# Wavelet Trees

02-714
Slides by Carl Kingsford

Following: Navarro,
Wavelet Trees for All, CPM 2012, pp. 2-26

# Operations on strings

- $\text{rank}_c(S, i)$ := the number of char $c$ at or before position $i$ in $S$.

- $\text{select}_c(S, j)$ := the position of the $j^{th}$ occurrence of $c$ in S.

- $S[i]$ = "access character $i$"

Note: $\text{rank}_c(S, \text{select}_c(S, j)) = j$, so rank and select are inverses of each other.

**Goal**: rank, select, access in quickly while using small space.

# Operations on bit vectors

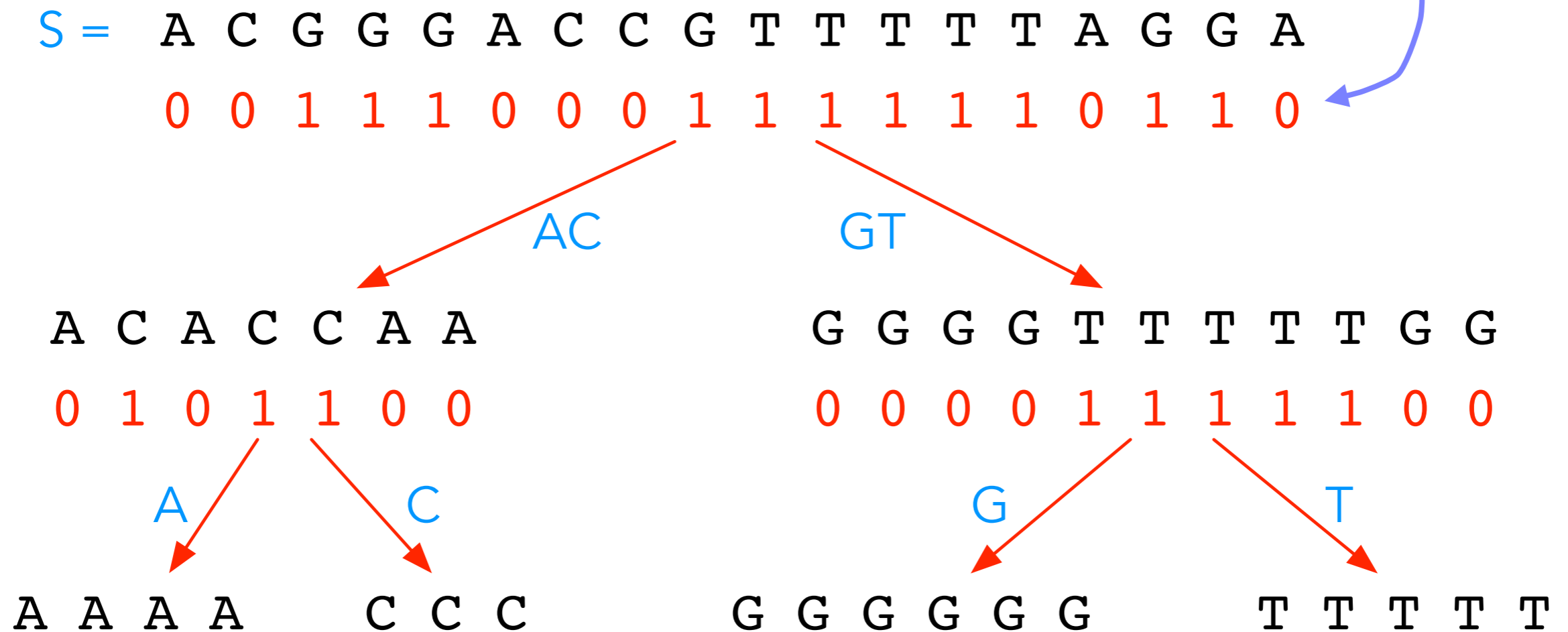Our operations will depend on similar operations on bit vectors:

- $\text{rank}_1(S, i)$ := the number of 1 bits at or before position $i$ in $S$.

- $\text{select}_1(S, j)$ := the position of the $j^{\text{th}}$ 1 bit in S.

- $\text{rank}_0(S, i)$ and $\text{select}_0(S, j)$ are defined analogously.

$S[i]$ = "access bit i" = $\text{rank}_1(S, i) - \text{rank}_1(S, i - 1)$

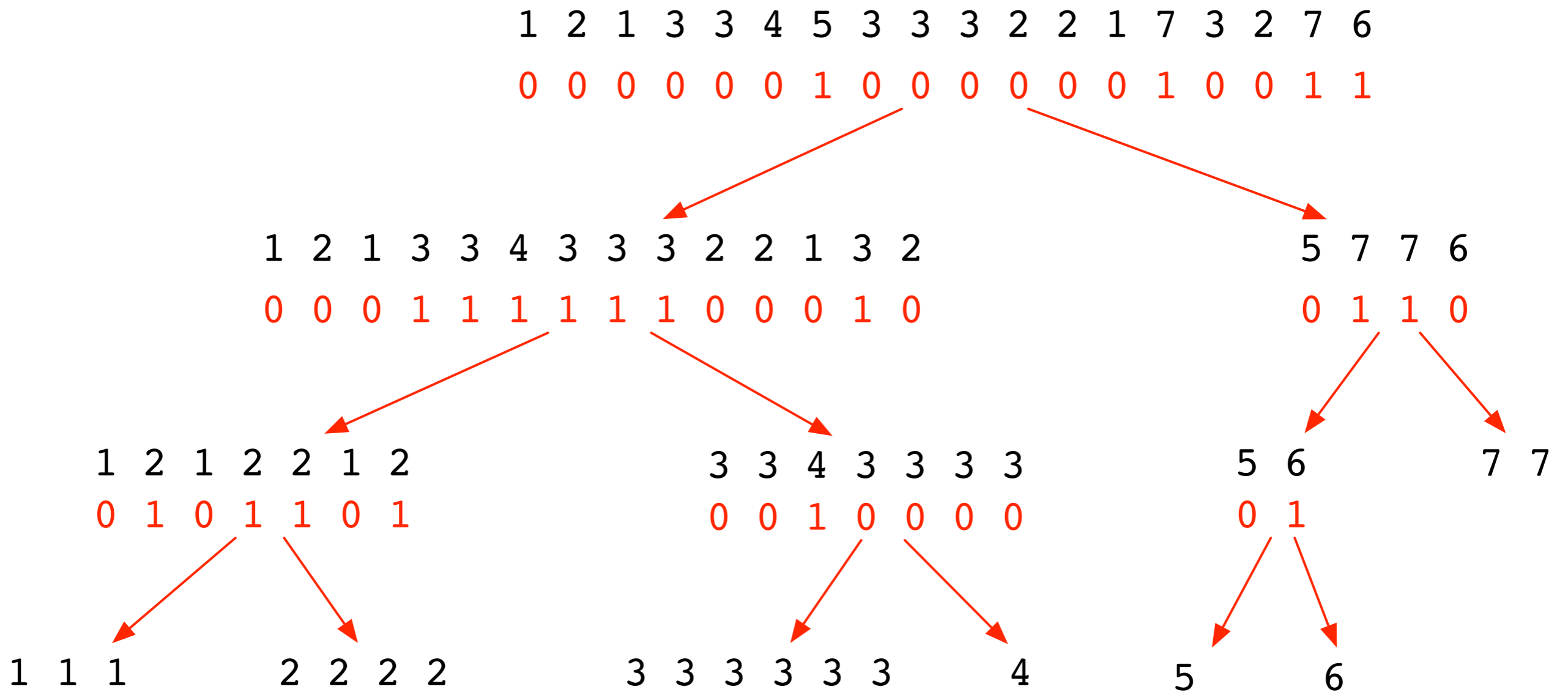We will see later how to implement these operations to run in O(1) time for binary vectors.

# Wavelet Tree

0: letter ∈ first half of alphabet
1: letter ∈ first half of alphabet

S =  A C G G G A C C G T T T T T A G G A
     0 0 1 1 1 0 0 0 1 1 1 1 1 1 0 1 1 0

AC

GT

A C A C C A A
0 1 0 1 1 0 0

G G G G T T T T T G G
0 0 0 0 1 1 1 1 1 0 0

A

C

G

T

A A A A      C C C

G G G G G G      T T T T T

# Wavelet Tree, Example 2

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

1 2 1 3 3 4 3 3 3 2 2 1 3 2
0 0 0 1 1 1 1 1 1 0 0 0 1 0

5 7 7 6
0 1 1 0

1 2 1 2 2 1 2
0 1 0 1 1 0 1

3 3 4 3 3 3 3
0 0 1 0 0 0 0

5 6
0 1

7 7

1 1 1

2 2 2 2

3 3 3 3 3 3

4

5

6

# S[i]

Go left if bit is 0, go right if bit is 1.

Use $rank_{01}()$ to map a bit at a node to the right place in the children.

*i*

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

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

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

```
  1 1 1      2 2 2 2          3 3 3 3 3 3      4        5    6
```
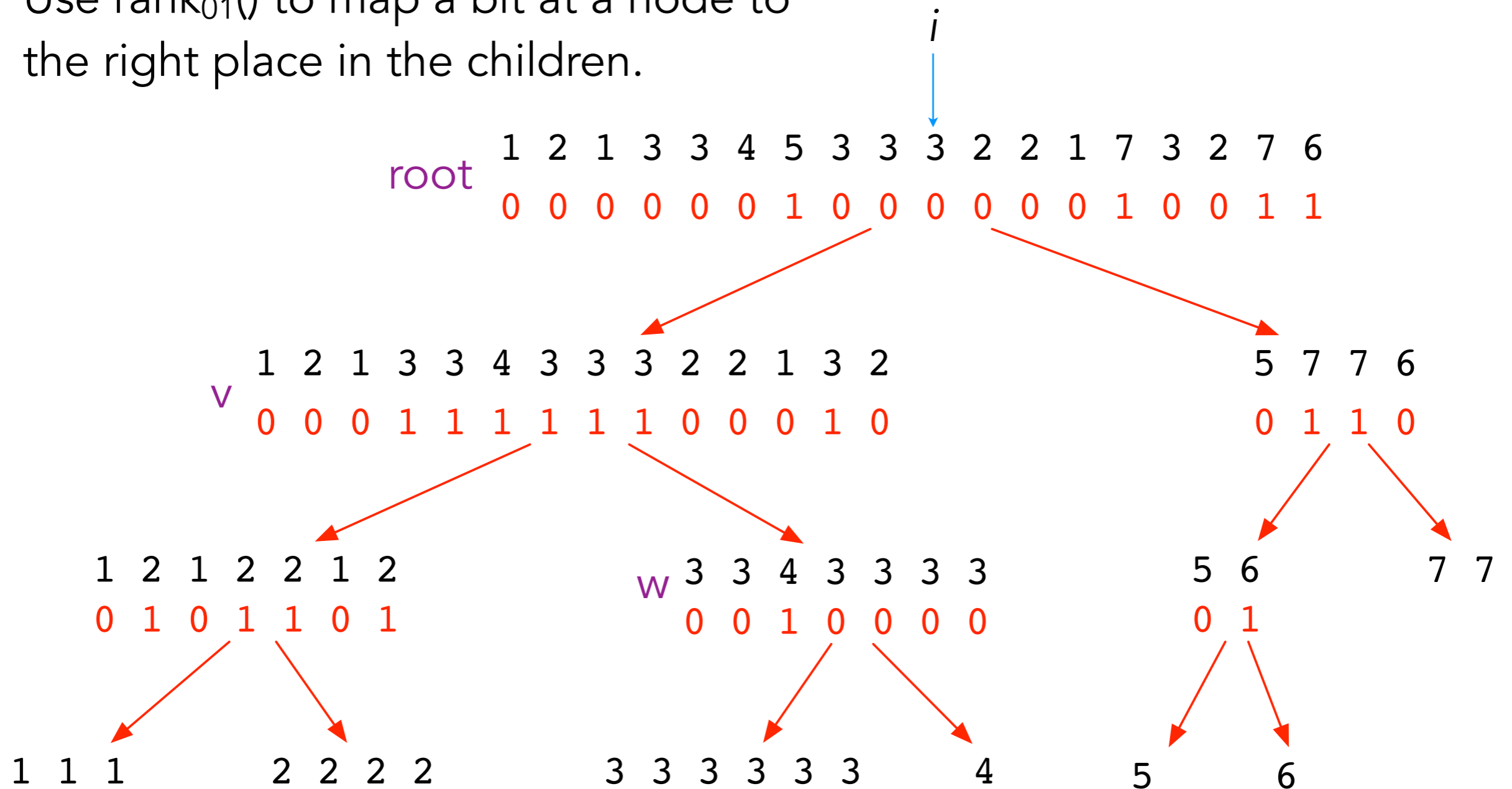
# S[i]

Go left if bit is 0, go right if bit is 1.

Use $\text{rank}_{01}()$ to map a bit at a node to the right place in the children.

*i*

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

$i_v = \text{rank}_0(\text{root}, i)$

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

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

1 1 1          2 2 2 2          3 3 3 3 3 3          4          5          6

# S[i]

Go left if bit is 0, go right if bit is 1.

Use $\text{rank}_{01}()$ to map a bit at a node to the right place in the children.

*i*

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

$i_v = \text{rank}_0(\text{root}, i)$

v
1 2 1 3 3 4 3 3 3 2 2 1 3 2
0 0 0 1 1 1 1 1 1 0 0 0 1 0

5 7 7 6
0 1 1 0

$i_w = \text{rank}_1(v, i_v)$

1 2 1 2 2 1 2
0 1 0 1 1 0 1

w
3 3 4 3 3 3 3
0 0 1 0 0 0 0

5 6
0 1

7 7

1 1 1      2 2 2 2

3 3 3 3 3 3      4

5      6

# S[i]

Go left if bit is 0, go right if bit is 1.

Use $rank_{01}()$ to map a bit at a node to the right place in the children.

$i$

root

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

$i_v = rank_0(root, i)$

v

1 2 1 3 3 4 3 3 3 2 2 1 3 2
0 0 0 1 1 1 1 1 1 0 0 0 1 0

5 7 7 6
0 1 1 0

$i_w = rank_1(v, i_v)$

1 2 1 2 2 1 2
0 1 0 1 1 0 1

w

3 3 4 3 3 3 3
0 0 1 0 0 0 0

5 6
0 1

7 7

1 1 1

2 2 2 2

3 3 3 3 3 3

4

5

6

$i_3 = rank_0(w, i_w) = 5$

# rank$_c$(S,i)

Go left if **c is in first half of alphabet**, go right if c is in second half of alphabet.

Use rank$_{01}$() to map a bit at a node to the right place in the children.

*i*

rank$_5$(S,*i*)

root

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

1 2 1 3 3 4 3 3 3 2 2 1 3 2
0 0 0 1 1 1 1 1 1 0 0 0 1 0

5 7 7 6 V
0 1 1 0

1 2 1 2 2 1 2
0 1 0 1 1 0 1

3 3 4 3 3 3 3
0 0 1 0 0 0 0

5 6 W
0 1

7 7

1 1 1    2 2 2 2

3 3 3 3 3 3    4

5    6

# rank$_c$(S,i)

Go left if **c is in first half of alphabet**, go right if c is in second half of alphabet.

Use rank$_{01}$() to map a bit at a node to the right place in the children.

*i*

rank$_5$(S,*i*)

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

$i_v$ = rank$_1$(root, i)

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

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

```
1 1 1     2 2 2 2          3 3 3 3 3 3     4          5     6
```

# rank$_c$(S,i)

Go left if **c is in first half of alphabet**, go right if c is in second half of alphabet.

Use rank$_{01}$() to map a bit at a node to the right place in the children.

$i$

rank$_5$(S,$i$)

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

$i_v$ = rank$_1$(root, i)

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

$i_w$ = rank$_0$(v, i$_v$)

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

```
 1 1 1      2 2 2 2      3 3 3 3 3 3      4      5      6
```

# rank$_c$(S,i)

Go left if **c is in first half of alphabet**, go right if c is in second half of alphabet.

Use rank$_{01}$() to map a bit at a node to the right place in the children.

*i*

rank$_5$(S,*i*)

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

root

0 0 0 0 0 0 1 0 0 0 0 0 0 1 0 0 1 1

$i_v$ = rank$_1$(root, i)

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

0 0 0 1 1 1 1 1 1 0 0 0 1 0

5 7 7 6$_v$

0 1 1 0

$i_w$ = rank$_0$(v, $i_v$)

1 2 1 2 2 1 2

0 1 0 1 1 0 1

3 3 4 3 3 3 3

0 0 1 0 0 0 0

5 6 $_w$

0 1

7 7

1 1 1

2 2 2 2

3 3 3 3 3 3

4

5

6

$i_3$ = rank$_0$(w, $i_w$) = 1

# select$_c$(S,j)

Start at position $j$ in leaf corresponding to $c$.
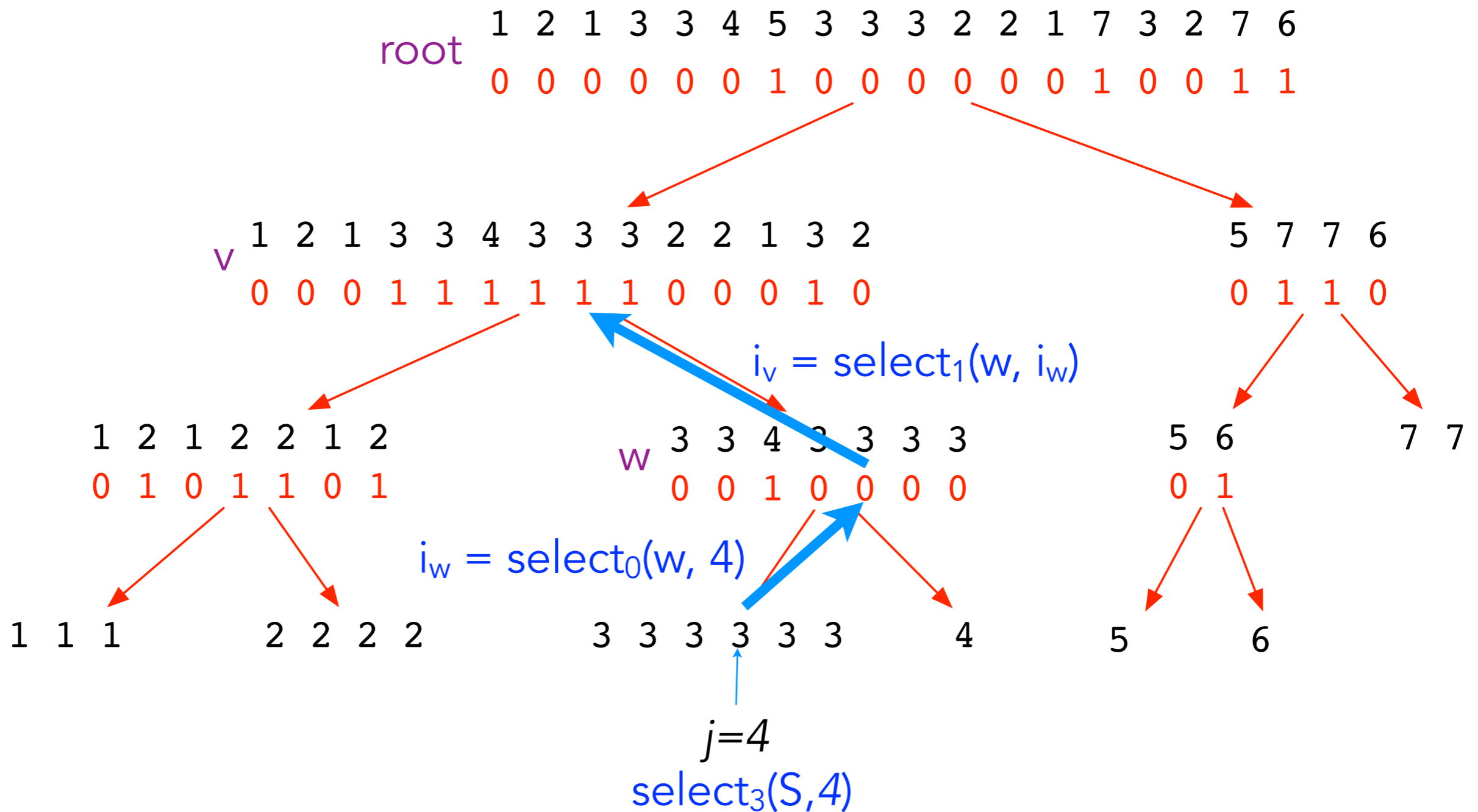
Repeat: new $i$ = select$_b$($p$, $i$) where $p$ = parent of current node, $b$ = 0 if current node is left child of $p$, 1 if right child

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

v
1 2 1 3 3 4 3 3 3 2 2 1 3 2
0 0 0 1 1 1 1 1 1 0 0 0 1 0

5 7 7 6
0 1 1 0

1 2 1 2 2 1 2
0 1 0 1 1 0 1

w
3 3 4 3 3 3 3
0 0 1 0 0 0 0

5 6
0 1

7 7

1 1 1

2 2 2 2

3 3 3 3 3 3

4

5

6

$j=4$

select$_3$(S,4)

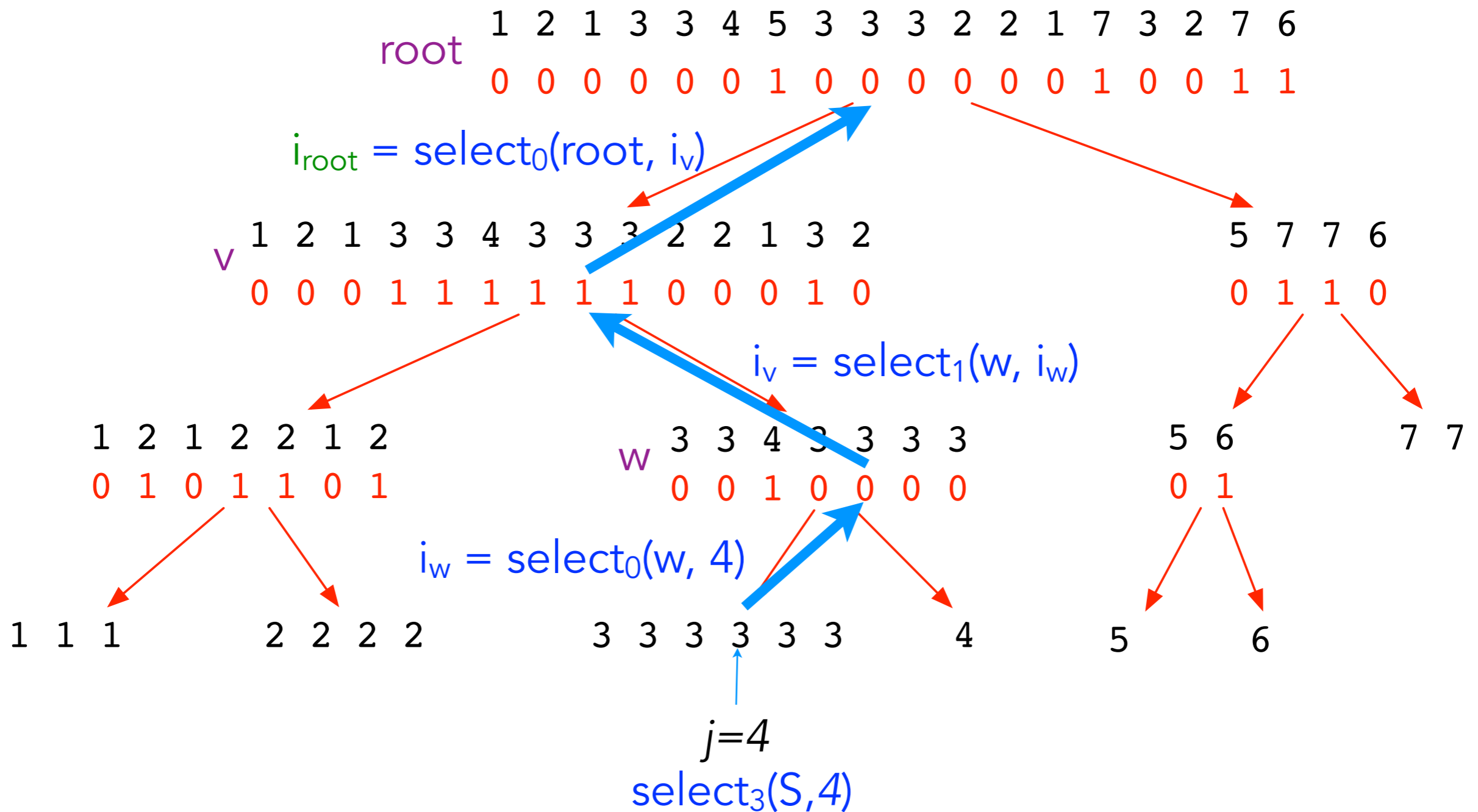# select$_c$(S,j)

Start at position *j* in leaf corresponding to *c*.

Repeat: new *i* = select$_b$(*p, i*) where *p* = parent of current node, *b* = 0 if current node is left child of *p*, 1 if right child

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

v
1 2 1 3 3 4 3 3 3 2 2 1 3 2
0 0 0 1 1 1 1 1 1 0 0 0 1 0

5 7 7 6
0 1 1 0

1 2 1 2 2 1 2
0 1 0 1 1 0 1

w
3 3 4 3 3 3 3
0 0 1 0 0 0 0

5 6
0 1

7 7

1 1 1

2 2 2 2

$i_w$ = select$_0$(w, 4)

3 3 3 3 3 3

4

5

6

*j=4*

select$_3$(S,4)

# select$_c$(S,j)

Start at position *j* in leaf corresponding to *c*.

Repeat: new *i* = select$_b$(*p, i*) where *p* = parent of current node, *b* = 0 if current node is left child of *p*, 1 if right child

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

v
1 2 1 3 3 4 3 3 3 2 2 1 3 2
0 0 0 1 1 1 1 1 1 0 0 0 1 0

5 7 7 6
0 1 1 0

$i_v$ = select$_1$(w, $i_w$)

1 2 1 2 2 1 2
0 1 0 1 1 0 1

w
3 3 4 3 3 3 3
0 0 1 0 0 0 0

5 6
0 1

7 7

$i_w$ = select$_0$(w, 4)

1 1 1      2 2 2 2

3 3 3 3 3 3      4

5      6

*j=4*

select$_3$(S,4)

# select_c(S,j)

Start at position *j* in leaf corresponding to *c*.

Repeat: new $i = select_b(p, i)$ where $p$ = parent of current node, $b = 0$ if current node is left child of *p*, 1 if right child



root

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

$i_{root} = select_0(root, i_v)$

v

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

$i_v = select_1(w, i_w)$

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

$i_w = select_0(w, 4)$

1 1 1      2 2 2 2      3 3 3 3 3 3      4      5      6

*j=4*

select_3(S,4)

# Running Times

Tree height = log |$\sum$|, where $\sum$ is the alphabet.

rank, select, access follow a root-leaf path in the tree, taking O(1) time at each node.

Therefore, rank, select, access take O(log |$\sum$|) to run.

If alphabet size is constant, this is O(1).

# Tree Shape

- Don't have to use balanced tree shape

- Can instead encode using, say, Huffman code tree shape
  - makes accesses to frequent characters faster
  - gives good space usage even without compressed bit vectors.

- Doesn't have to be binary tree
  - by selecting optimal branching factor, can get query time to O(log n / log log n).

# Application: Inverted Indices

```
 1   If you really want to hear about it, the first thing you'll
13   probably want to know is where I was born, and what my lousy
26   childhood was like, and how my parents were occupied and all
37   before they had me, and all that David Copperfield kind of crap,
49   but I don't feel like going into it, if you want to know the
63   truth.
```

Traditional inverted index represents the document as an array of lists:

| | |
|---|---|
| If | 1, 57 |
| you | 2, 58 |
| really | 3 |
| want | 4,14,59 |
| and | 22,29,35,41 |
| know | 16,61 |
| truth | 63 |

⋮

Good for searching for word $w$

**Hard to compute S[i]**

# Application: Inverted Indices, 2

- Represent text as a string of word ids.

- Store as a wavelet tree.

- S[$i$] now O(log |$\sum$|) = O(log $n$)

- select$_w$(S, $j$) now gives the position of the jth occurrence of word $w$ in time O(log $n$).

# Application: Document Retrieval

Given a collection of documents $D_1, ..., D_m$, answer the following types of queries quickly:

- In which documents does word w appear?

Represent documents as strings of word ids.

Concatenate the documents together, separated by a $ word that does not occur elsewhere.



```
i = 1
repeat:
    pᵢ = selectc(i)          # ith occurrence
    dᵢ = rank$(pᵢ) + 1       # document containing it
    print dᵢ
    p' = select$(dᵢ)         # end of document
    i = rankc(p') + 1        # find 1st occurrence after end of doc
```
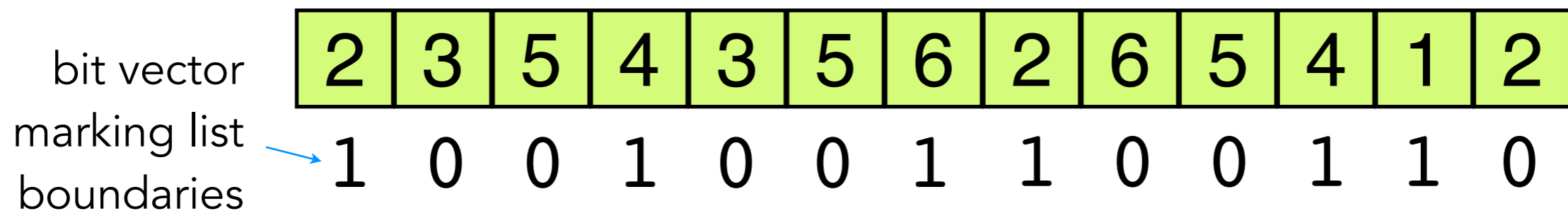
# Application: Graphs

Given a directed graph G, answer the following queries quickly:

- successor(u, i) := the ith vertex v such that edge (v,u) exists.

- predecessor(u, i) := the jth vertex v such that edge (u,v) exists.

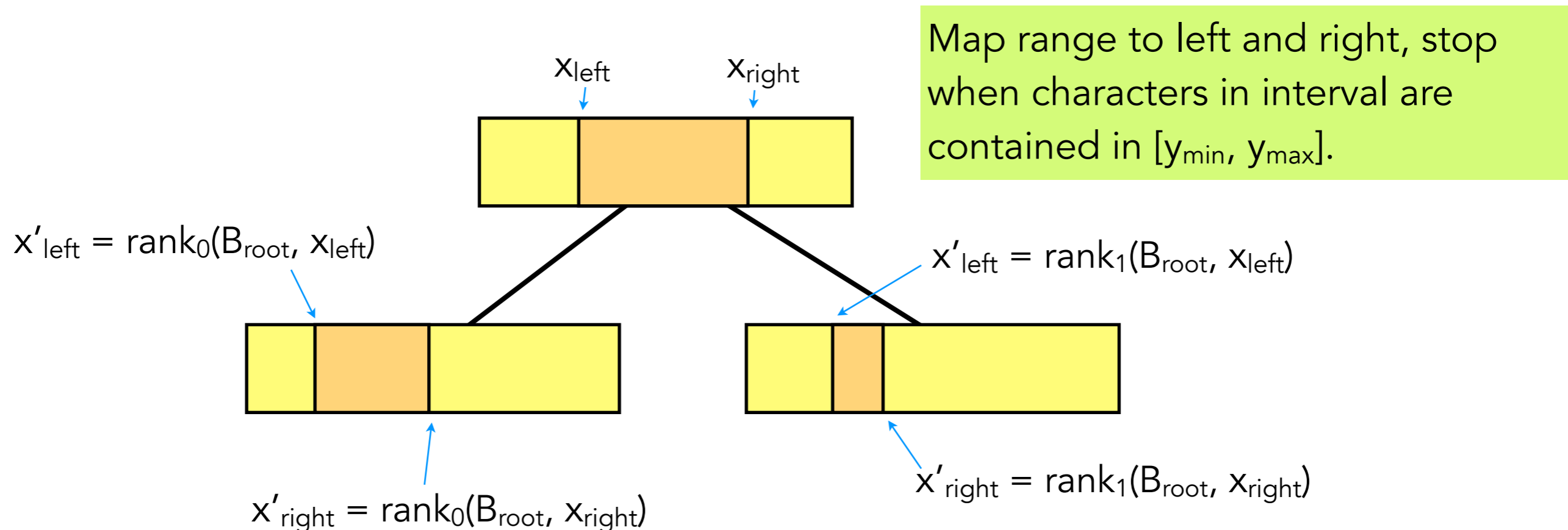Represent G as a concatenation of adjacencies lists, and store in wavelet tree:

bit vector marking list boundaries →

| 2 | 3 | 5 | 4 | 3 | 5 | 6 | 2 | 6 | 5 | 4 | 1 | 2 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|

1 0 0 1 0 0 1 1 0 0 1 1 0

- successor(u, i): $p$ = select$_1$(B,$u$); return S[$p+i$-1]

- predecessor(u, i): $p$ = select$_u$(S, $i$); return rank$_1$(B, $p$)

$|\sum| = n$, so these take $O(\log n)$ time.

# Application: Grid of Points

Suppose you have points $(x_1, y_1), ..., (x_n, y_n)$ on an m × m grid and you want to answer range queries quickly:

- Which points fall inside rectangle $[x_{min}, x_{max}]$ by $[y_{min}, y_{max}]$?

Sort points by x-coordinate, consider $y_{\pi(1)}, y_{\pi(4)}, y_{\pi(3)}, ..., y_{\pi(n)}$ as a string and store in a wavelet tree.



Map range to left and right, stop when characters in interval are contained in $[y_{min}, y_{max}]$.

$x_{left}$    $x_{right}$

$x'_{left} = rank_0(B_{root}, x_{left})$

$x'_{left} = rank_1(B_{root}, x_{left})$

$x'_{right} = rank_0(B_{root}, x_{right})$

$x'_{right} = rank_1(B_{root}, x_{right})$

# Summary

- Wavelet trees compactly store strings.

- Allowing access almost as fast as for a plain array.

- And allowing for fast rank and select queries too.

- Lots of applications and extensions, often storing things not normally thought of as strings.