# Wheeler graphs, part 5: Data structures

Ben Langmead

**JOHNS HOPKINS**
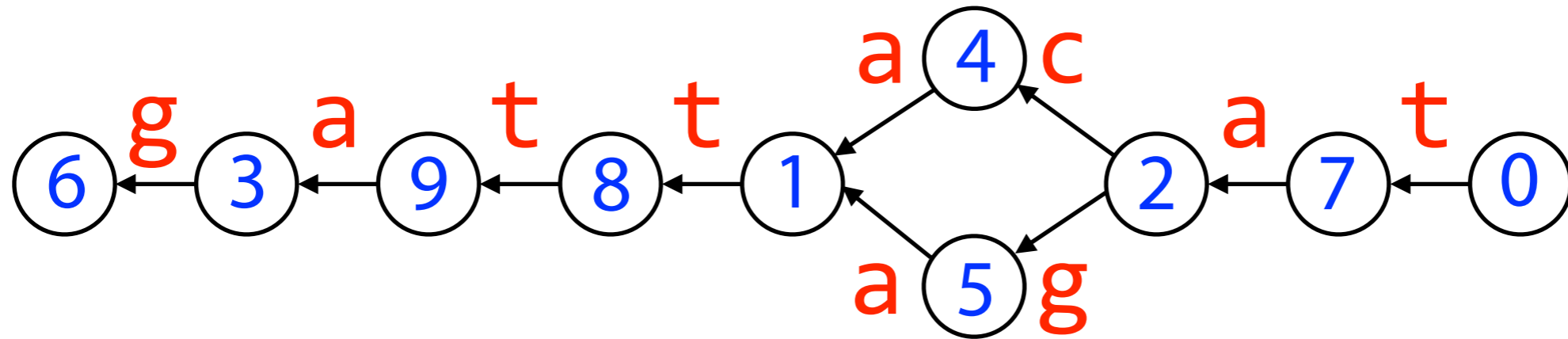WHITING SCHOOL
*of* ENGINEERING

Department of Computer Science
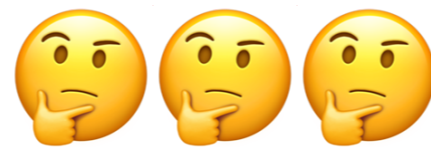
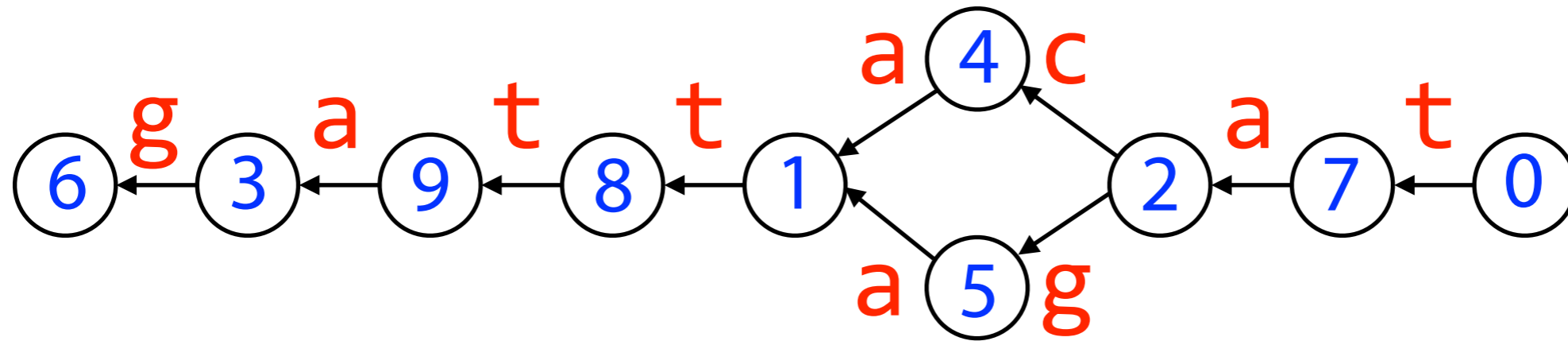# Wheeler graphs



Can we represent a Wheeler graph with **bitvectors**?

🤔🤔🤔

# Wheeler graphs



Idea 1: Encode in- and outdegree of each node in unary

| # | Unary |
|---|-------|
|   |       |

# Wheeler graphs



Idea 1: Encode in- and outdegree of each node in unary

| # | Unary |
|---|-------|
| 0 | 1 |
| 1 | 01 |
| 2 | 001 |
| 3 | 0001 |

# Wheeler graphs



Idea 1: Encode in- and outdegree
of each node in unary

Idea 2: Concatenate in order by node

Outdegree bitvector $O =$

| # | Unary |
|---|-------|
| 0 | 1 |
| 1 | 01 |
| 2 | 001 |
| 3 | 0001 |

# Wheeler graphs



Idea 1: Encode in- and outdegree of each node in unary

Idea 2: Concatenate in order by node

| # | Unary |
|---|-------|
| 0 | 1 |
| 1 | 01 |
| 2 | 001 |
| 3 | 0001 |

Outdegree bitvector $O = $ `0101001010101011010101`

Nodes:

# Wheeler graphs



Idea 1: Encode in- and outdegree of each node in unary

Idea 2: Concatenate in order by node

| # | Unary |
|---|-------|
| 0 | 1 |
| 1 | 01 |
| 2 | 001 |
| 3 | 0001 |

Outdegree bitvector $O = $ 01010010101010110101 01

Nodes:  0  1  2   3  4  5  67  8  9

# Wheeler graphs



$I =$

# Wheeler graphs



$$I = 1001010101010101010101$$

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

# Wheeler graphs



Idea 3: Encode edge labels corresponding to 0s in $O$

$$O = \texttt{010100010101011010101}$$

$$L =$$

# Wheeler graphs



Idea 3: Encode edge labels corresponding to 0s in $O$

$$O = 010100010101011010101$$

$$L = \text{t t cg g a a a t a}$$

# Wheeler graphs



Idea 3: Encode edge labels corresponding to 0s in $O$

$$O = 010100010101011010101$$

$$
\begin{array}{cccccccccc}
| & | & | & | & | & | & | & | & | & | \\
t & t & c & g & g & a & a & a & t & a
\end{array}
$$

$$L = \text{ttcggaaata}$$

# Wheeler graphs



$I = $ `10010101010101010101`

$O = $ `01010010101010110101`

$L = $ ttcggaaata

How long is $I$?    (# edges) + (# nodes) bits

How long is $O$?    (# edges) + (# nodes) bits

How long is $L$?    (# edges) chars

# Wheeler graphs



I: 100101010101010101010101

O: 010100101010110101010101

L: ttcggaaata

I: 1010101010101010101

O: 010101010110101010101

L: ttcgaata

$$BWT(T) = \text{ttcga\$ata}$$

$L$ is like BWT; $I$ & $O$ are specifically for graph structure

# Wheeler graphs



How to find indegree of node $i = 3$ ?

$$I = 1001010101010101010101$$

# Wheeler graphs



How to find indegree of node $i = 3$ ?

$$I = 1001010\underline{0}1010101010101$$

$$I.\text{select}_1(3) - I.\text{select}_1(2) - 1$$

Similar for outdegree

# Wheeler graphs



How to get labels of edges outgoing from node $i = 2$ ?

$$L = \texttt{ttcggaaata}$$
$$O = \texttt{0101001010101101010101}$$

# Wheeler graphs



How to get labels of edges outgoing from node $i = 2$ ?

$$L = \texttt{ttcggaaata}$$
$$O = \texttt{01010010101011010101}$$

# Wheeler graphs



How to get labels of edges outgoing from node $i = 2$ ?

$$L = \text{ttcggaaata}$$

$$O = \text{0101001010101011010101}$$

$$O \cdot \text{select}_1(1) = 3 \quad O \cdot \text{select}_1(2) = 6$$

# Wheeler graphs



How to get labels of edges outgoing from node $i = 2$ ?

$L = $ ttcggaaata

$O = $ 010100010101011010101

$O . \text{select}_1(1) = 3$     $O . \text{select}_1(2) = 6$

Extract $L[\text{rank}_0(3) = 2 : \text{rank}_0(6) = 4] = $ cg

# Wheeler graphs



How do we use these bitvectors for *matching*?

🤔🤔🤔

# Wheeler graphs

P = aba

**FM Index** match query loop:

| F | | | | | | L |
|---|---|---|---|---|---|---|
| $ | a | b | a | a | b | $a_0$ |
| $a_0$ | $ | a | b | a | a | $b_0$ |
| $a_1$ | a | b | a | $ | a | $b_1$ |
| $a_2$ | b | a | $ | a | b | $a_1$ |
| $a_3$ | b | a | a | b | a | $ |
| $b_0$ | a | $ | a | b | a | $a_2$ |
| $b_1$ | a | a | b | a | $ | $a_3$ |

*Skip:* $C[c]$          *+LF:* $L.\mathrm{rank}_c(\ldots)$

# Wheeler graphs

P = ab$a$

**FM Index** match query loop:

| F | | | | | | L |
|---|---|---|---|---|---|---|
| **\$** | a | b | a | a | b | **a₀** |
| **a₀** | \$ | a | b | a | a | **b₀** |
| **a₁** | a | b | a | \$ | a | **b₁** |
| **a₂** | b | a | \$ | a | b | **a₁** |
| **a₃** | b | a | a | b | a | **\$** |
| **b₀** | a | \$ | a | b | a | **a₂** |
| **b₁** | a | a | b | a | \$ | **a₃** |

First $c$ = a

*Skip:* $C[\text{c}]$

$+LF:$ $L.\text{rank}_c(\dots)$

# Wheeler graphs

P = a$\textcolor{orange}{ba}$

**FM Index** match query loop:

$$F \qquad\qquad\qquad\qquad\qquad L$$

| | | | | | | |
|---|---|---|---|---|---|---|
| $\$$ | a | b | a | a | b | $a_0$ |
| $a_0$ | $\$$ | a | b | a | a | $b_0$ |
| $a_1$ | a | b | a | $\$$ | a | $b_1$ |
| $a_2$ | b | a | $\$$ | a | b | $a_1$ |
| $a_3$ | b | a | a | b | a | $\$$ |
| $b_0$ | a | $\$$ | a | b | a | $a_2$ |
| $b_1$ | a | a | b | a | $\$$ | $a_3$ |

*Skip:* $C[\text{c}]$

$+ LF:$ $L.\text{rank}_{\text{c}}(\ldots)$

Next $c = \textcolor{orange}{b}$

# Wheeler graphs

$F$        $L$

P = aba

**FM Index** match query loop:

| $F$ | | | | | | $L$ |
|---|---|---|---|---|---|---|
| **\$** | a | b | a | a | b | $a_0$ |
| $a_0$ | \$ | a | b | a | a | $b_0$ |
| $a_1$ | a | b | a | \$ | a | $b_1$ |
| $a_2$ | b | \$ | a | b | | $a_1$ |
| $a_3$ | b | a | a | b | a | \$ |
| $b_0$ | a | \$ | a | b | a | $a_2$ |
| $b_1$ | a | a | b | a | \$ | $a_3$ |

*Skip:* $C[c]$

*+LF:* $L \, . \, \mathrm{rank}_c(\ldots)$

# Wheeler graphs

P = aba

**FM Index** match query loop:

|     | F |   |   |   |   | L |
|-----|---|---|---|---|---|---|
| $ | a | b | a | a | b | $a_0$ |
| $a_0$ | $ | a | b | a | a | $b_0$ |
| $a_1$ | a | b | a | $ | a | $b_1$ |
| $a_2$ | b | a | $ | a | b | $a_1$ |
| $a_3$ | b | a | a | b | a | $ |
| $b_0$ | a | $ | a | b | a | $a_2$ |
| $b_1$ | a | a | b | a | $ | $a_3$ |

*Skip:* $C[\text{c}]$

$+ LF:$ $L . \text{rank}_\text{c}( \ldots )$

Next $c =$ a

# Wheeler graphs

$$F \qquad\qquad L$$

P = aba

| | | | | | | |
|---|---|---|---|---|---|---|
| **$**  | a | b | a | a | b | **$a_0$** |
| **$a_0$** | $ | a | b | a | a | **$b_0$** |
| **$a_1$** | a | b | a | $ | a | **$b_1$** |
| **$a_2$** | b | a | $ | a | b | **$a_1$** |
| **$a_3$** | b | a | a | b | a | **$** |
| **$b_0$** | a | $ | a | b | a | **$a_2$** |
| **$b_1$** | a | a | b | a | $ | **$a_3$** |

**FM Index** match query loop:

*Skip:* $C[c]$

$+LF:$ $L\,.\,\mathrm{rank}_c(\ldots)$

# Wheeler graphs



**Wheeler graph** match
query loop:

$I$ : 10010101010101010101

$O$ : 01010010101011010101

$L$ : ttcggaaata

*Find range of characters in $L$:*

$$L \cdot \mathrm{rank}_c(\dots)$$

*Skip:* $C[c]$

*Find outgoing edges in $O$:*

$$O \cdot \mathrm{rank}_0(O \cdot \mathrm{select}_1(\dots))$$

*Follow incoming edges in $I$:*

$$I \cdot \mathrm{rank}_1(I \cdot \mathrm{select}_0(\dots))$$

# Wheeler graphs

P = aga



F: aaaacggttt

I: 100101010101010101 0101

O: 010100101010111010101

L: ttcggaaata

$$L \cdot \text{rank}_c(\dots)$$

$$C[c]$$

$$O \cdot \text{rank}_0(O \cdot \text{select}_1(\dots))$$

$$I \cdot \text{rank}_1(I \cdot \text{select}_0(\dots))$$

# Wheeler graphs



P = ag**a**

F: aaaacggttt

I: 10010101010101010101

O: 01010010101011010101

L: ttcggaaata

First $c$ = a

$$L . \mathrm{rank}_c( \ldots )$$

$$C[\mathrm{c}]$$

$$O . \mathrm{rank}_0(O . \mathrm{select}_1( \ldots ))$$

$$I . \mathrm{rank}_1(I . \mathrm{select}_0( \ldots ))$$

# Wheeler graphs

P = ag**a**



F: aaaacggttt

I: 100101010101010101010101

O: 010100101010110101010101

L: ttcggaaata

$$L . \mathrm{rank}_c( \dots )$$

$$C[\mathrm{c}]$$

$$O . \mathrm{rank}_0(O . \mathrm{select}_1( \dots ))$$

$$I . \mathrm{rank}_1(I . \mathrm{select}_0( \dots ))$$

# Wheeler graphs

P = ag**a**



F: aaaacggttt

I: 10010101010101010101

O: 01010010101011010101

L: ttcggaaata

$$L . \mathrm{rank}_c( \dots )$$

$$C[\mathtt{c}]$$

$$O . \mathrm{rank}_0( O . \mathrm{select}_1( \dots ))$$

$$I . \mathrm{rank}_1(I . \mathrm{select}_0( \dots ))$$

# Wheeler graphs



P = a*ga*

F: aaaacggttt

I: 100101010101010101010101

O: 0101001010101011010101

L: ttcggaaata

Next $c = $ g

$$L . \text{rank}_c( \ldots )$$

$C[c]$

$O . \text{rank}_0(O . \text{select}_1( \ldots ))$

$I . \text{rank}_1(I . \text{select}_0( \ldots ))$

# Wheeler graphs

P = a<span style="color:orange">ga</span>



F: aaaac<span style="color:magenta">gg</span>ttt
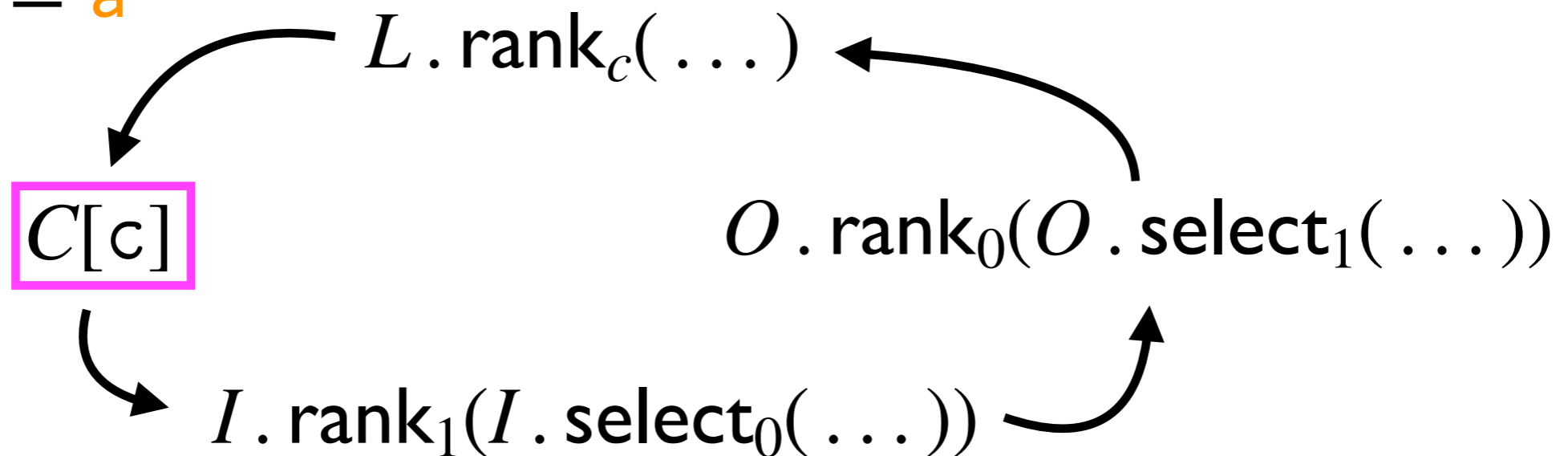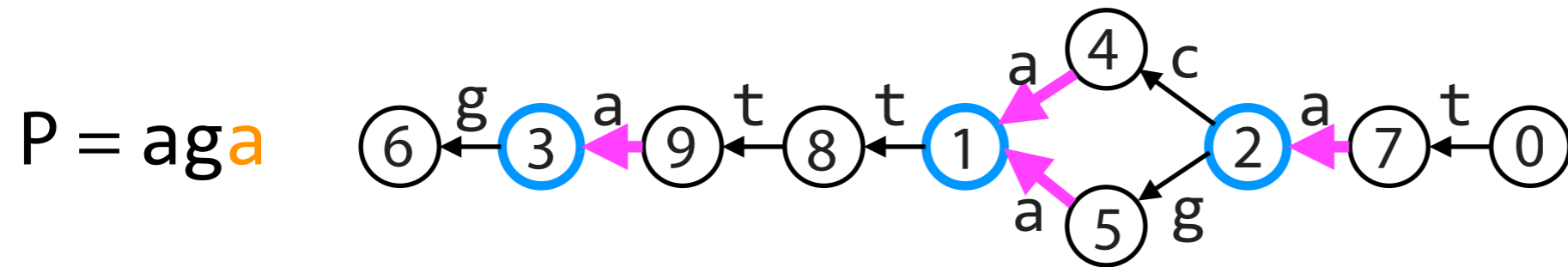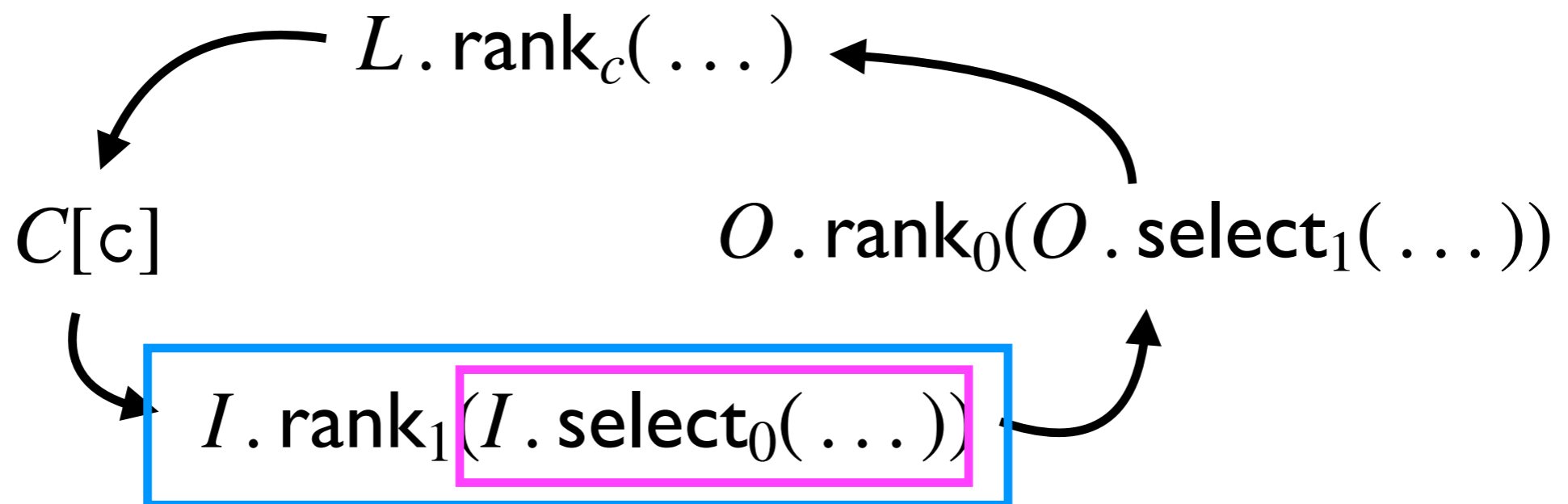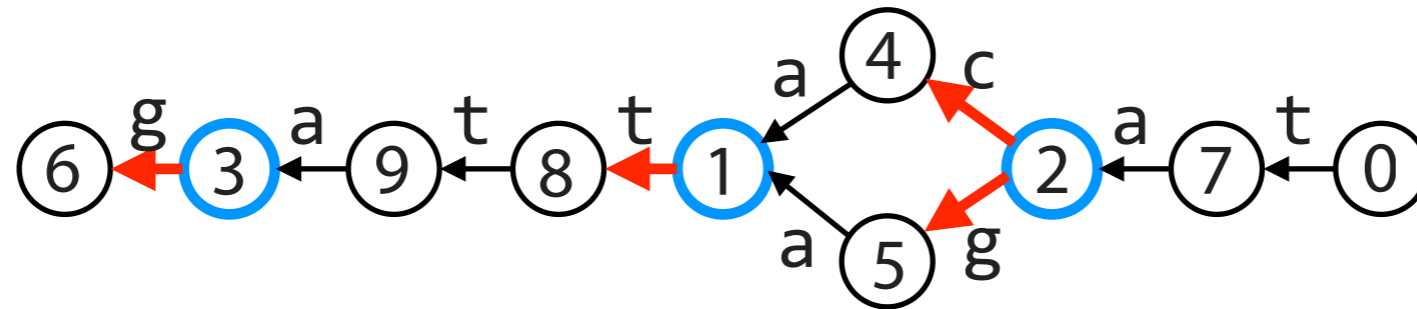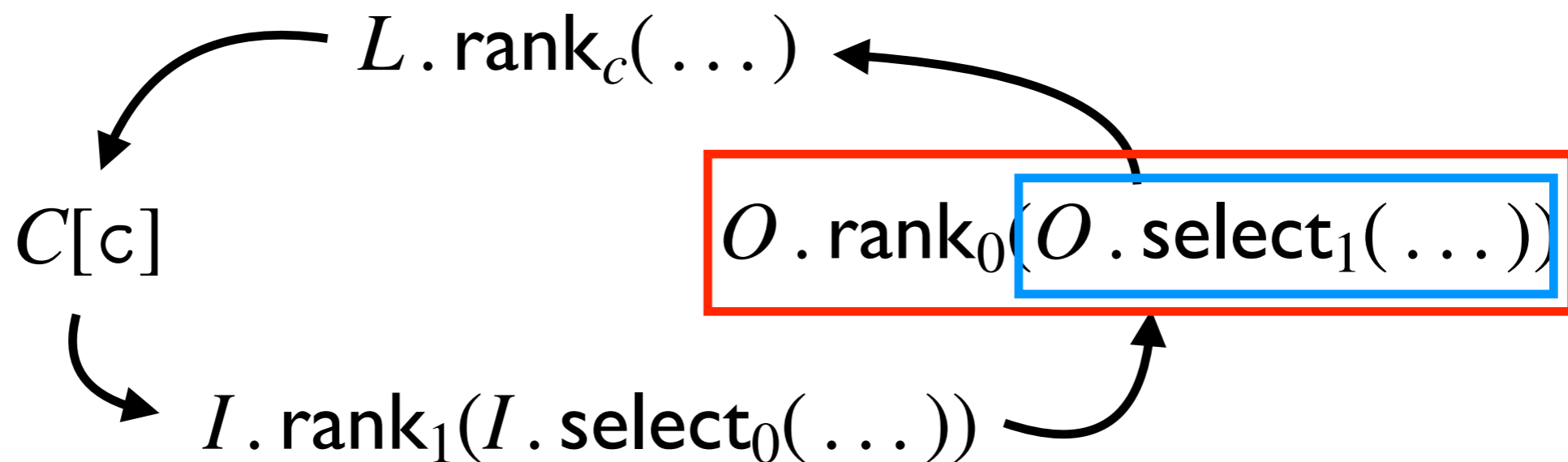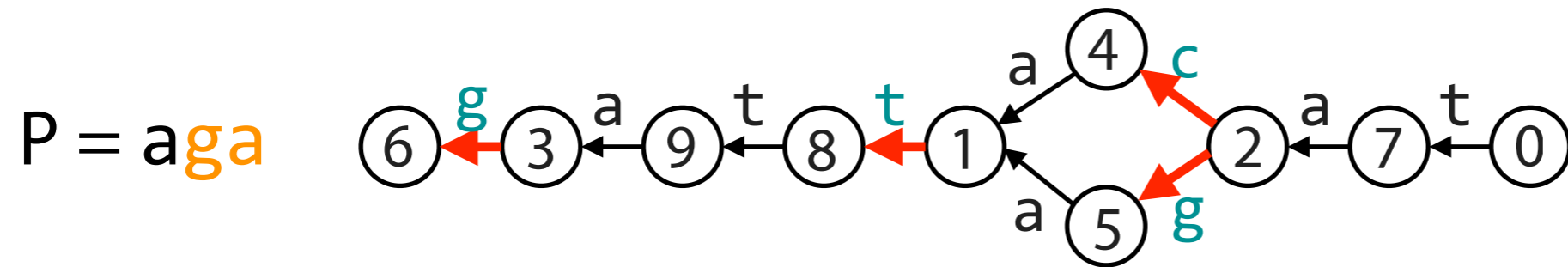
I: 10010101010101010101

O: 01010010101011010101

L: <span style="color:teal">ttcgg</span>aaata

$$L . \mathrm{rank}_c( \ldots )$$

$$C[\mathrm{c}]$$

$$O . \mathrm{rank}_0(O . \mathrm{select}_1( \ldots ))$$

$$I . \mathrm{rank}_1(I . \mathrm{select}_0( \ldots ))$$

# Wheeler graphs



P = a**ga**

F: aaaac**gg**ttt

I: 100101010101010101010101

O: 010100101010101011010101

L: ttcggaaata

$$L \cdot \text{rank}_c(\dots)$$

$$C[\text{c}]$$

$$O \cdot \text{rank}_0(O \cdot \text{select}_1(\dots))$$

$$I \cdot \text{rank}_1(I \cdot \text{select}_0(\dots))$$

# Wheeler graphs



$P = \text{a}\textcolor{orange}{\text{ga}}$

F: aaaacggttt

I: 10010101010101010101

O: 01010010101011010101

L: ttcggaaata

$$L . \text{rank}_c( \dots )$$

$$C[\text{c}]$$

$$I . \text{rank}_1(I . \text{select}_0( \dots ))$$

$$O . \text{rank}_0(O . \text{select}_1( \dots ))$$

# Wheeler graphs

P = aga



F: aaaacggttt

I: 10010101010101010101

O: 0101001010101011010101

L: ttcggaaata

Next $c$ = a

$L . \operatorname{rank}_c( \dots )$

$C[\text{c}]$

$I . \operatorname{rank}_1(I . \operatorname{select}_0( \dots ))$

$O . \operatorname{rank}_0(O . \operatorname{select}_1( \dots ))$

# Wheeler graphs

P = aga



F: aaaacggttt

I: 10010101010101010101

O: 01010010101011010101

L: ttcggaaata

$$L.\text{rank}_c(\ldots)$$

$$O.\text{rank}_0(O.\text{select}_1(\ldots))$$

$$C[c]$$

$$I.\text{rank}_1(I.\text{select}_0(\ldots))$$

# Wheeler graphs

P = aga



Final answer: 1 match, corresponding to this path

$$L . \mathrm{rank}_c( \dots )$$

$$C[c]$$

$$O . \mathrm{rank}_0(O . \mathrm{select}_1( \dots ))$$

$$I . \mathrm{rank}_1(I . \mathrm{select}_0( \dots ))$$

# Wheeler graph

Given character $c$ & next character $c'$, one step of matching process:

FM index:                    Wheeler graph:

$r_{top} \leftarrow C[c]$

$r_{bot} \leftarrow C[c+1]$

$m_{top} \leftarrow BWT(T).\mathsf{rank}_c(r_{top}, c')$

$m_{top} \leftarrow BWT(T).\mathsf{rank}_c(r_{bot}, c')$

$r_{top} \leftarrow C[c]$                        $r_{bot} \leftarrow C[c+1]$

$i_{top} \leftarrow I.\mathsf{select}_0(r_{top})$        $i_{bot} \leftarrow I.\mathsf{select}_0(r_{bot}-1)$

$j_{top} \leftarrow I.\mathsf{rank}_1(i_{top})$        $j_{bot} \leftarrow I.\mathsf{rank}_1(i_{bot})+1$

$k_{top} \leftarrow O.\mathsf{select}_1(j_{top}-1)$    $k_{bot} \leftarrow O.\mathsf{select}_1(j_{bot}-1)$

$\ell_{top} \leftarrow O.\mathsf{rank}_0(k_{top})$        $\ell_{bot} \leftarrow O.\mathsf{rank}_0(k_{bot})$

$m_{top} \leftarrow S.\mathsf{rank}_c(\ell_{top}, c')$    $m_{bot} \leftarrow S.\mathsf{rank}_c(\ell_{bot}, c')$

# Wheeler graph

## What takes **space**?

### FM index:

$$r_{top} \leftarrow C[c]$$
$$r_{bot} \leftarrow C[c+1]$$
$$m_{top} \leftarrow BWT(T) \,.\, \text{rank}_c(r_{top}, c')$$
$$m_{top} \leftarrow BWT(T) \,.\, \text{rank}_c(r_{bot}, c')$$

### Wheeler graph:

$$r_{top} \leftarrow C[c] \qquad\qquad r_{bot} \leftarrow C[c+1]$$
$$i_{top} \leftarrow I \,.\, \text{select}_0(r_{top}) \qquad i_{bot} \leftarrow I \,.\, \text{select}_0(r_{bot} - 1)$$
$$j_{top} \leftarrow I \,.\, \text{rank}_1(i_{top}) \qquad j_{bot} \leftarrow I \,.\, \text{rank}_1(i_{bot}) + 1$$
$$k_{top} \leftarrow O \,.\, \text{select}_1(j_{top} - 1) \qquad k_{bot} \leftarrow O \,.\, \text{select}_1(j_{bot} - 1)$$
$$\ell_{top} \leftarrow O \,.\, \text{rank}_0(k_{top}) \qquad \ell_{bot} \leftarrow O \,.\, \text{rank}_0(k_{bot})$$
$$m_{top} \leftarrow L \,.\, \text{rank}_c(\ell_{top}, c') \qquad m_{bot} \leftarrow L \,.\, \text{rank}_c(\ell_{bot}, c')$$

$C$ array: $\quad \sigma \log n$

$WT(BWT)$ rank: $\quad n \log \sigma + \breve{o}(n \log \sigma)$

*(units are bits)*

# Wheeler graph

## What takes **space**?

$r_{top} \leftarrow C[c]$

$r_{bot} \leftarrow C[c+1]$

$m_{top} \leftarrow BWT(T) . \text{rank}_c(r_{top}, c')$

$m_{top} \leftarrow BWT(T) . \text{rank}_c(r_{bot}, c')$

### Wheeler graph:

$r_{top} \leftarrow C[c]$      $r_{bot} \leftarrow C[c+1]$

$i_{top} \leftarrow I . \text{select}_0(r_{top})$      $i_{bot} \leftarrow I . \text{select}_0(r_{bot} - 1)$

$j_{top} \leftarrow I . \text{rank}_1(i_{top})$      $j_{bot} \leftarrow I . \text{rank}_1(i_{bot}) + 1$

$k_{top} \leftarrow O . \text{select}_1(j_{top} - 1)$      $k_{bot} \leftarrow O . \text{select}_1(j_{bot} - 1)$

$\ell_{top} \leftarrow O . \text{rank}_0(k_{top})$      $\ell_{bot} \leftarrow O . \text{rank}_0(k_{bot})$

$m_{top} \leftarrow L . \text{rank}_c(\ell_{top}, c')$      $m_{bot} \leftarrow L . \text{rank}_c(\ell_{bot}, c')$

$C$ array: $\sigma \log |E|$

$I$ rank+select: $|E| + |N| + \breve{o}(|E| + |N|)$

$O$ rank+select: $|E| + |N| + \breve{o}(|E| + |N|)$

$WT(L)$ rank: $|E| \log \sigma + \breve{o}(|E| \log \sigma)$

*(units are bits)*

# Wheeler graph

## What takes **space**?

FM index:

$C$ array:  $\sigma \log n$

$WT(BWT)$ rank:  $n \log \sigma + \breve{o}(n \log \sigma)$

Wheeler graph:

$C$ array:  $\sigma \log |E|$

$I$ rank+select: $|E| + |N| + \breve{o}(|E| + |N|)$

$O$ rank+select: $|E| + |N| + \breve{o}(|E| + |N|)$

$WT(L)$ rank: $|E| \log \sigma + \breve{o}(|E| \log \sigma)$

# Wheeler graph

What takes **time**?

FM index:

$$r_{top} \leftarrow C[c]$$

$$r_{bot} \leftarrow C[c+1]$$

$$m_{top} \leftarrow BWT(T)\,.\,\text{rank}_c(r_{top}, c')$$

$$m_{top} \leftarrow BWT(T)\,.\,\text{rank}_c(r_{bot}, c')$$

Wheeler graph:

$$r_{top} \leftarrow C[c] \qquad\qquad r_{bot} \leftarrow C[c+1]$$

$$i_{top} \leftarrow I\,.\,\text{select}_0(r_{top}) \qquad i_{bot} \leftarrow I\,.\,\text{select}_0(r_{bot}-1)$$

$$j_{top} \leftarrow I\,.\,\text{rank}_1(i_{top}) \qquad j_{bot} \leftarrow I\,.\,\text{rank}_1(i_{bot})+1$$

$$k_{top} \leftarrow O\,.\,\text{select}_1(j_{top}-1) \qquad k_{bot} \leftarrow O\,.\,\text{select}_1(j_{bot}-1)$$

$$\ell_{top} \leftarrow O\,.\,\text{rank}_0(k_{top}) \qquad \ell_{bot} \leftarrow O\,.\,\text{rank}_0(k_{bot})$$

$$m_{top} \leftarrow L\,.\,\text{rank}_c(\ell_{top}, c') \qquad m_{bot} \leftarrow L\,.\,\text{rank}_c(\ell_{bot}, c')$$

Ranks on wavelet trees:  $O(\log \sigma)$

Ranks and selects on bitvectors:  $O(1)$