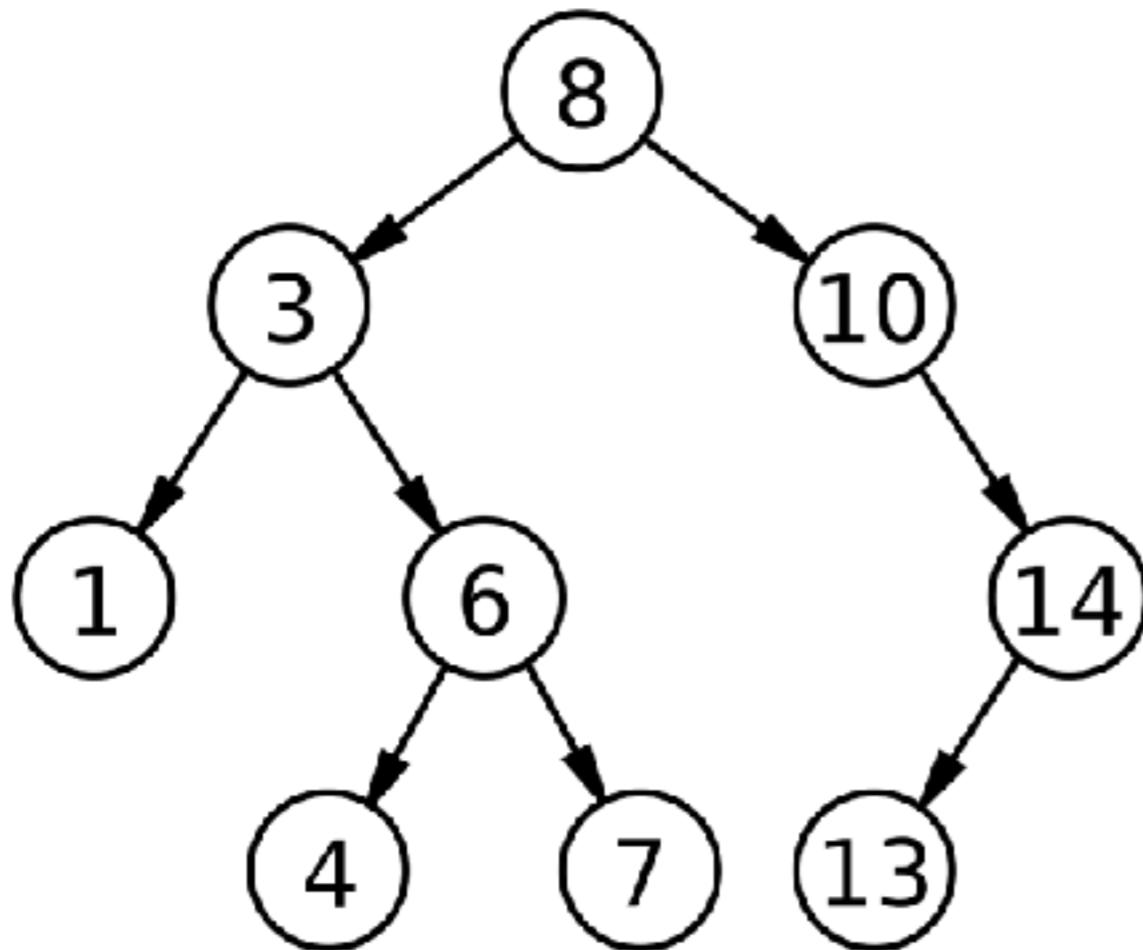
Department of Computer Science



Please sign guestbook ([www.langmead-lab.org/teaching-materials](http://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](mailto:ben.langmead@gmail.com)).

# Trees

We know trees that partition “value space”



Different idea:  
partition  
*alphabet* space

# Wavelet trees

mississippi

Partition alphabet into

# Wavelet trees

mississippi

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

mississippi

# Wavelet trees

mississippi

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

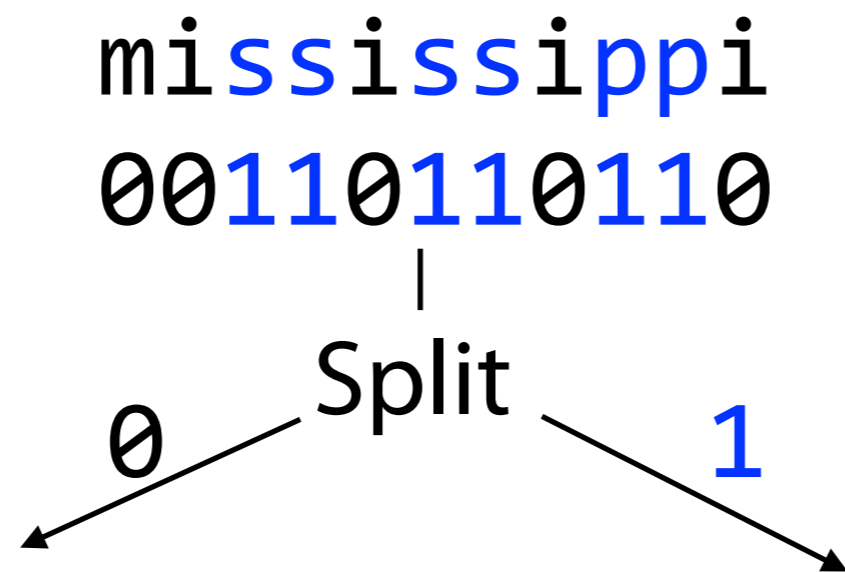
mississippi

00110110110

# Wavelet trees

mississippi

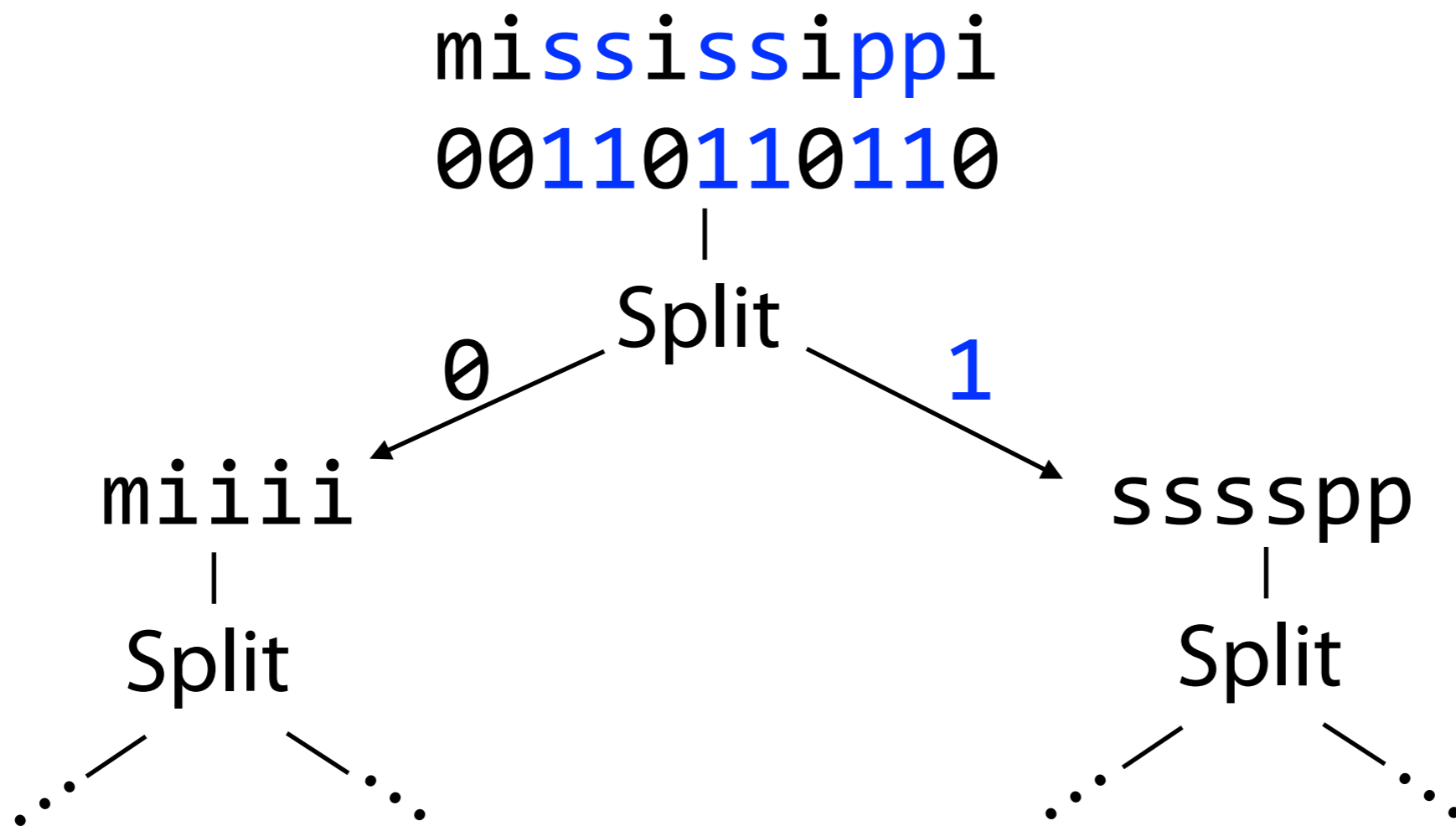
Partition alphabet into  $\{i, m\}$ ,  $\{p, s\}$



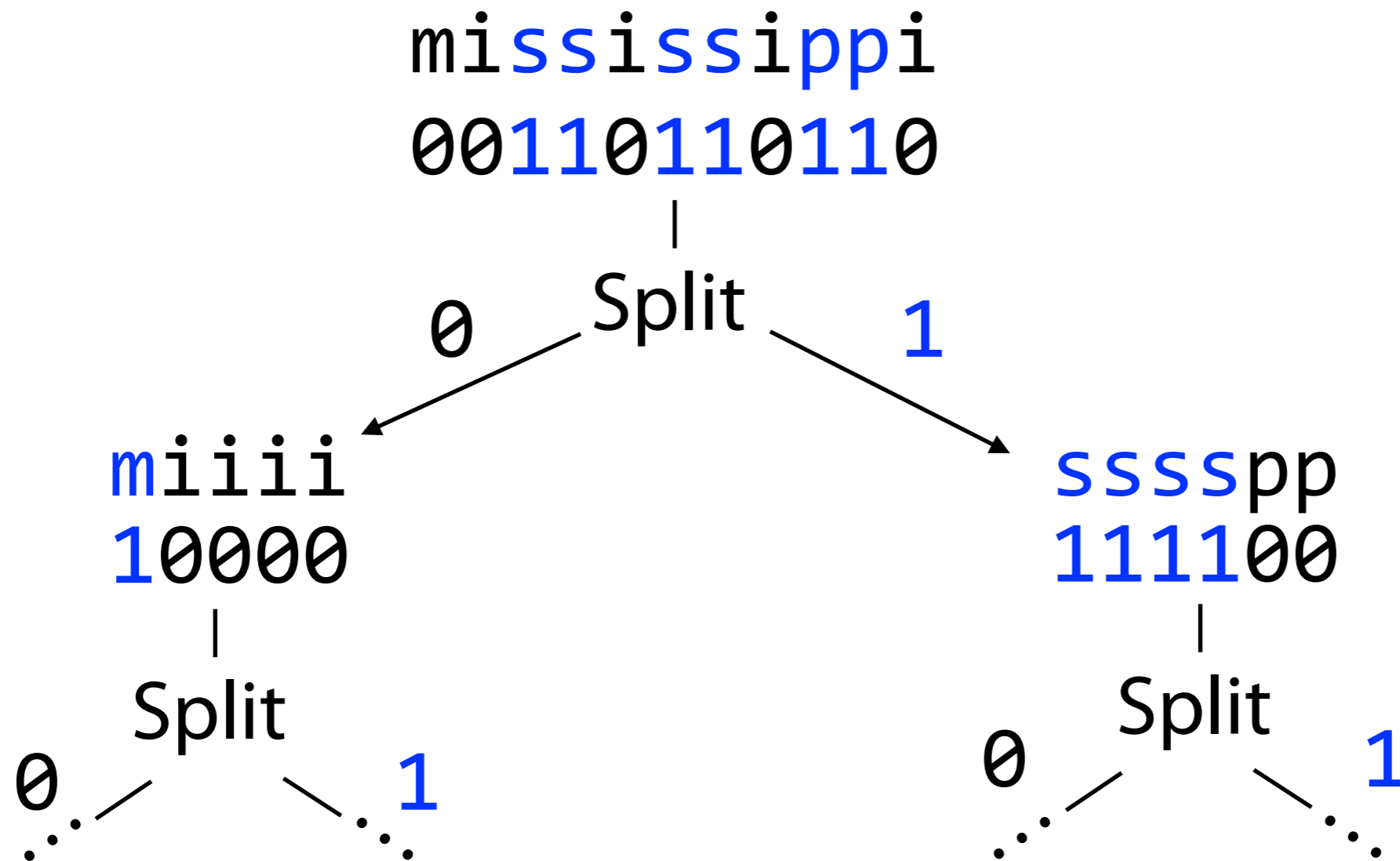
# Wavelet trees

mississippi

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



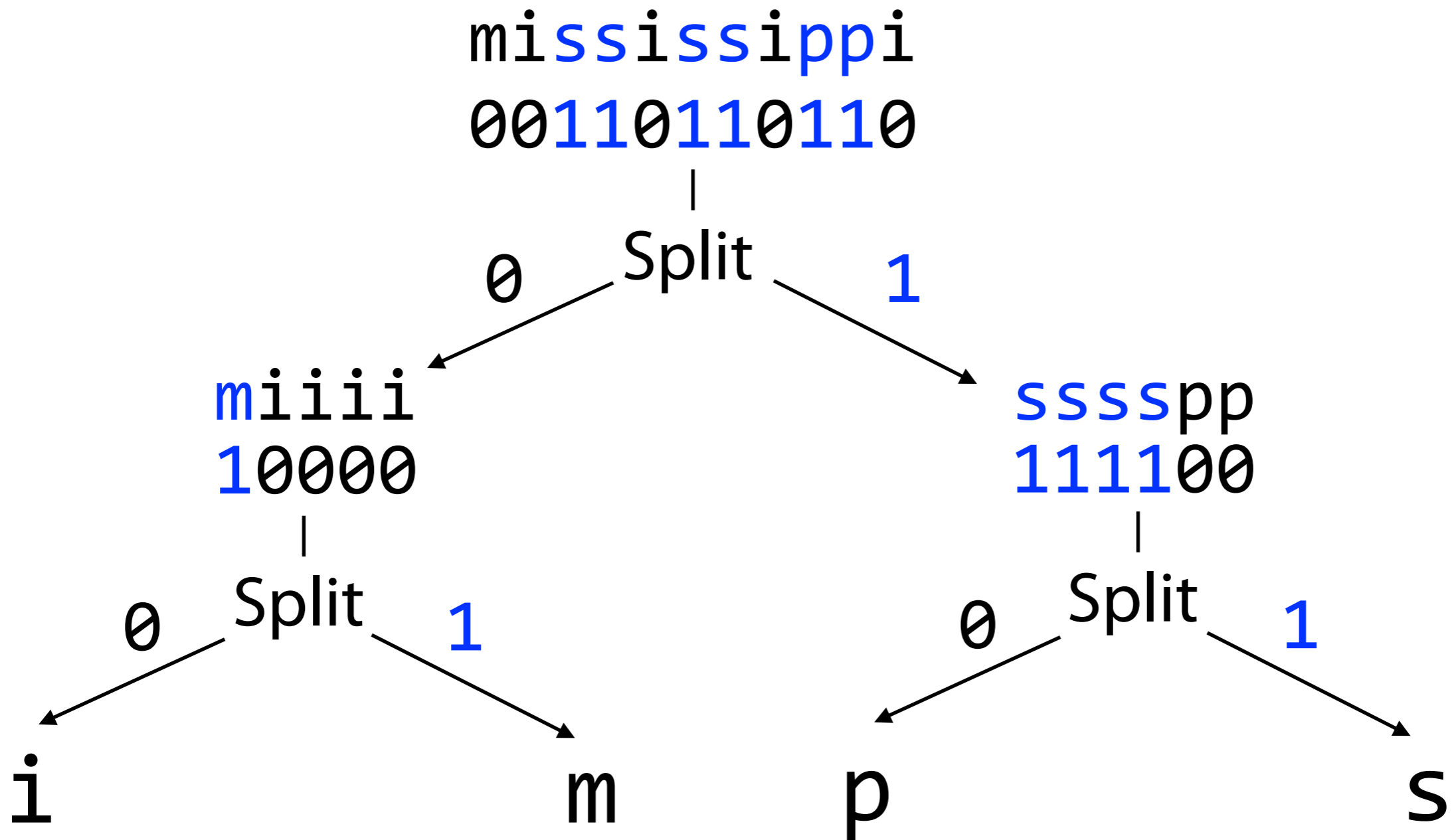
# Wavelet trees



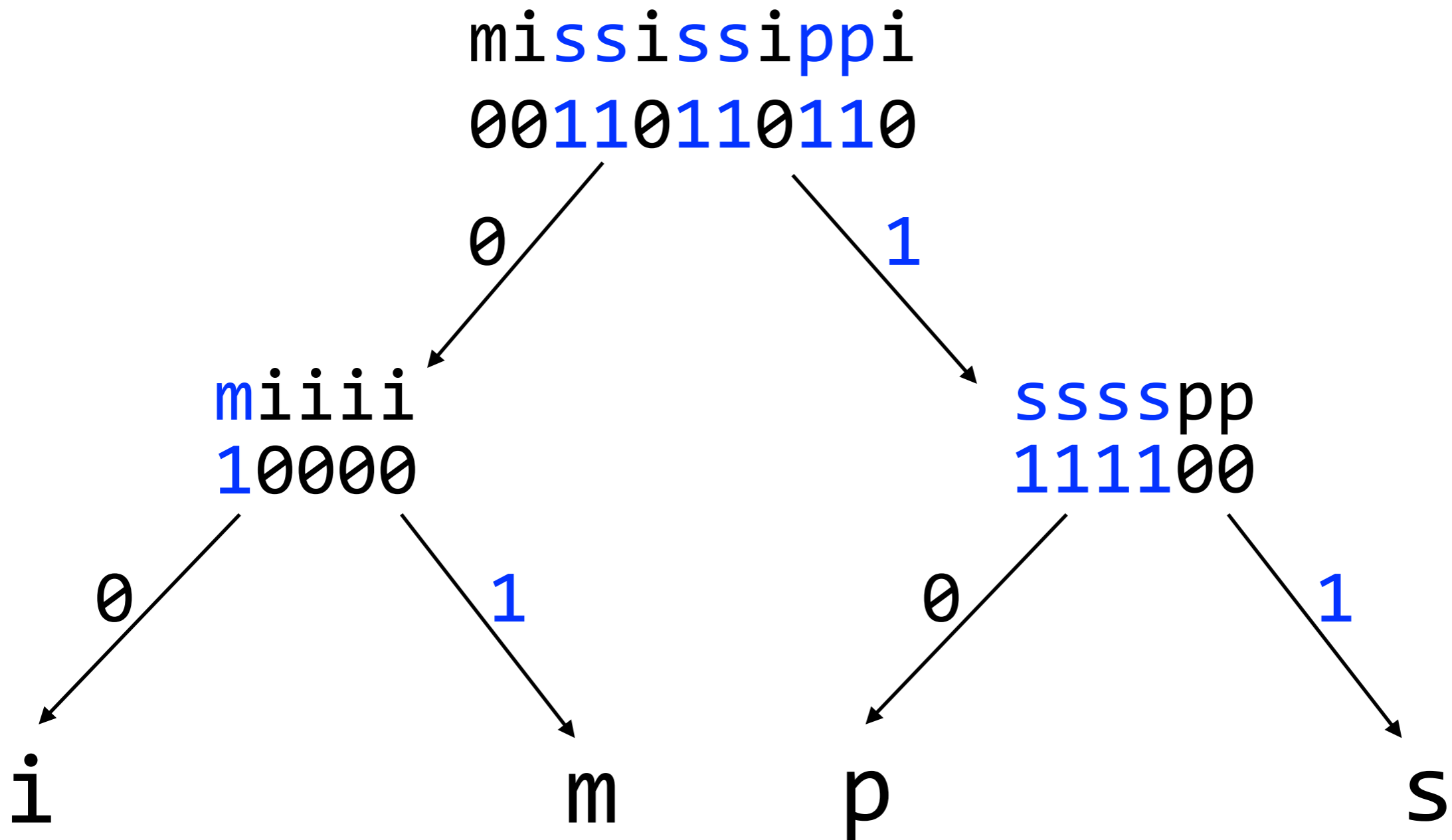
What goes in this next layer?



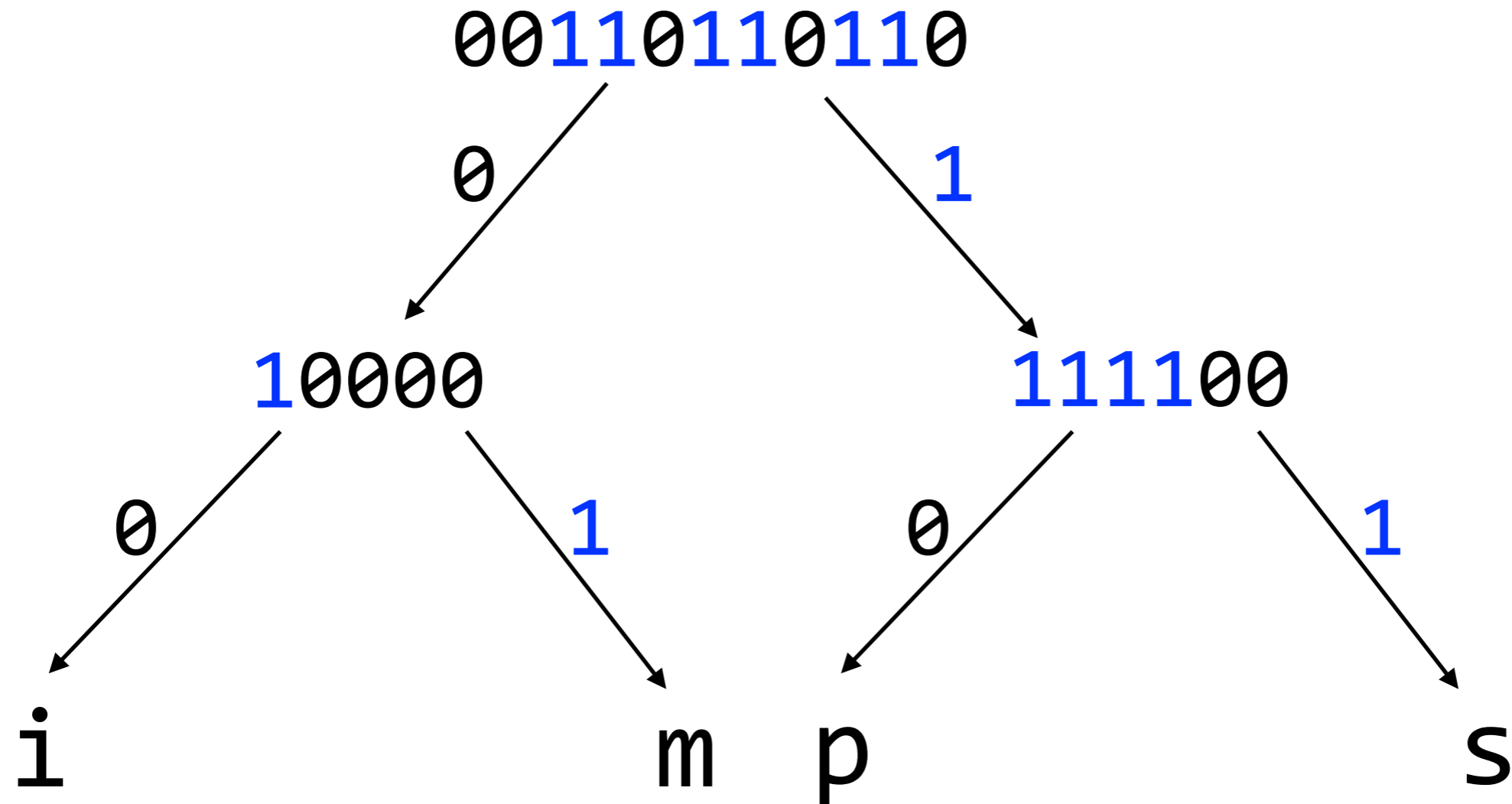
# Wavelet trees



# 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  $\Sigma$  and  $c \in \Sigma$

# Wavelet trees

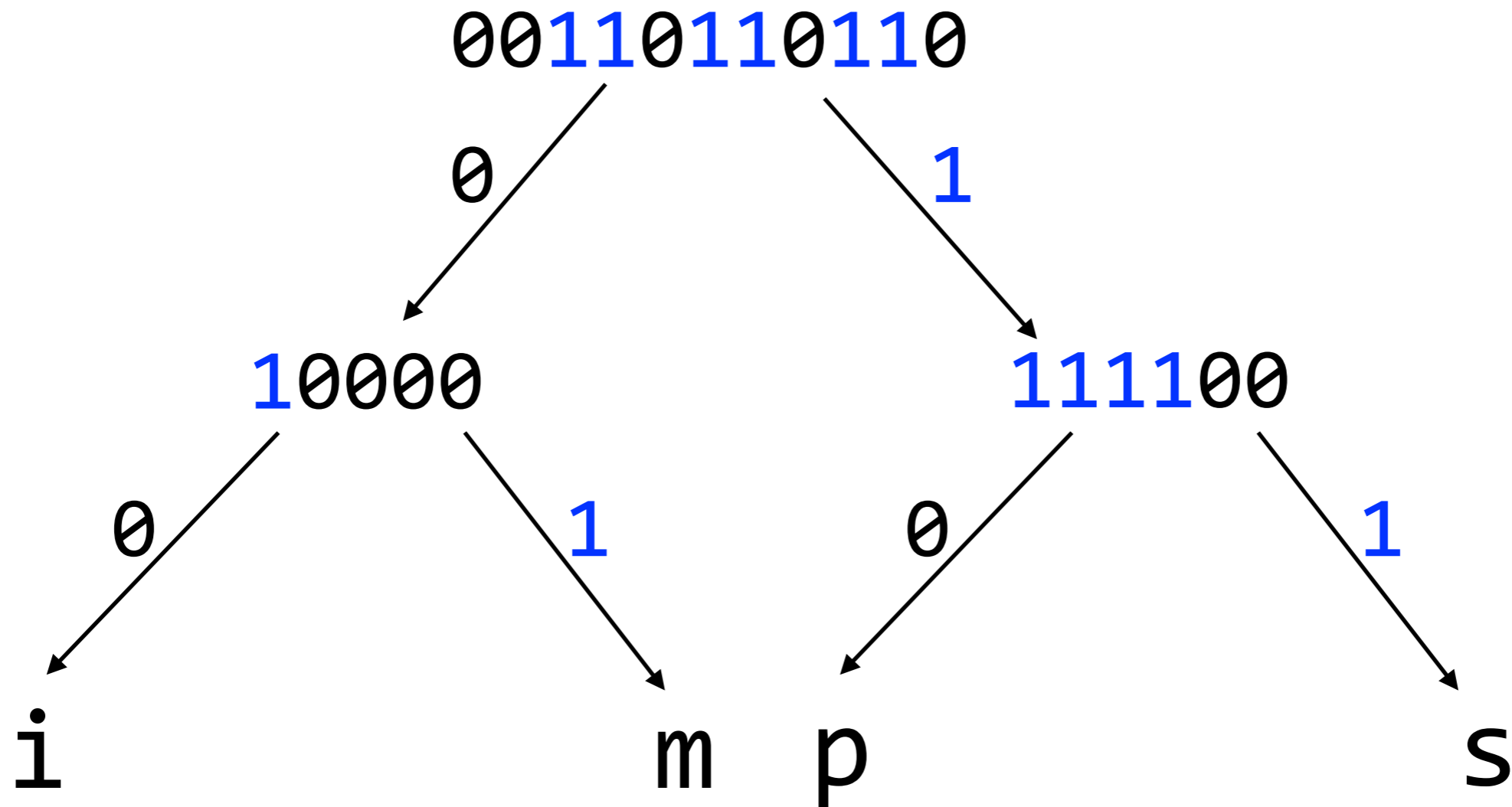
$$S . \text{access}(i) = S[i]$$

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

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

# Wavelet trees

$S$ . **access**(4)

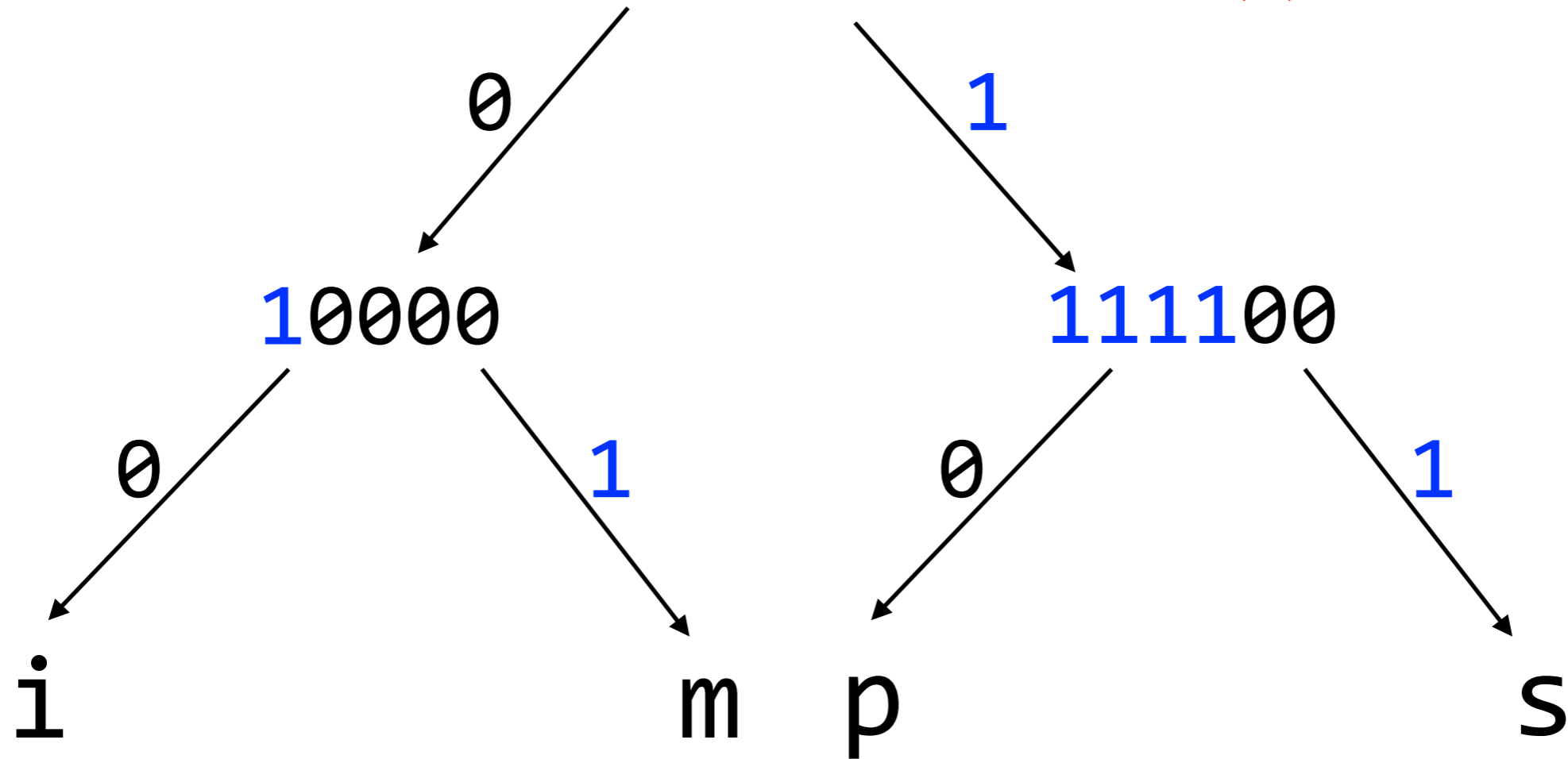


# Wavelet trees

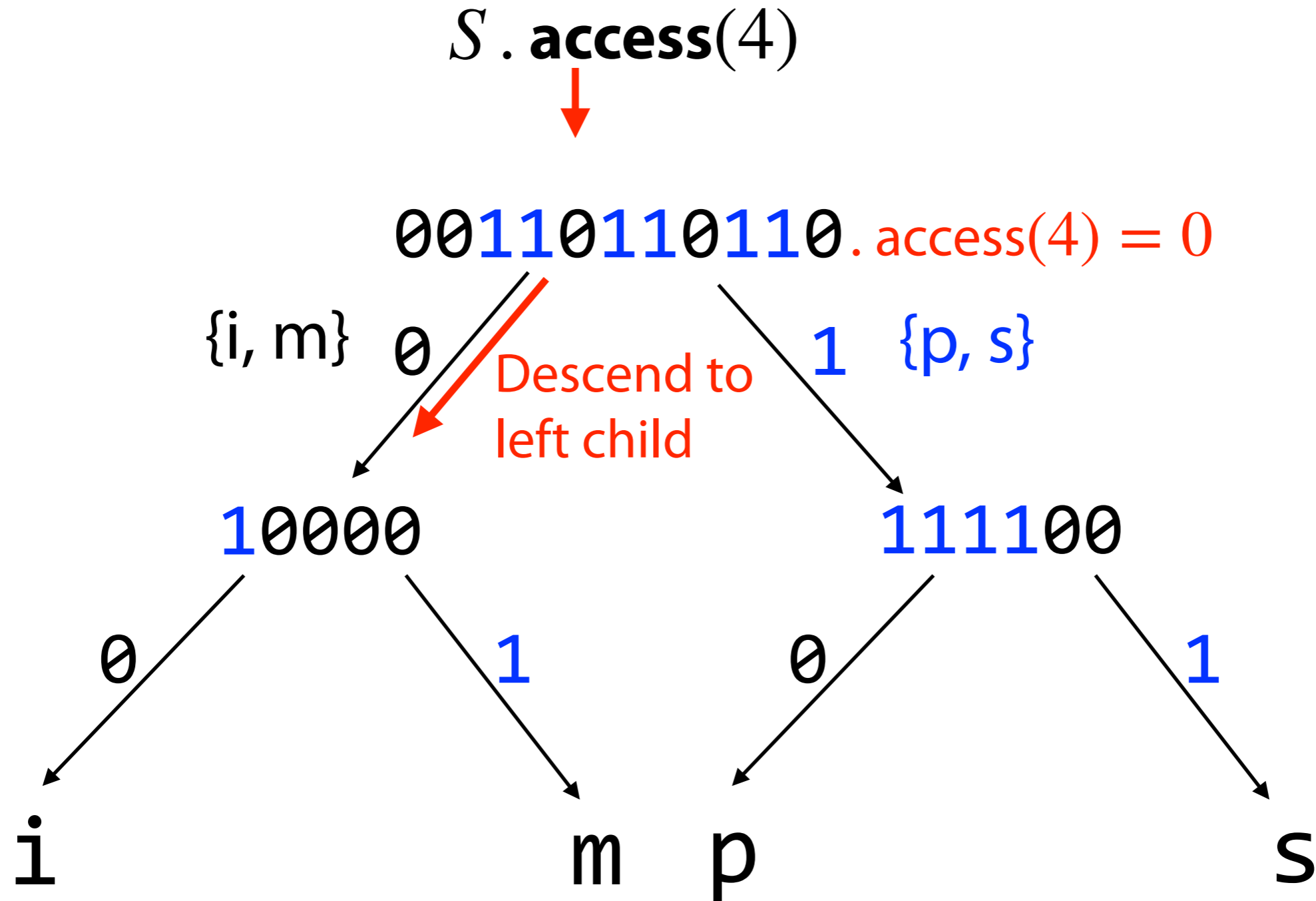
$S$ . **access**(4)



00**110110110**. **access**(4) = 0



# Wavelet trees





# Wavelet trees

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



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

$\{i, m\}$   $\emptyset$

Descend to  
left child

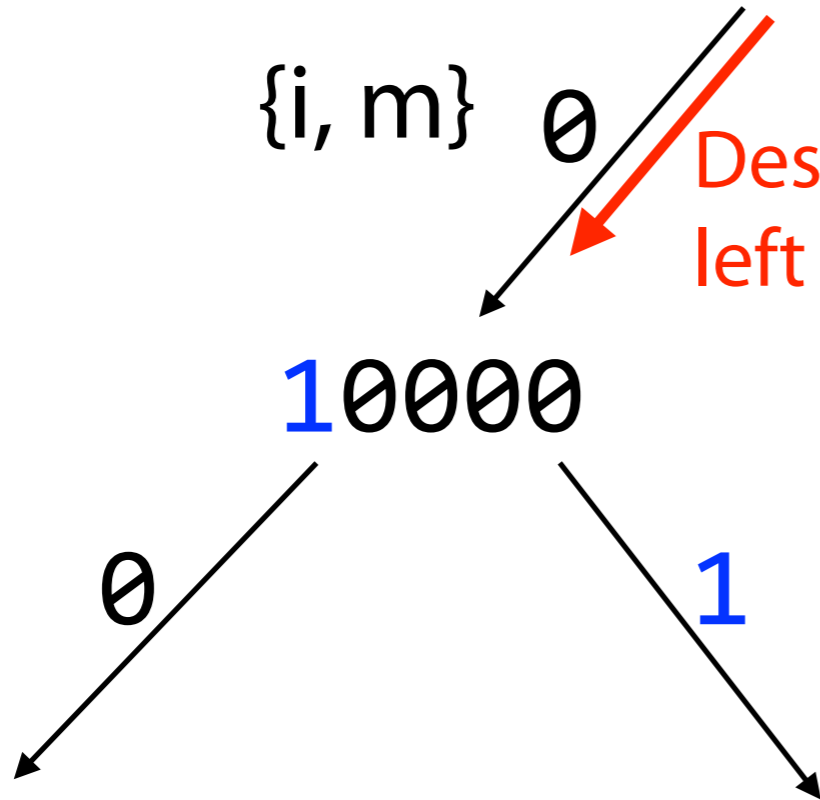
$10000$

$\emptyset$

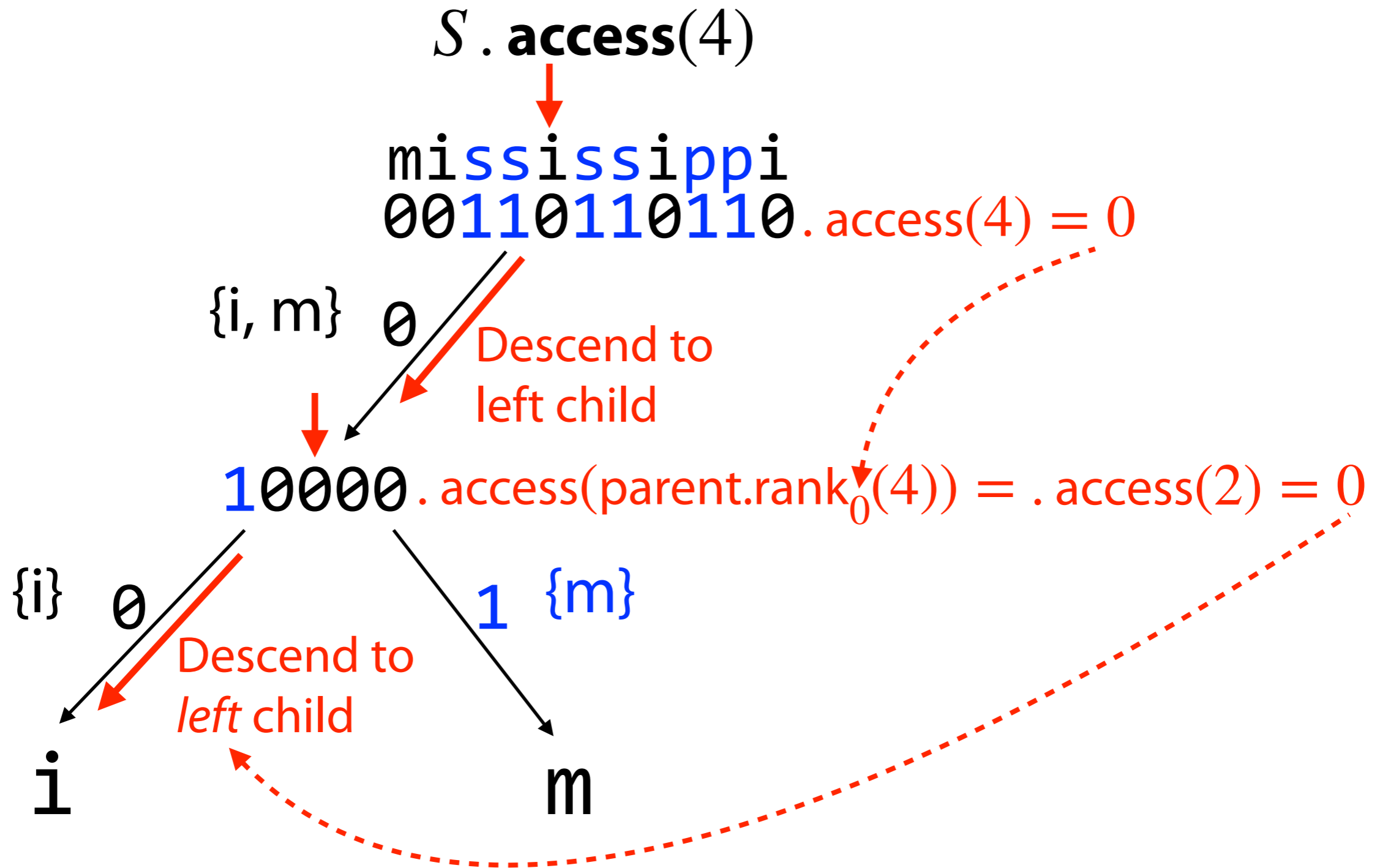
$1$

$i$

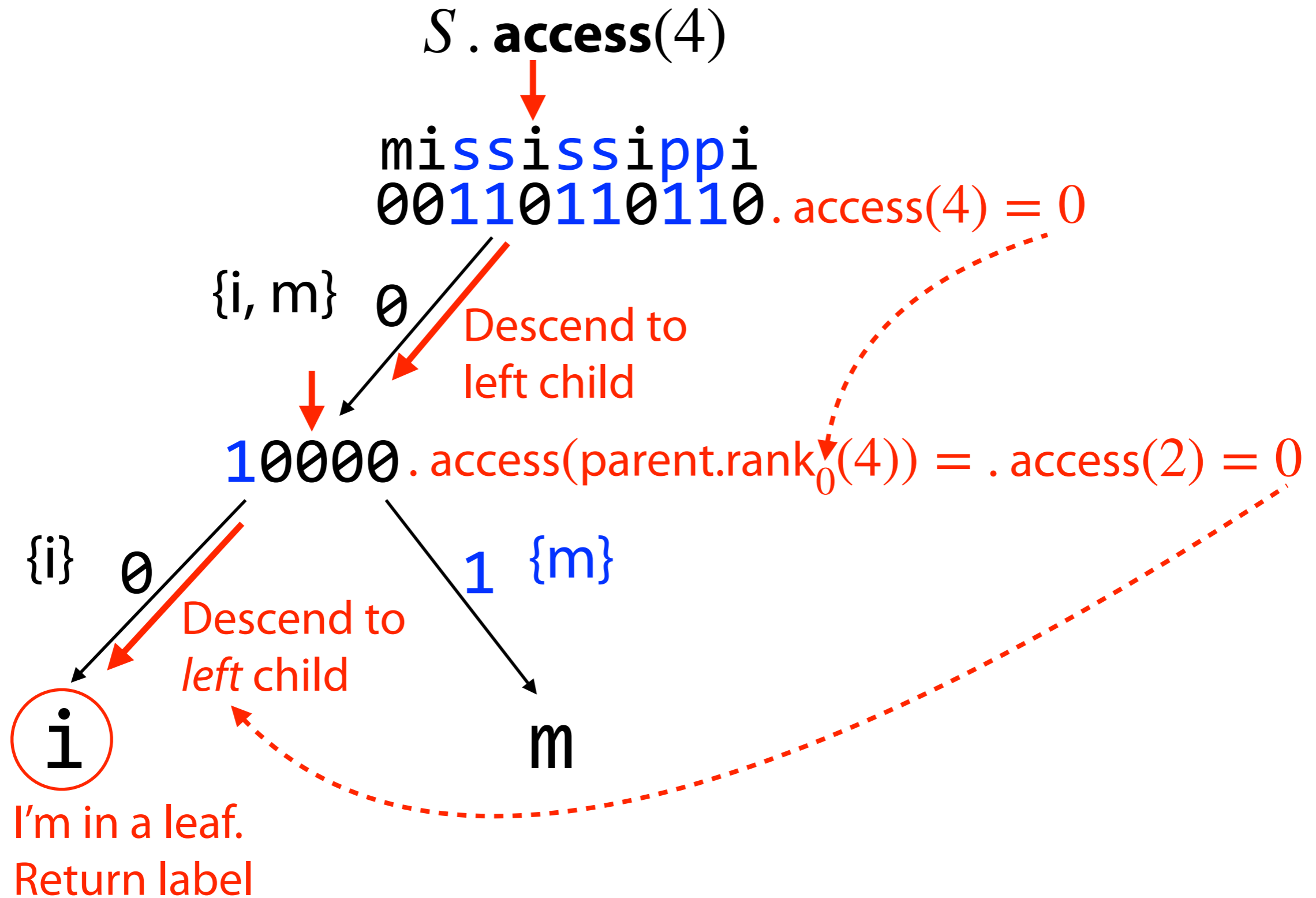
$m$



# Wavelet trees



# Wavelet trees



# Wavelet trees

Wavelet tree access( $i$ ):



$N \leftarrow \text{root}$

while  $N$  is not leaf

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

$b \leftarrow B.\text{access}(i)$

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

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

return  $N.\text{label}$

# Wavelet trees

Wavelet tree access( $i$ ):



$N \leftarrow \text{root}$


while  $N$  is not leaf

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

$b \leftarrow B.\text{access}(i)$

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

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



left or right child?

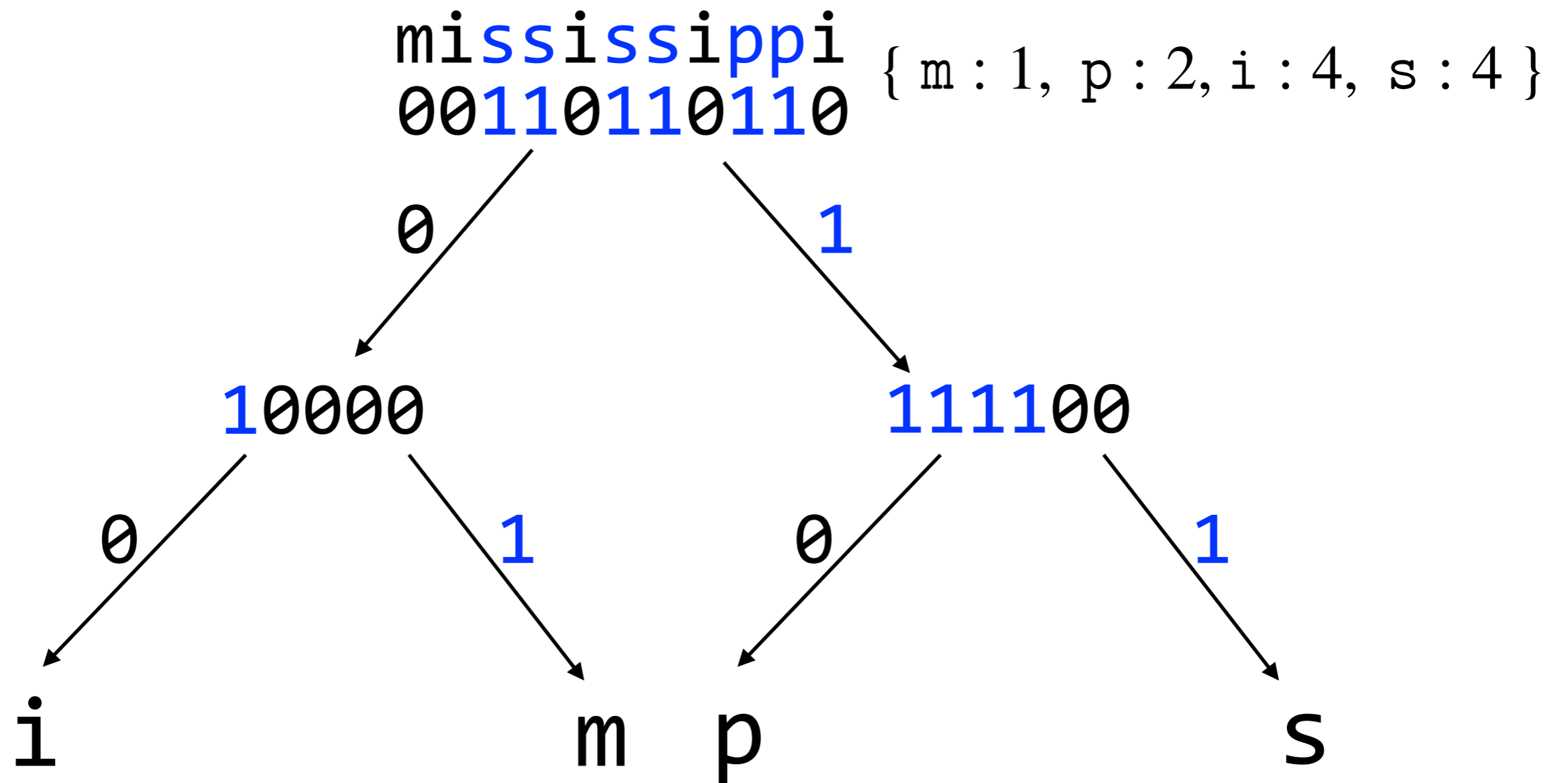
new offset



return  $N.\text{label}$

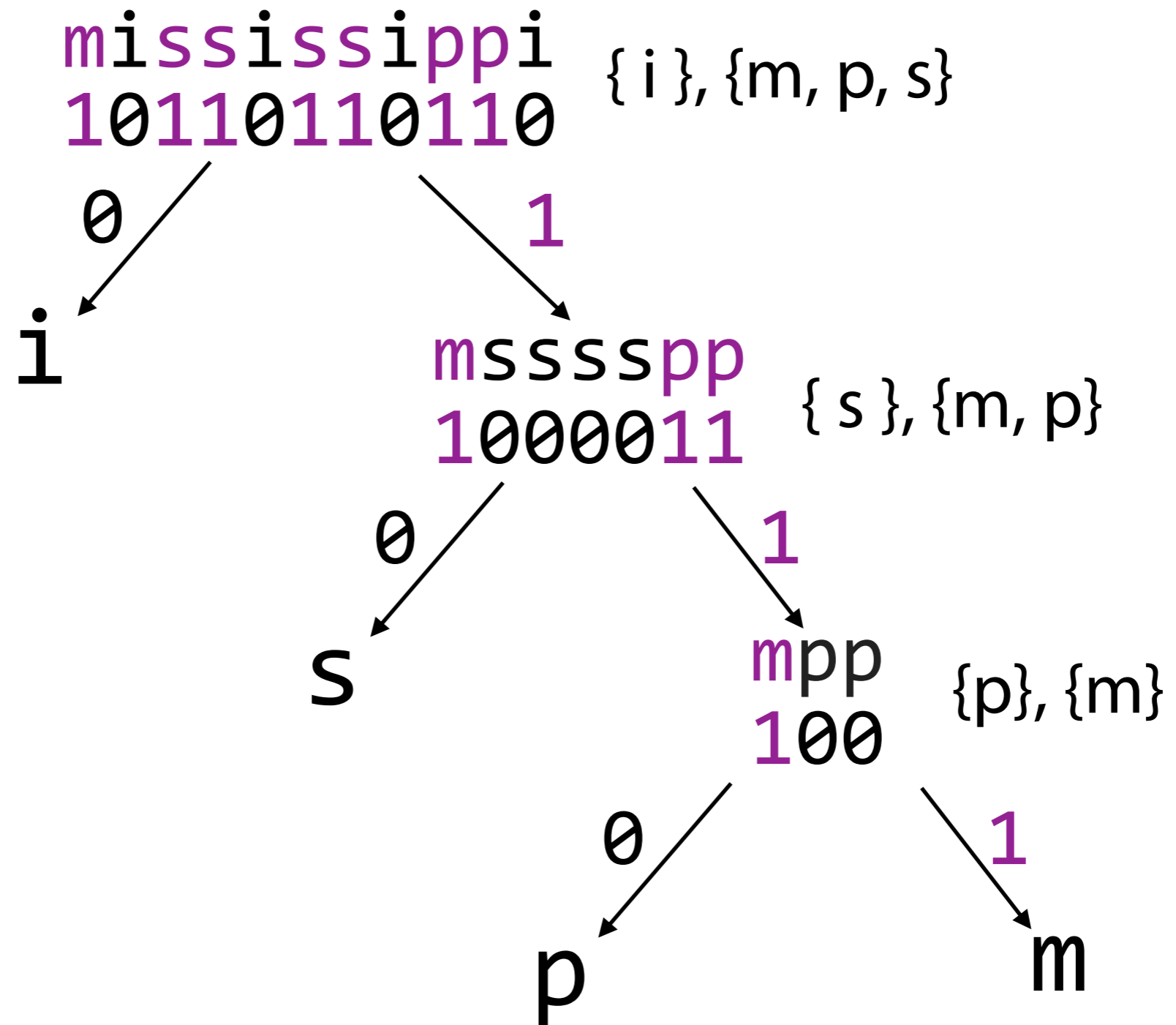
# Wavelet trees

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

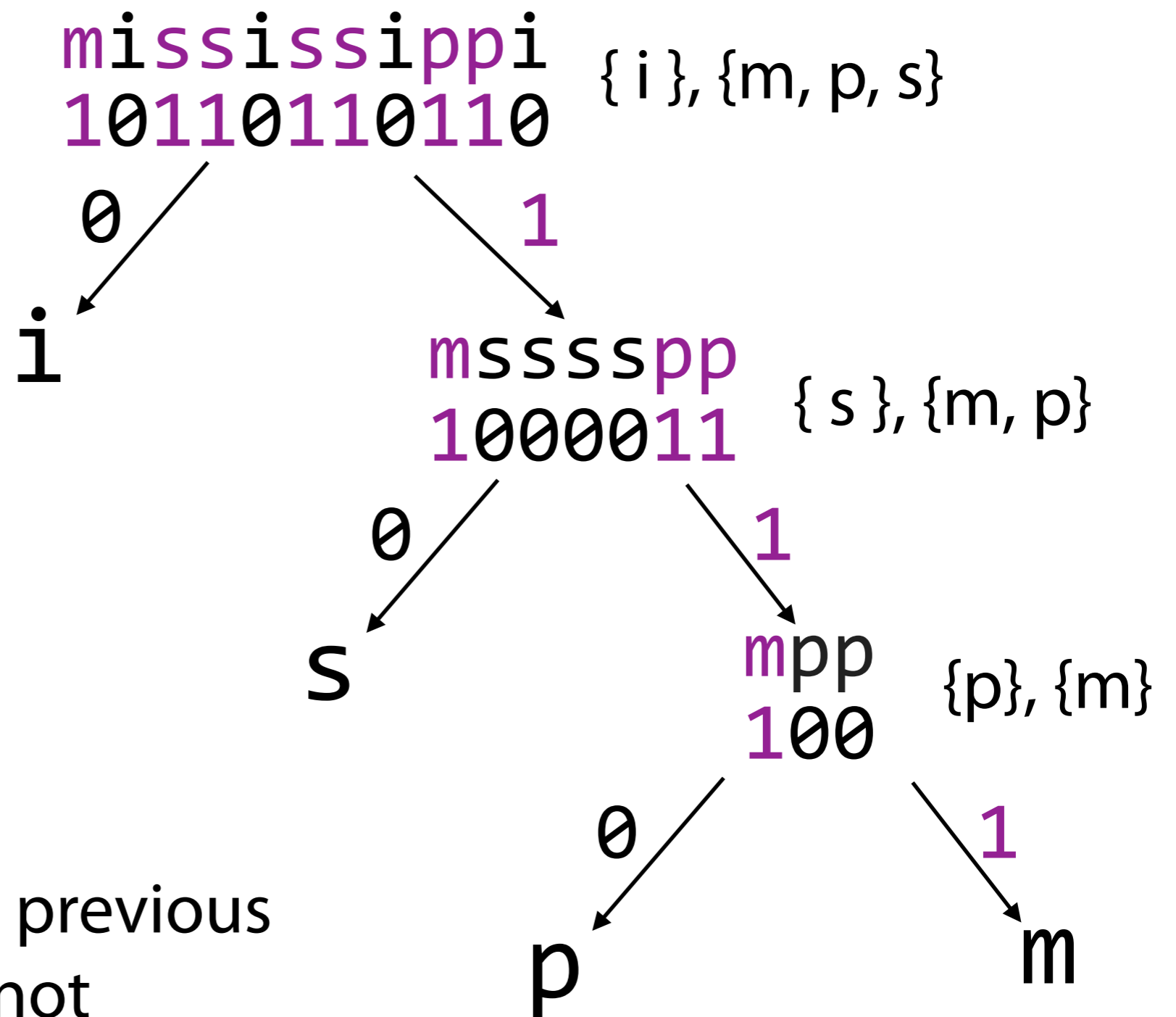
$$C(i) = 0$$

$$C(s) = 10$$

$$C(p) = 110$$

$$C(m) = 111$$

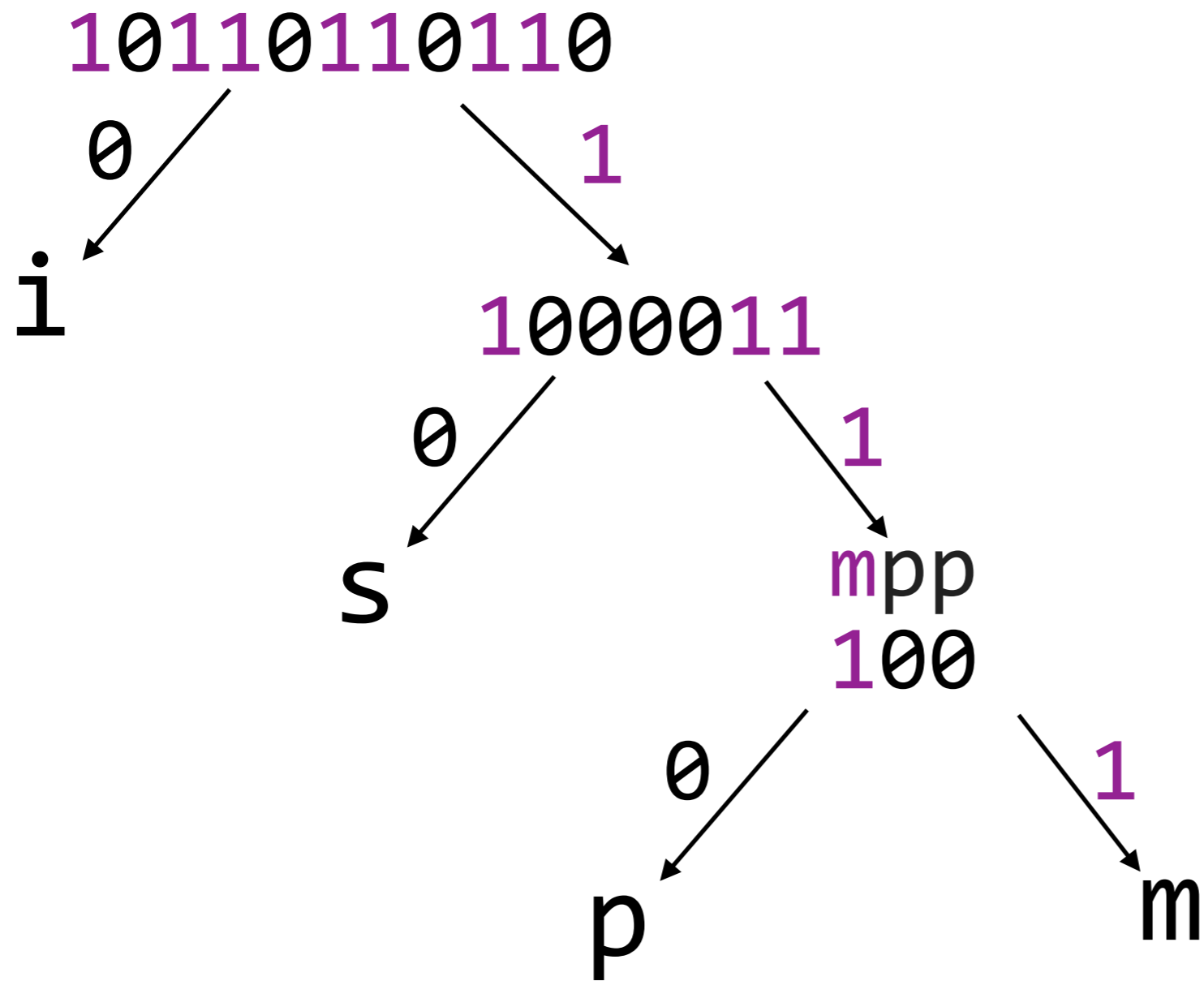
This tree is Huffman; previous  
(balanced) tree was not





# Wavelet trees

mississippi



# 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

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

Rank can ask about *any* character  $c$  at *any* offset  $i$

Path we take **depends on**  $c$

Apart from that, will be similar to access...

# Wavelet trees

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



00**11**0**11**0**11**0

0

1

**1**0000

**1111**00

0

1

0

1

**i**

**m**

**p**

**s**

# Wavelet trees

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



00**11**0**11**0**11**0

0

1

**1**0000

**1111**000

0

1

0

1

**i**

**m**

**p**

**s**

$C(i) = 00$

$C(m) = 01$

$C(p) = 10$

$C(s) = 11$

# Wavelet trees

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



00**11**0**11**0**11**0

0

1

**1**0000

**1111**00

0

1

0

1

**i**

**m**

**p**

**s**

$C(i) = 00$

$C(m) = 01$

$C(p) = 10$

$C(s) = 11$

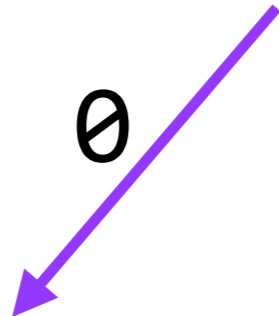
# Wavelet trees

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



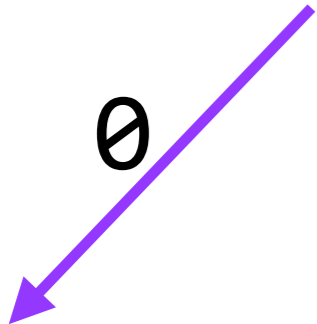
00110110110

0



10000

0



$i$

$$C(i) = 00$$

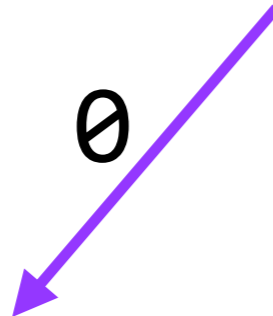
# Wavelet trees

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



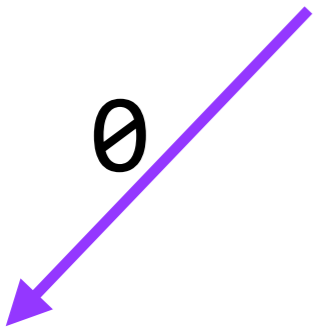
00**11**0**11**0**11**0 .  $\text{rank}_0(6) = 3$

0



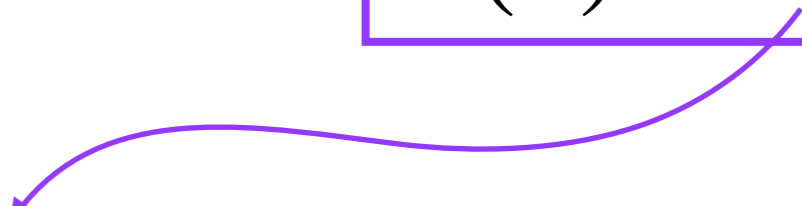
**1**0000

0



**i**

$C(i) = 00$





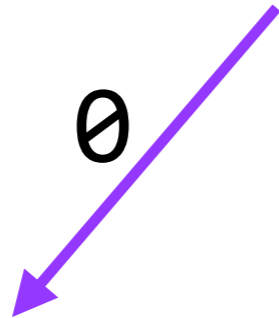
# Wavelet trees

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



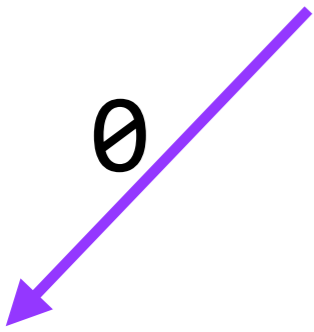
00**11**0**11**0**11**0 .rank<sub>0</sub>(6) = 3

0



**1**0000

0



**i**

$$C(i) = 00$$

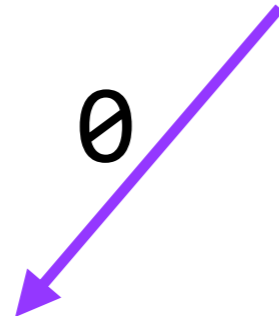
# Wavelet trees

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



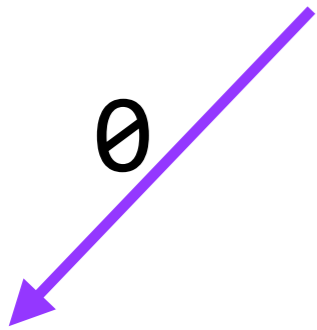
00**11**0**11**0**11**0 .rank<sub>0</sub>(6) = 3

0



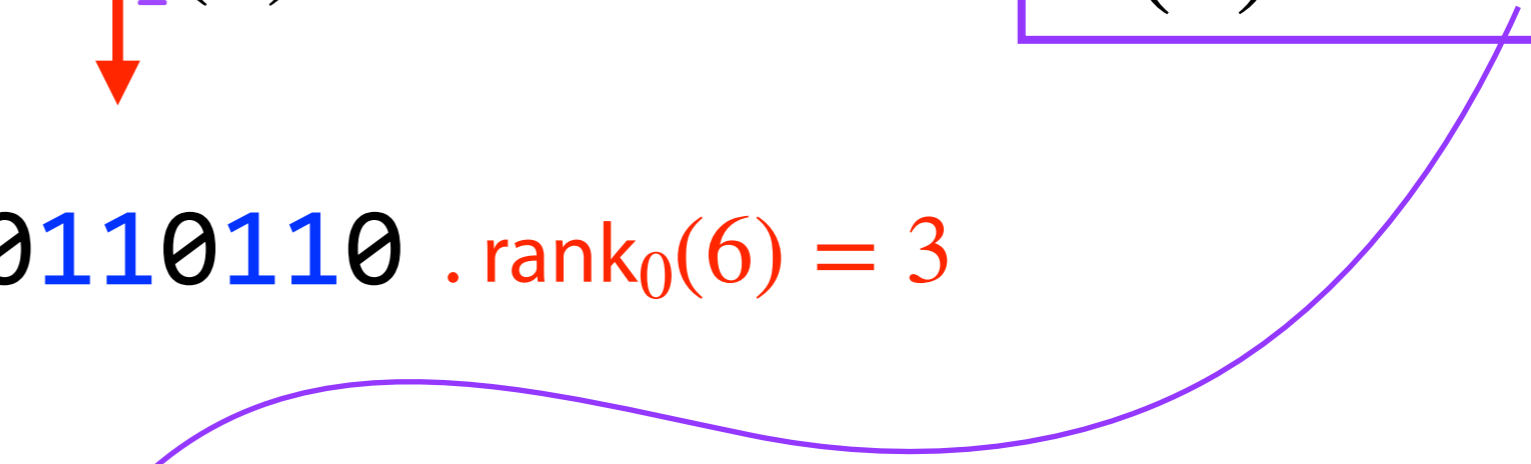
**1**0000 .rank<sub>0</sub>(3) = 2

0



**i**

$C(i) = 00$



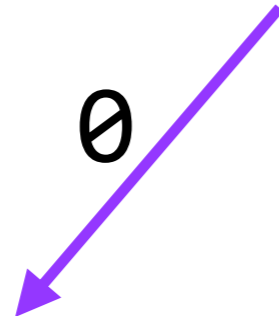
# Wavelet trees

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



00**11**0**11**0**11**0 .  $\text{rank}_0(6) = 3$

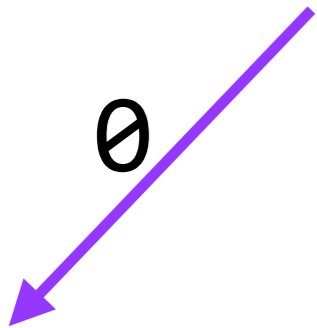
0



**1**0000 .  $\text{rank}_0(3) = 2$

Answer:

0



**i**

$$C(i) = 00$$

# Wavelet trees

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

$$C(i) = 00$$

00**11**0**11**0**11**0 .rank<sub>0</sub>(6) = 3

0

**1**0000 .rank<sub>0</sub>(3) = 2

Answer: 2

0

**i**

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

Overall answer equals answer from final rank query

# Wavelet trees

Wavelet tree rank <sub>$x$</sub> ( $i$ ):



$N \leftarrow \text{root}$

$k \leftarrow 0$

while  $N$  is not leaf

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

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

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

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

$k \leftarrow k + 1$

return  $i$

# Wavelet trees

Wavelet tree rank <sub>$x$</sub> ( $i$ ):



$N \leftarrow \text{root}$

$k \leftarrow 0$

while  $N$  is not leaf

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

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

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

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

$k \leftarrow k + 1$

return  $i$

left or right child? depends on *code*

# Wavelet trees

Wavelet tree rank <sub>$x$</sub> ( $i$ ):



$N \leftarrow \text{root}$

$k \leftarrow 0$

while  $N$  is not leaf

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

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

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

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



$k \leftarrow k + 1$

return  $i$

left or right child? depends on *code*

# Wavelet trees

Wavelet tree rank <sub>$x$</sub> ( $i$ ):



$N \leftarrow \text{root}$

$k \leftarrow 0$

while  $N$  is not leaf

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

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

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

left or right child? depends on *code*

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

new offset

$k \leftarrow k + 1$

on to next bit in code

return  $i$



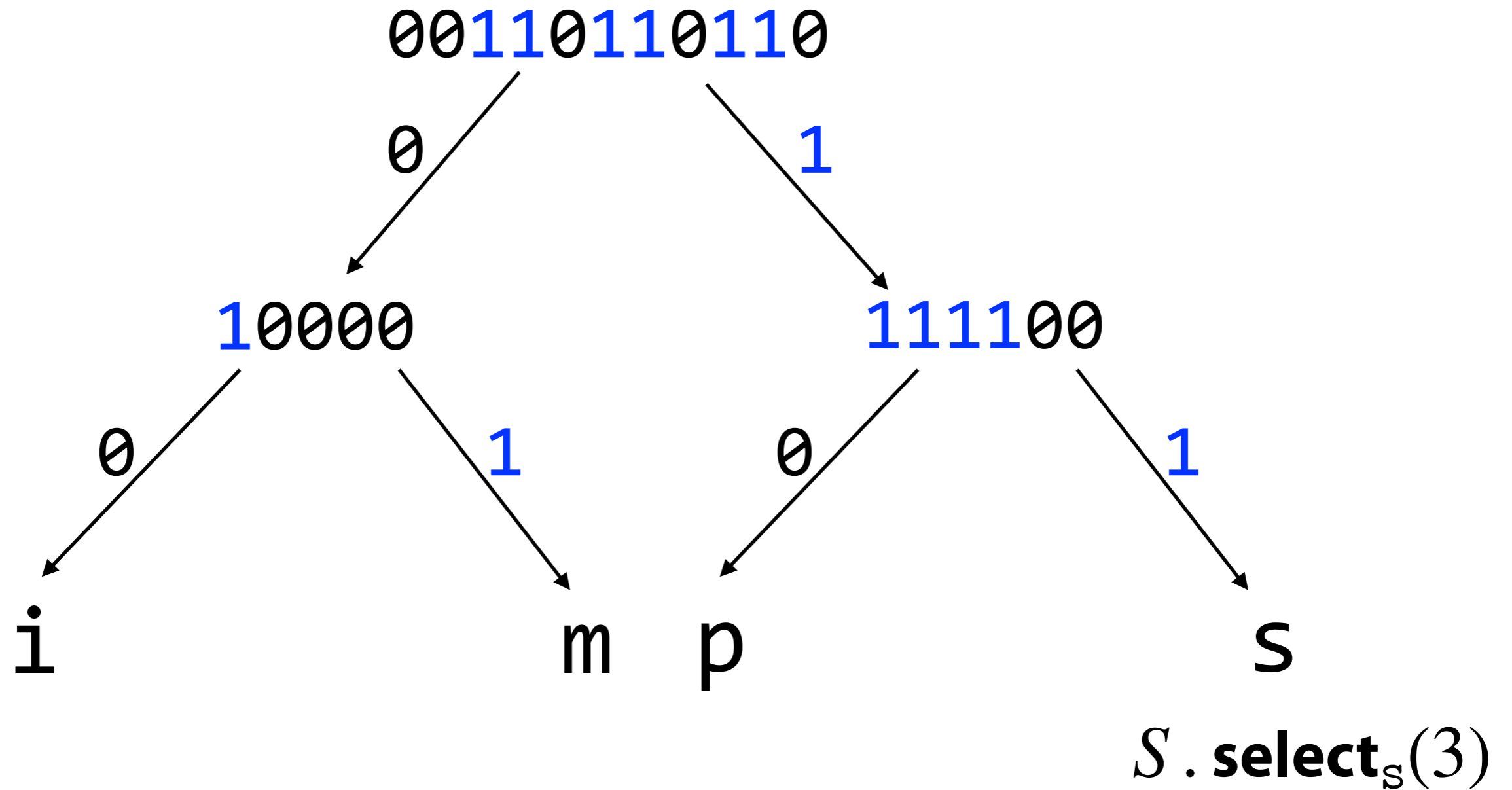
# Wavelet trees

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

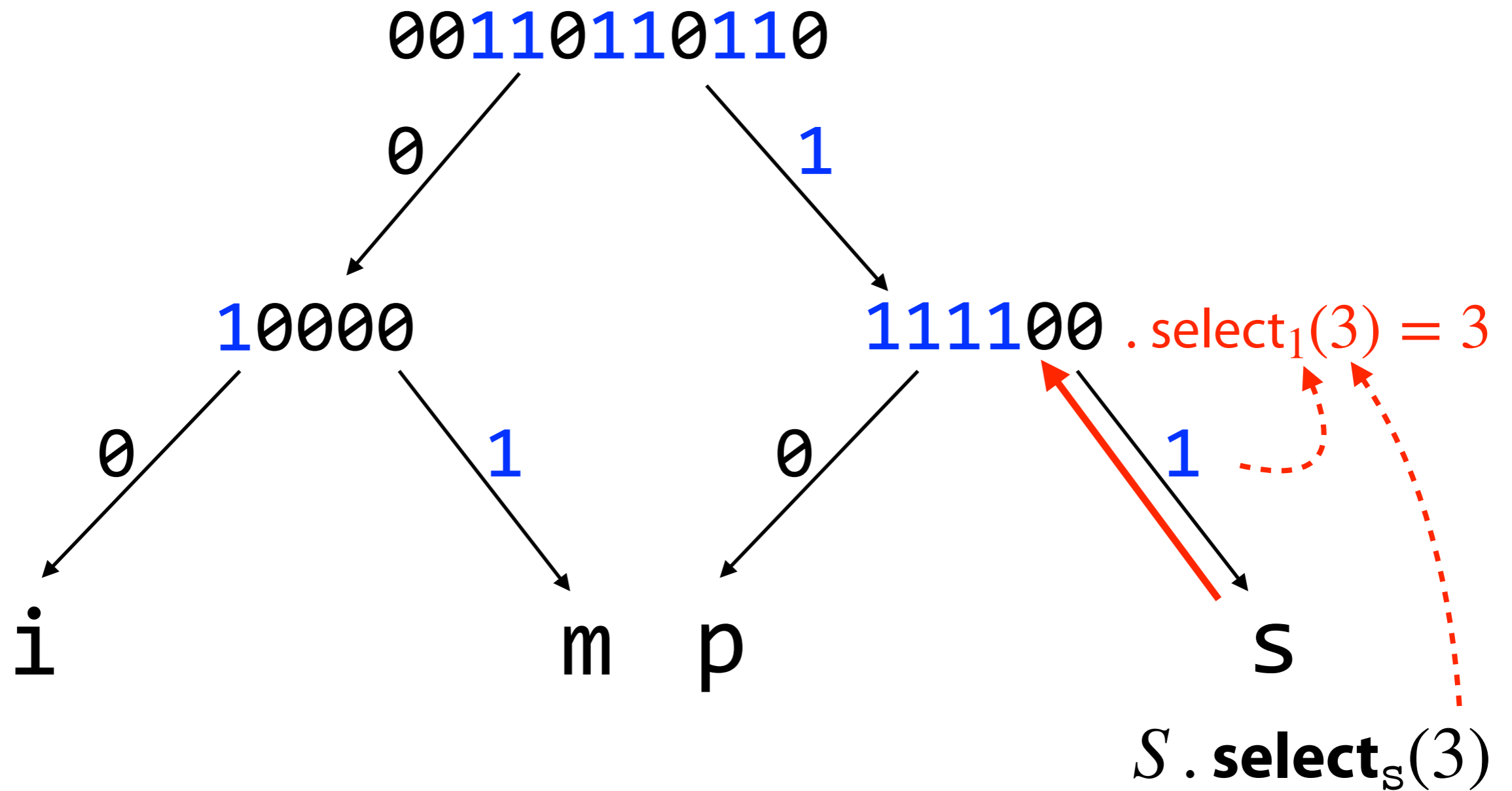
$$S . \text{rank}_c(i) = \sum_{j=0}^{i-1} \begin{cases} 1 & \text{if } S[j] = c \\ 0 & \text{otherwise} \end{cases}$$

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

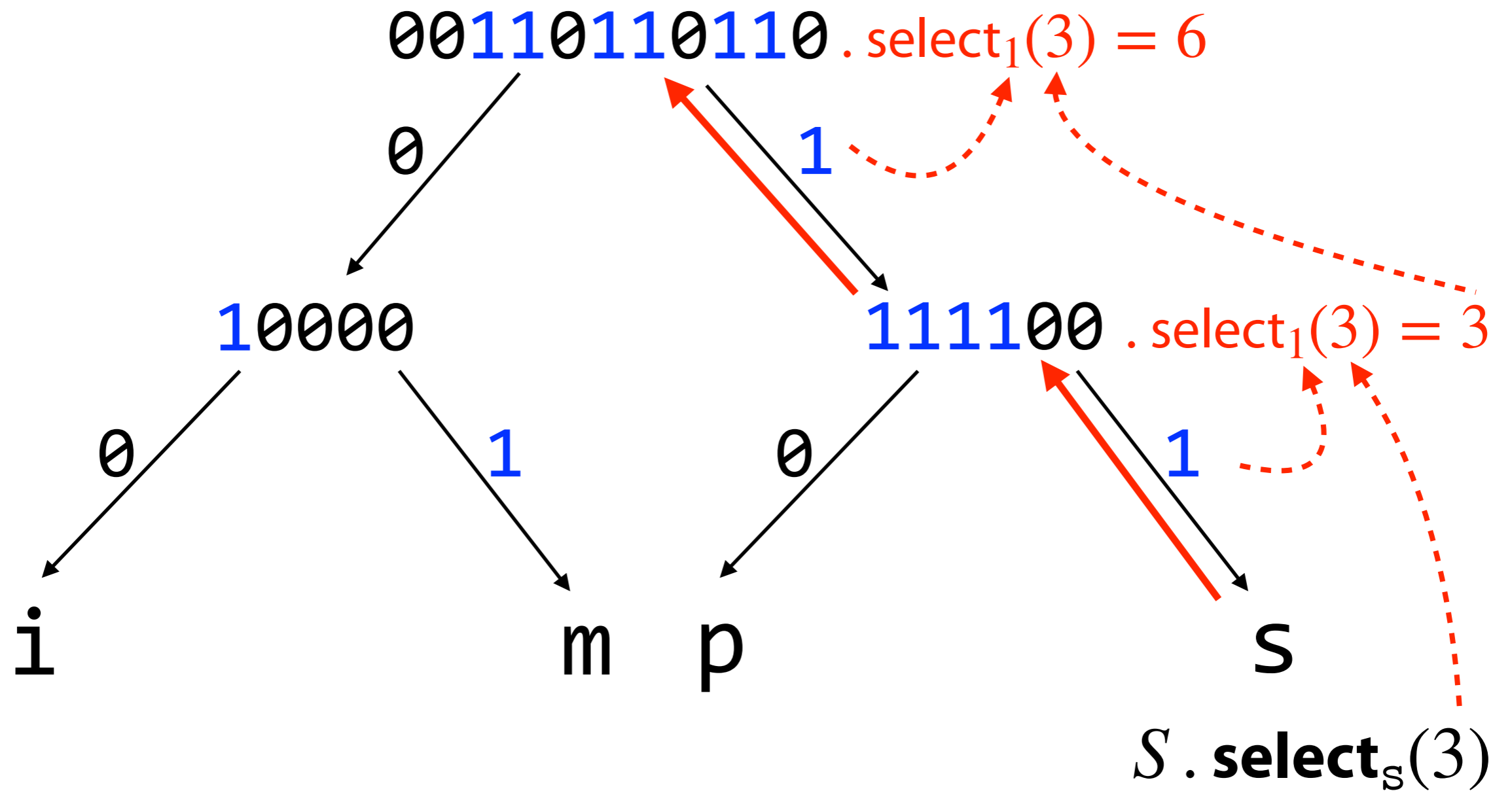
# Wavelet trees



# Wavelet trees

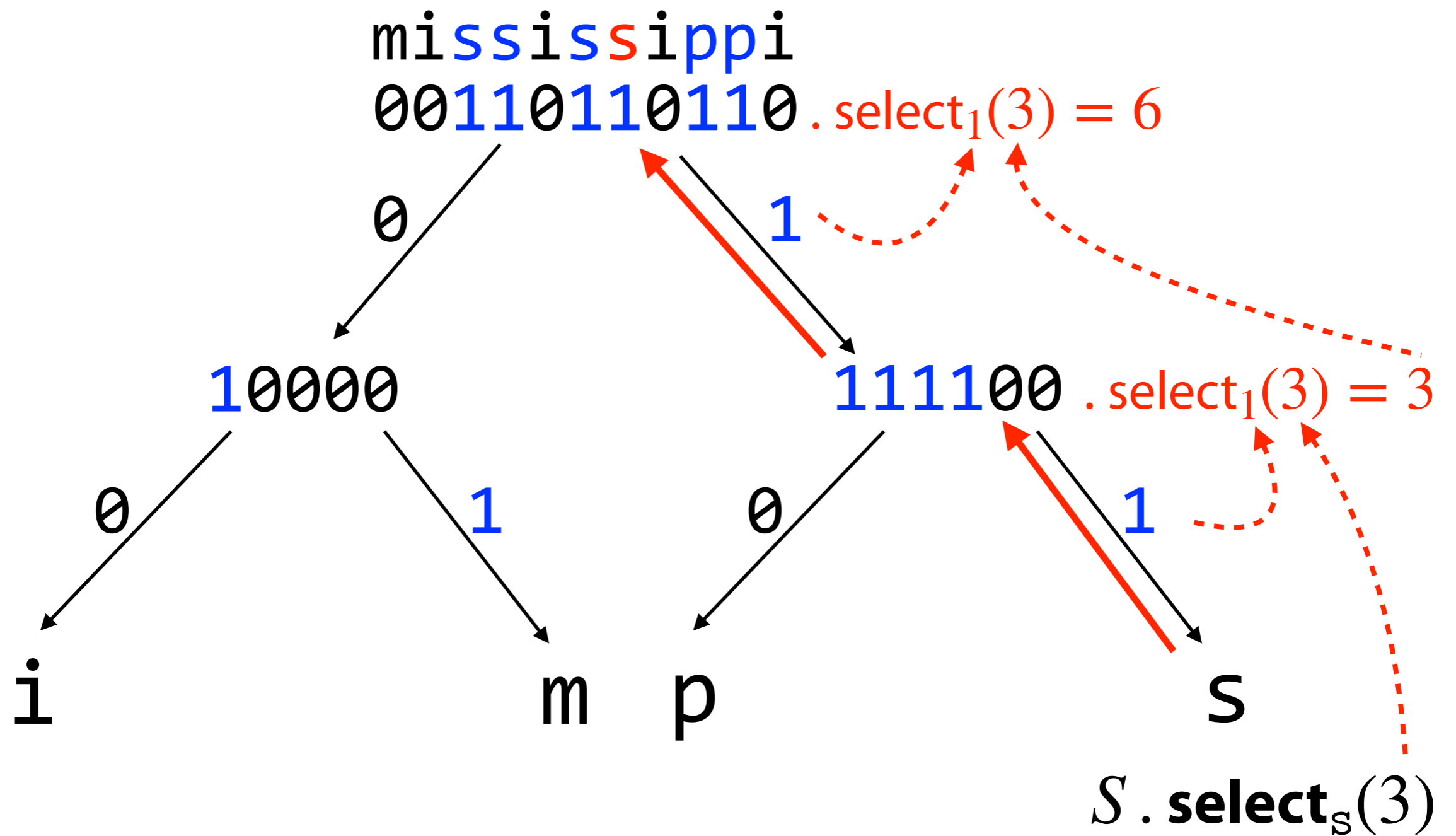


# Wavelet trees



# Wavelet trees

Answer: 6



# Wavelet trees

Wavelet tree select<sub>*x*</sub>(*i*):



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

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

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

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

$N \leftarrow \text{leaf}(x)$  descend all the way first

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

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

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

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

descend all the way first

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

while  $N$  is not root

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

step upward

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

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

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

$k \leftarrow k - 1$

return  $i$



# Wavelet trees

Wavelet tree select <sub>$x$</sub> ( $i$ ):



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

descend all the way first



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

while  $N$  is not root

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

step upward



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

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

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

offset in prev. round becomes

rank in next



$k \leftarrow k - 1$

return  $i$

# Wavelet trees

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

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

descend all the way first

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

while  $N$  is not root

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

step upward

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

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

offset in prev. round becomes

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

rank in next

$k \leftarrow k - 1$

move to previous bit in code

return  $i$