# Wheeler graphs, part 2

Ben Langmead

JOHNS HOPKINS
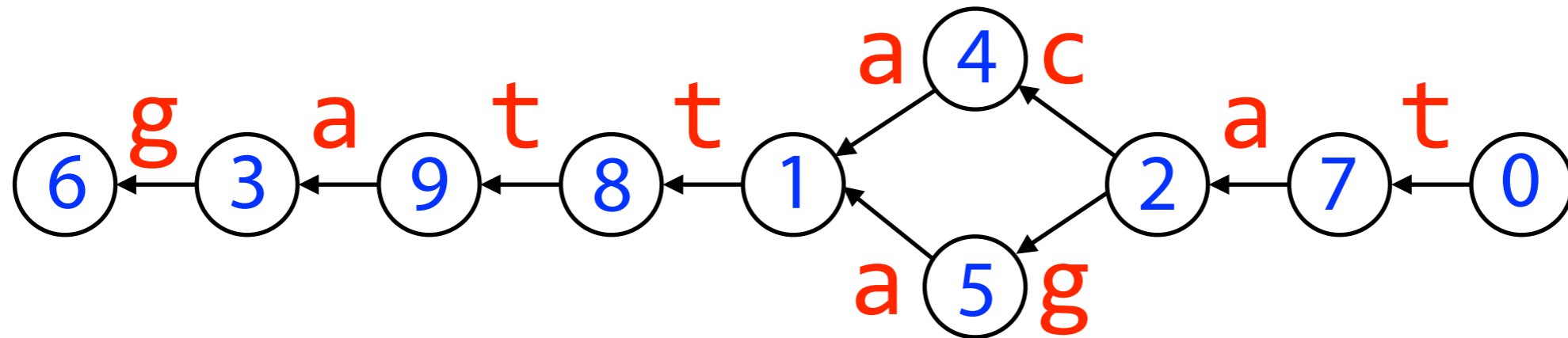WHITING SCHOOL
*of* ENGINEERING
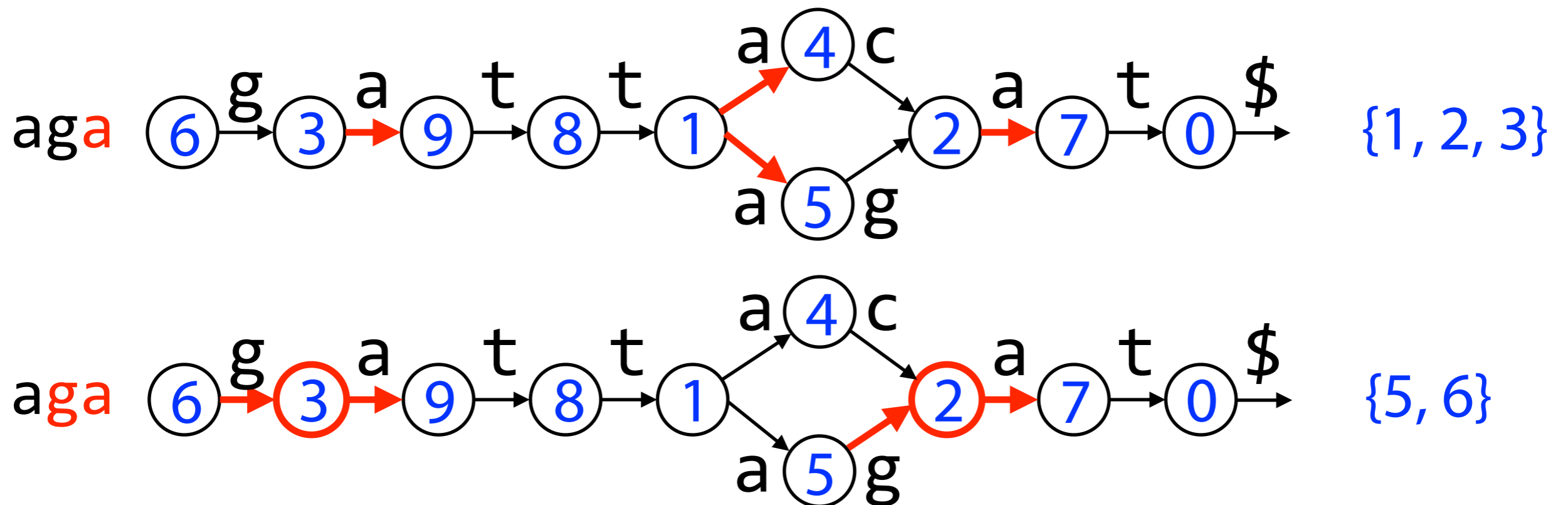
## Department of Computer Science

# Wheeler graphs



Can we show Wheeler Graphs are "special," giving us analogs to FM index queries?

We start with "path coherence," the graph version of **consecutivity**
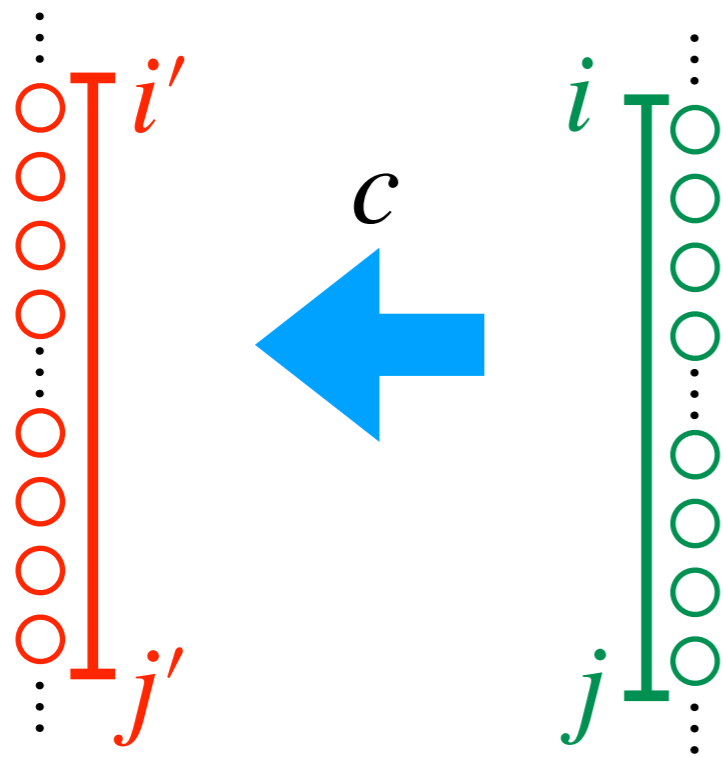
# Wheeler graphs

A graph is **path coherent** if there is a total order of the nodes such that:

For any consecutive range $[i, j]$ of nodes and string $\alpha$, the nodes reached by following edges matching $\alpha$ also form a consecutive range

# Wheeler graphs

Consider a single step where our <span style="color:green">initial set of nodes</span> are in consecutive range $[i, j]$ and, after advancing on a single character $c \in \Sigma$, $[i', j']$ is the smallest range containing our <span style="color:red">next set of nodes</span>
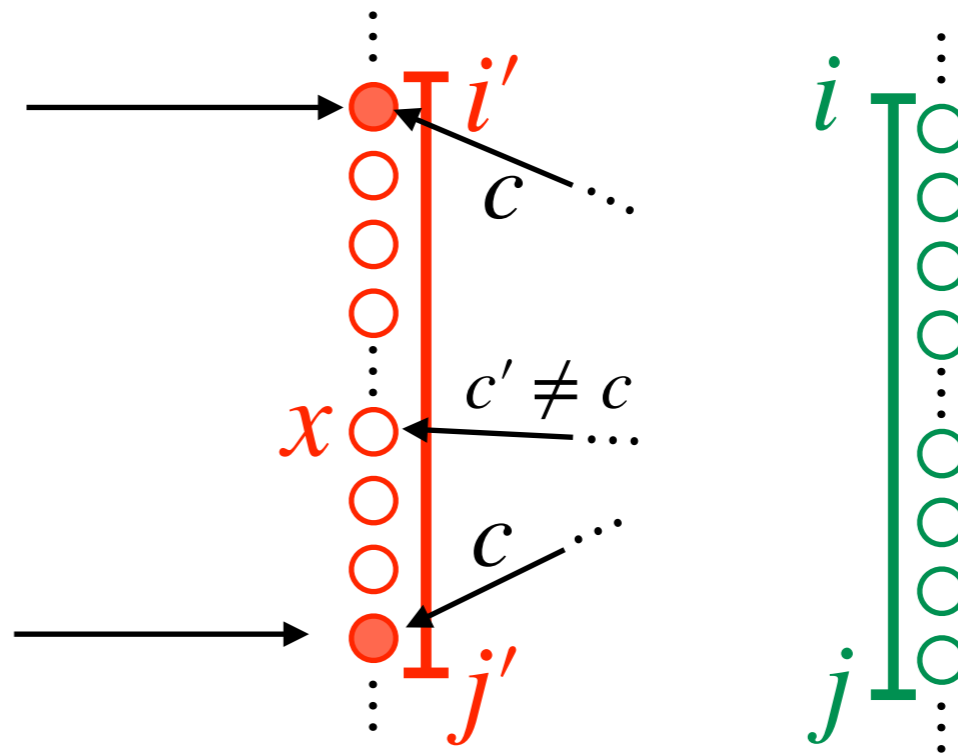


We want to show that the nodes in $[i', j']$ consist *only* of the $c$-successors of nodes $[i, j]$

# Wheeler graphs

As defined, $i'$ is reachable
via an edge labeled $c$
from a node in $[i, j]$

Same for $j'$



Consider node $x$, where $i' < x < j'$ with incoming edge
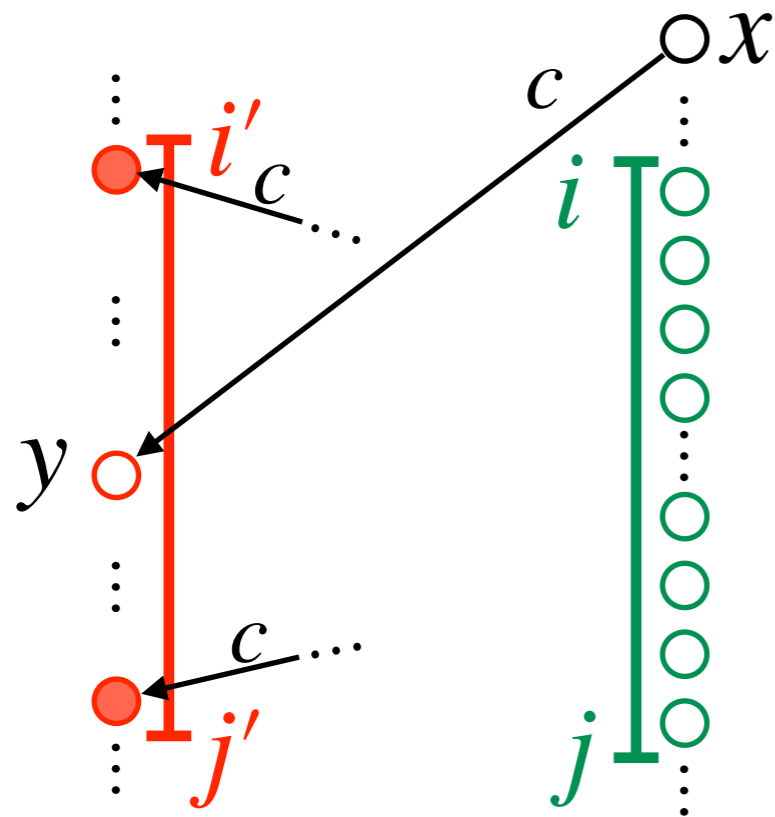labeled $c'$. Suppose $c' \neq c$.

Recall: $a \prec a' \implies v < v'$

Since $x \not\prec i'$, we have $c' \not\prec c$
Since $j' \not\prec x$, we have $c \not\prec c'$

We have $c' \succeq c$, $c \succeq c'$, and $c' \neq c$, giving a contradiction

# Wheeler graphs



Could node $x \notin [i, j]$ be a $c-$predecessor of a node $y, i' < y < j'$?

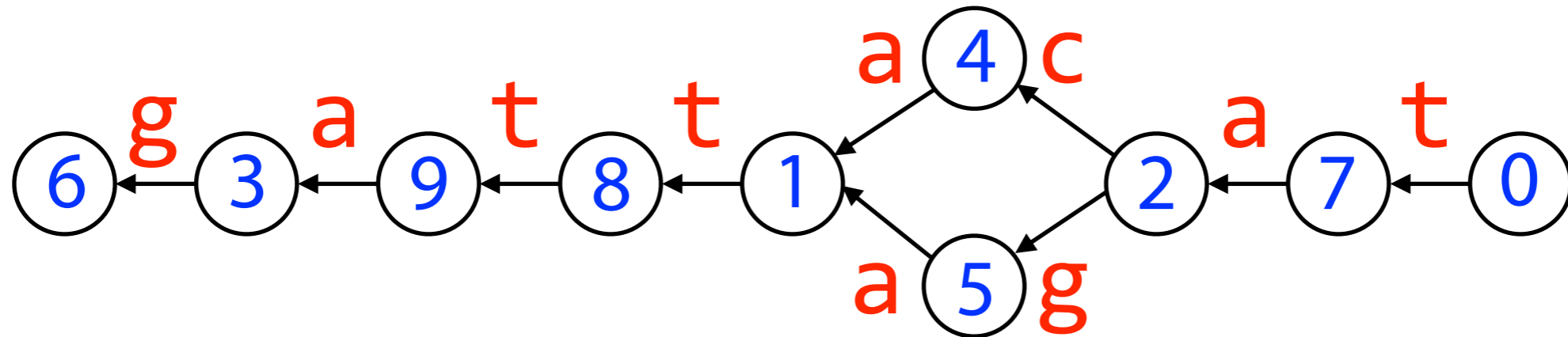**No.** Proof is by contradiction, similar to previous argument and using:

$$(a = a') \wedge (u < u') \implies v \leq v'$$

# Wheeler graphs

For any consecutive range $[i, j]$ of nodes and string $\alpha$, the nodes reached by following edges matching $\alpha$ also form a consecutive range

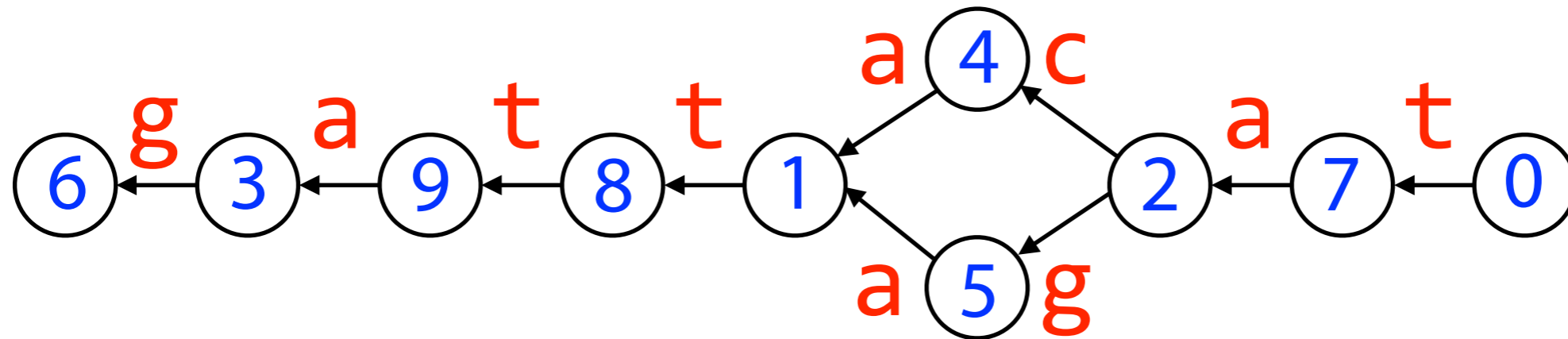Proof takes previous arguments and extends them to string $\alpha$ inductively

# Wheeler graphs



How would we represent a Wheeler graph with bitvectors?

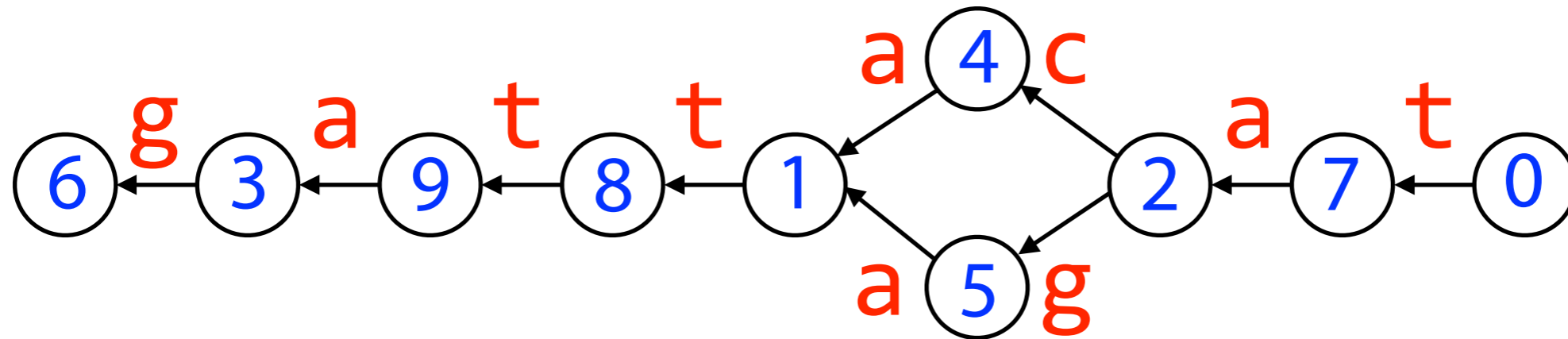Need to represent *structure* as well as *node and edge labels*

# Wheeler graphs



Idea 1: Encode in- and outdegree of each node in unary
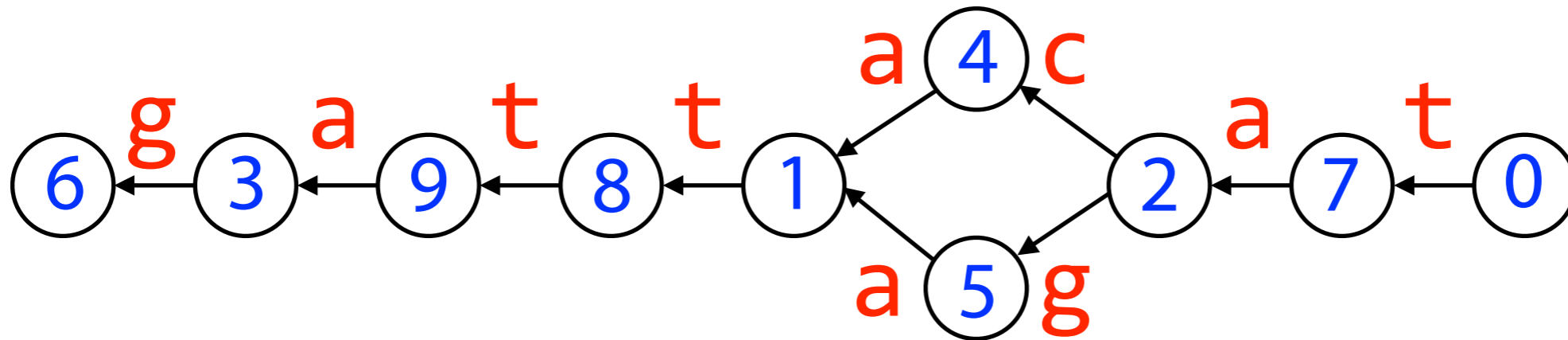
Idea 2: Concatenate in order by node

Here, $O = $ 0101001010101011010101

# Wheeler graphs



$$I = 10010101010101010101$$

0   1  2  3  4  5  6  7  8  9

# Wheeler graphs



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

$$O = \texttt{010100010101011010101}$$

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

# Wheeler graphs



$$I = \text{10010101010101010101}$$

$$O = \text{01010010101010110101011}$$

$$L = \text{ttcggaaata}$$

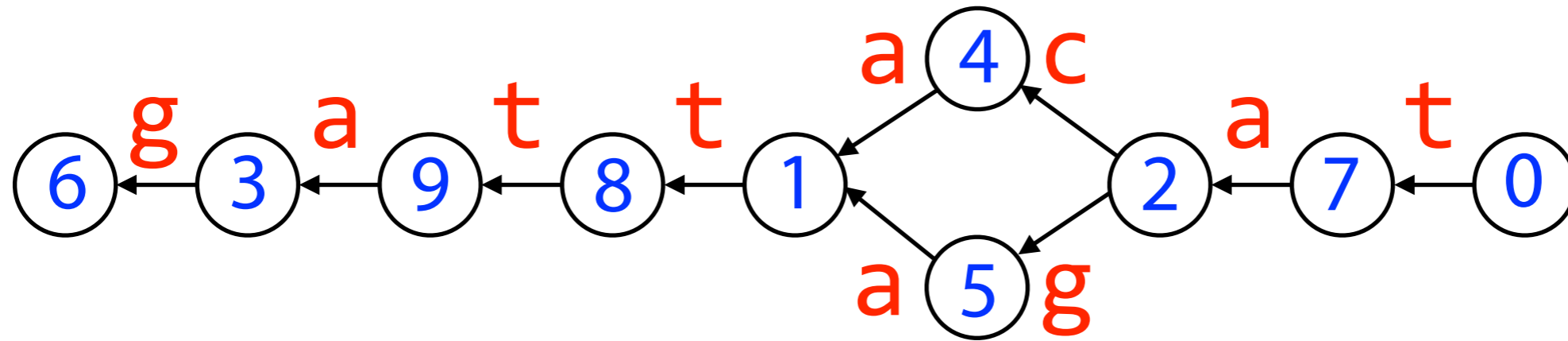| | |
|---|---|
| How long is $I$? | (# edges) + (# nodes) bits |
| How long is $O$? | (# edges) + (# nodes) bits |
| How long is $L$? | (# edges) chars |

# Wheeler graphs



$I = 10010101010101010101$
$O = 01010010101010110101$
$L = $ ttcggaaata

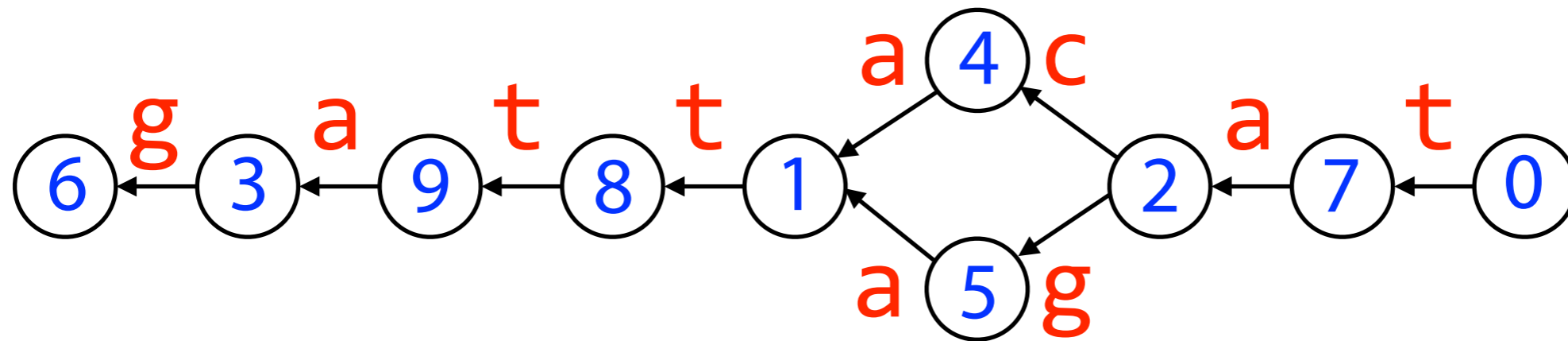How to find in- or outdegree of a node?

# Wheeler graphs



$$I = \texttt{10010101010101010101}$$

$$G.\,\text{indegree}(i) = I.\,\text{select}_1(i) - I.\,\text{select}_1(i-1) - 1$$

$$(\text{let } I.\,\text{select}(\text{-1}) = 0)$$

$$G.\,\text{indegree}(1) = I.\,\text{select}_1(1) - I.\,\text{select}_1(0) - 1$$
$$= 3 - 0 - 1 = 2$$

# Wheeler graphs



$$O = \texttt{010100101010110101}$$

$$G.\,\text{outdegree}(i) = O.\,\text{select}_1(i) - O.\,\text{select}_1(i-1) - 1$$

$$(\text{let } O.\,\text{select}(\text{-1}) = 0)$$

$$G.\,\text{outdegree}(2) = O.\,\text{select}_1(2) - O.\,\text{select}_1(1) - 1$$

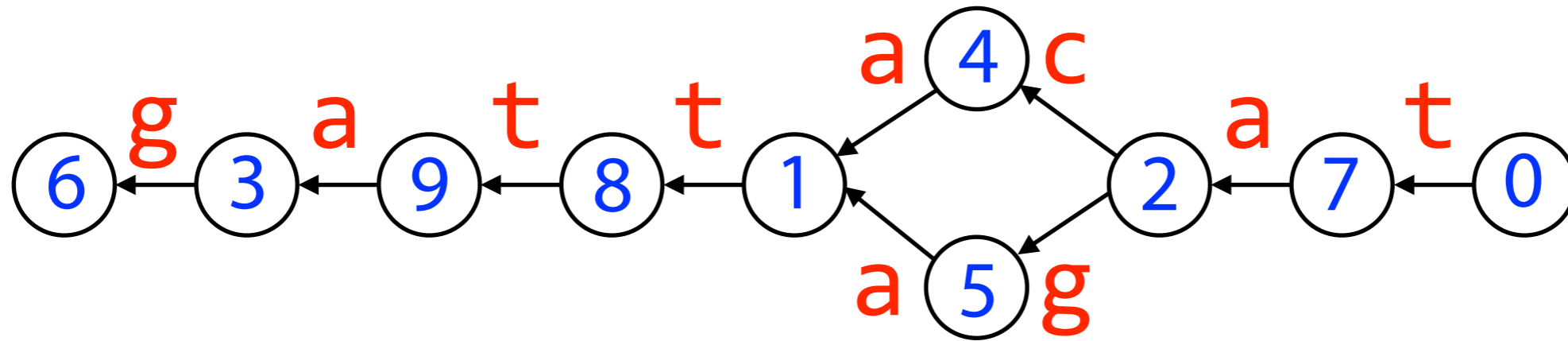$$= 6 - 3 - 1 = 2$$

# Wheeler graphs



$$I = 1001010101010101010101$$
$$O = 0101001010101011010101$$
$$L = \text{ttcggaaata}$$

How to access the labels of the edges outgoing from a node?
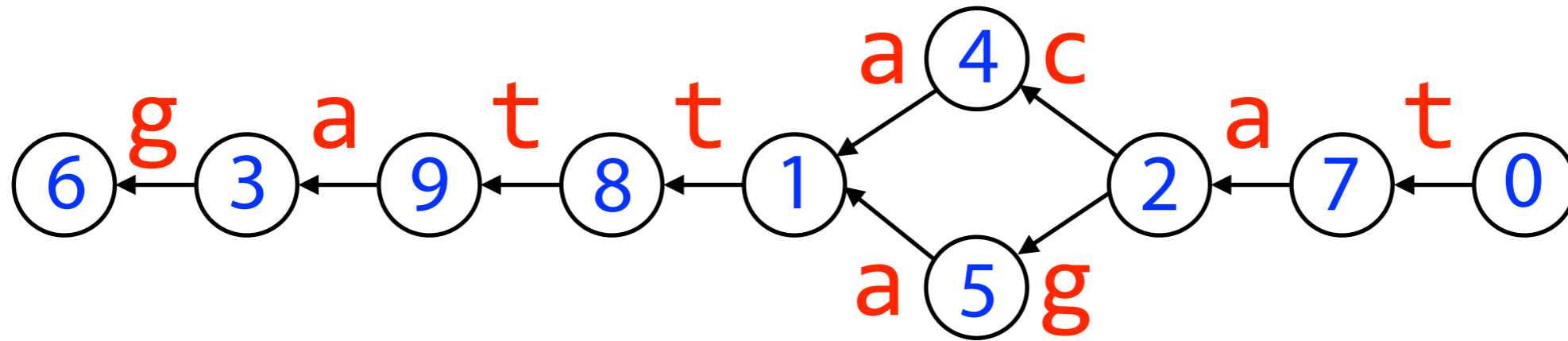
# Wheeler graphs



$$O = \texttt{010100010101011010101}$$
$$L = \texttt{ttcggaaata}$$

$$\text{off} = O \,.\, \text{rank}_0(O \,.\, \text{select}_1(i - 1))$$

# Wheeler graphs



$$O = \texttt{010100010101011010101}$$
$$L = \texttt{ttcggaaata}$$

num $= G \,.\, \mathrm{outdegree}(i)$

off $= O \,.\, \mathrm{rank}_0(O \,.\, \mathrm{select}_1(i-1))$

for $j$ in $\{0, 1, \ldots, \text{num-1}\}$:

$L \,.\, \mathrm{access}(\mathrm{off} + j)$

# Wheeler graphs



$$O = 010100010101011010101$$
$$L = \texttt{ttcggaaata}$$

num $= G\,.\,\text{outdegree}(2) = \textbf{2}$

off $= O\,.\,\text{rank}_0(O\,.\,\text{select}_1(2-1))$

for $j$ in $\{0, 1, \ldots, \text{num-1}\}$:

$\quad L\,.\,\text{access}(\text{off} + j)$

# Wheeler graphs



$O = \texttt{010100010101011010101}$

$L = \texttt{ttcggaaata}$

$\text{num} = 2$

$\longrightarrow \text{off} = O \,.\, \text{rank}_0(O \,.\, \text{select}_1(2 - 1) = 3)$

$\text{for } j \text{ in } \{0, 1, \ldots, \text{num-1}\}:$

$L \,.\, \text{access}(\text{off} + j)$

# Wheeler graphs



$O =$ 010100010101011010101

$L =$ ttcggaaata

num $= 2$

off $= O \,.\, \mathrm{rank}_0(3) = 2$

for $j$ in $\{0, 1, \ldots, \text{num-1}\}$:

$L \,.\, \mathrm{access}(\text{off} + j)$

# Wheeler graphs



$$O = \texttt{010100010101011010101}$$
$$L = \texttt{ttcggaaata}$$

num $= 2$

off $= 2$

$\rightarrow$ for $j$ in $\{0, 1\}$:

$L.\text{access}(\text{off} + j)$

# Wheeler graphs



$O = $ 010100010101011010101

$L = $ ttcggaaata
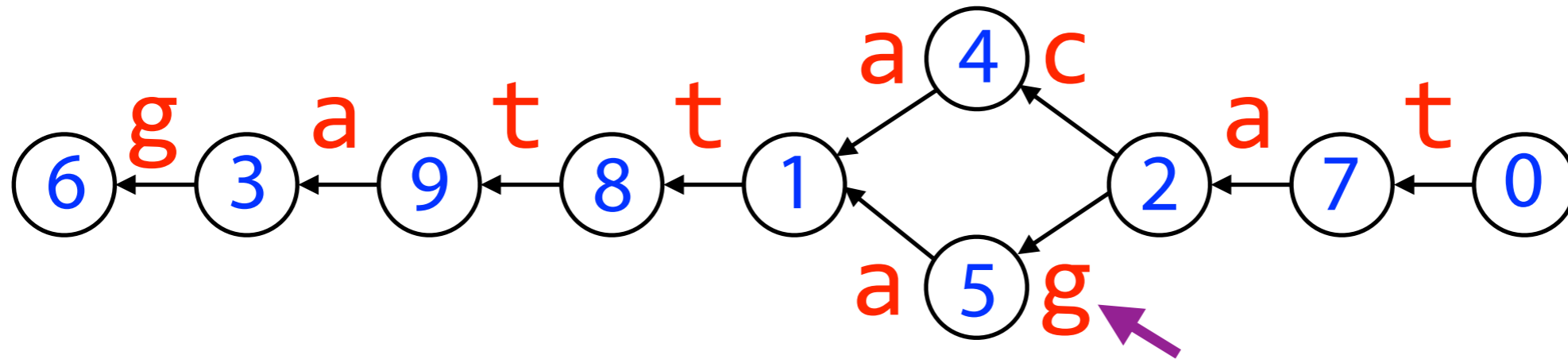
num $ = 2$

off $ = 2$

for $j$ in $\{0, 1\}$:
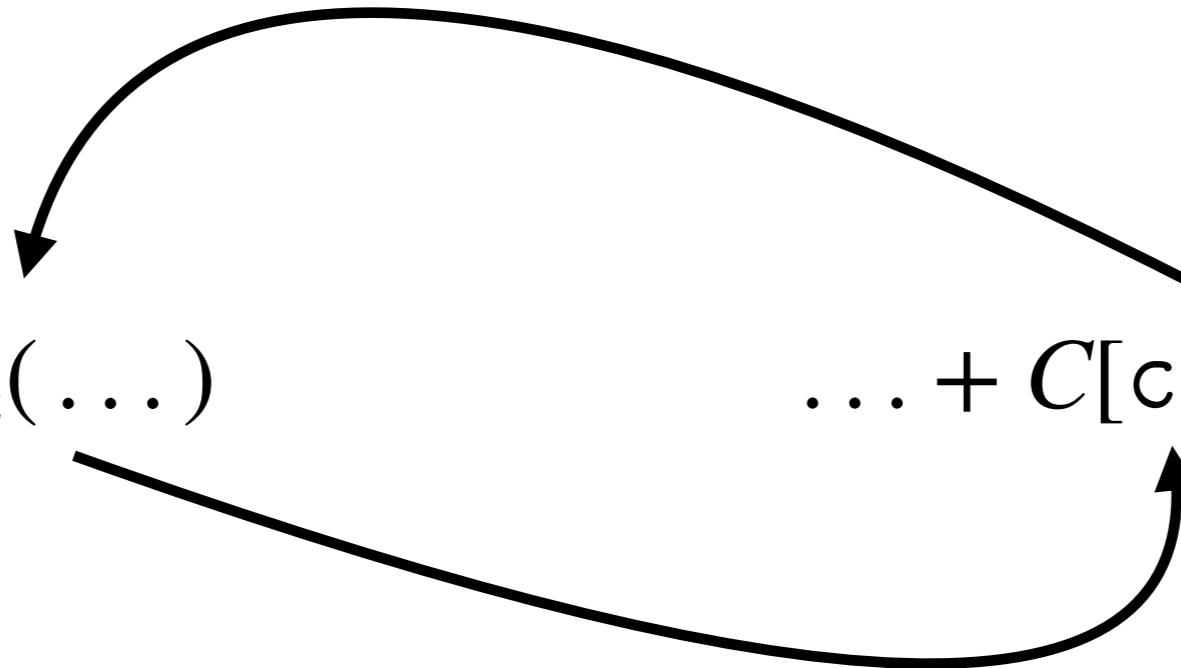
$L . \text{access}(2 + 0) = $ c

# Wheeler graphs



$$O = 0101001010101011010101$$
$$L = \text{ttcggaaata}$$

$$\text{num} = 2$$

$$\text{off} = 2$$

for $j$ in $\{0, 1\}$:

$L$ . access$(2 + 1) = \text{g}$

# Wheeler graphs

**FM Index** match query loop:

$$C[\text{c}], C[\text{c} + 1]$$

$$L.\text{rank}_{\text{c}}(\dots) \qquad \dots + C[\text{c}])$$

$F$                      $L$

| $\$$ | 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 | $\$$ | $a_3$ |

# Wheeler graphs



**Wheeler graph** match query:

$I$ : `10010101010101010101`

$O$ : `01010010101011010101`

$L$ : ttcggaaata

$C[\mathtt{c}], C[\mathtt{c}+1]$

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

$I . \mathsf{rank}_1(I . \mathsf{select}_0(\ldots + C[\mathtt{c}]))$

$L . \mathsf{rank}_\mathtt{c}(\ldots)$