



Short read alignment

Andrey Prjibelski

Center for Algorithmic Biotechnology, SPbU

Alignment

AACGCTAACGGTAA
AACCGCGAACTAA

Alignment

AACGCTAACGGTAA

AACCGCGAACTAA



AAC - GCTAACGGTAA

AACCGCGAAC - - TAA

Short read alignment

Find the read in the genome

Short read alignment

- Challenges?

Short read alignment

- **Challenges**

- Small length
- Gigabytes of data
- Different sequencing errors
- SNPs
- Genomic repeats

- **Tools**

- Bowtie2, BWA MEM, minimap2 (Genomic)
- HiSat2, STAR (RNA-Seq)
- and many more

Bowtie

Mapping Illumina reads

Burrows-Wheeler transform

a c a a c g

Burrows-Wheeler transform

a c a a c g \$

Burrows-Wheeler transform

a c a a c g \$
\$ a c a a c g

Burrows-Wheeler transform

a c a a c g \$

\$ a c a a c g

g \$ a c a a c

Burrows-Wheeler transform

a c a a c g \$
\$ a c a a c g
g \$ a c a a c
c g \$ a c a a
a c g \$ a c a
a a c g \$ a c
c a a c g \$ a

Burrows-Wheeler transform

\$ a c a a c g
a a c g \$ a c
a c a a c g \$
a c g \$ a c a
c a a c g \$ a
c g \$ a c a a
g \$ a c a a c

Burrows-Wheeler transform

\$ a c a a c g
a a c g \$ a c
a c a a c g \$
a c g \$ a c a
c a a c g \$ a
c g \$ a c a a
g \$ a c a a c

Burrows-Wheeler transform

a c a a c g \$



g c \$ a a a c

Burrows-Wheeler transform

a c a a c g \$
↑ ?
g c \$ a a a c

Burrows-Wheeler transform

\$ a c a a c g
a a c g \$ a c
a c a a c g \$
a c g \$ a c a
c a a c g \$ a
c g \$ a c a a
g \$ a c a a c

Burrows-Wheeler transform

\$ a c a a c g
a a c g \$ a c
a c a a c g \$
a c g \$ a c a
c a a c g \$ a
c g \$ a c a a
g \$ a c a a c

Burrows-Wheeler transform

g \$ a c a a c
c a a c g \$ a
\$ a c a a c g
a a c g \$ a c
a c a a c g \$
a c g \$ a c a
c g \$ a c a a

Burrows-Wheeler transform

\$ a c a a c g
a a c g \$ a c
a c a a c g \$
a c g \$ a c a
c a a c g \$ a
c g \$ a c a a
g \$ a c a a c

Burrows-Wheeler transform

g	\$ a c a a c g
c	a a c g \$ a c
\$	a c a a c g \$
a	a c g \$ a c a
a	c a a c g \$ a
a	c g \$ a c a a
c	g \$ a c a a c

Burrows-Wheeler transform

\$	\$ a c a a c g
a	a a c g \$ a c
a	a c a a c g \$
a	a c g \$ a c a
c	c a a c g \$ a
c	c g \$ a c a a
g	g \$ a c a a c

Burrows-Wheeler transform

\$
a
a
a
c
c
g

\$	a	c	a	a	c	g
a	a	c	g	\$	a	c
a	c	a	a	c	g	\$
a	c	g	\$	a	c	a
c	a	a	c	g	\$	a
c	g	\$	a	c	a	a
g	\$	a	c	a	a	c

Burrows-Wheeler transform

g \$

c a

\$ a

a a

a c

a c

c g

\$ a c a a c **g**

a a c g \$ a **c**

a c a a c g **\$**

a c g \$ a c **a**

c a a c g \$ **a**

c g \$ a c a **a**

g \$ a c a a **c**

Burrows-Wheeler transform

\$ a

a a

a c

a c

c a

c g

g \$

\$ a c a a c g

a a c g \$ a c

a c a a c g \$

a c g \$ a c a

c a a c g \$ a

c g \$ a c a a

g \$ a c a a c

Burrows-Wheeler transform

\$ a

a a

a c

a c

c a

c g

g \$

\$ a c a a c g

a a c g \$ a c

a c a a c g \$

a c g \$ a c a

c a a c g \$ a

c g \$ a c a a

g \$ a c a a c

Burrows-Wheeler transform

g \$ a

c a a

\$ a c

a a c

a c a

a c g

c g \$

\$ a c a a c **g**

a a c g \$ a **c**

a c a a c g **\$**

a c g \$ a c **a**

c a a c g \$ **a**

c g \$ a c a **a**

g \$ a c a a **c**

Burrows-Wheeler transform

\$ a c

a a c

a c a

a c g

c a a

c g \$

g \$ a

\$ a c a a c **g**

a a c g **\$** a **c**

a c a a c g **\$**

a c g **\$** a c **a**

c a a c g **\$** **a**

c g **\$** a c a **a**

g **\$** a c a a **c**

Burrows-Wheeler transform

\$ a c

a a c

a c a

a c g

c a a

c g \$

g \$ a

\$ a c a a c g

a a c g \$ a c

a c a a c g \$

a c g \$ a c a

c a a c g \$ a

c g \$ a c a a

g \$ a c a a c

Burrows-Wheeler transform

g \$ a c

c a a c

\$ a c a

a a c g

a c a a

a c g \$

c g \$ a

\$ a c a a c **g**

a a c g \$ a **c**

a c a a c g **\$**

a c g \$ a c **a**

c a a c g \$ **a**

c g \$ a c a **a**

g \$ a c a a **c**

Burrows-Wheeler transform

\$ a c a

a a c g

a c a a

a c g \$

c a a c

c g \$ a

g \$ a c

\$ a c a a c g

a a c g \$ a c

a c a a c g \$

a c g \$ a c a

c a a c g \$ a

c g \$ a c a a

g \$ a c a a c

Burrows-Wheeler transform

g \$ a c a

c a a c g

\$ a c a a

a a c g \$

a c a a c

a c g \$ a

c g \$ a c

\$ a c a a c **g**

a a c g \$ a **c**

a c a a c g **\$**

a c g \$ a c **a**

c a a c g \$ **a**

c g \$ a c a **a**

g \$ a c a a **c**

Burrows-Wheeler transform

\$ a c a a
a a c g \$
a c a a c
a c g \$ a
c a a c g
c g \$ a c
g \$ a c a

\$ a c a a c g
a a c g \$ a c
a c a a c g \$
a c g \$ a c a
c a a c g \$ a
c g \$ a c a a
g \$ a c a a c

Burrows-Wheeler transform

g \$ a c a a

c a a c g \$

\$ a c a a c

a a c g \$ a

a c a a c g

a c g \$ a c

c g \$ a c a

\$ a c a a c **g**

a a c g \$ a **c**

a c a a c g **\$**

a c g \$ a c **a**

c a a c g \$ **a**

c g \$ a c a **a**

g \$ a c a a **c**

Burrows-Wheeler transform

\$ a c a a c
a a c g \$ a
a c a a c g
a c g \$ a c
c a a c g \$
c g \$ a c a
g \$ a c a a

\$ a c a a c g
a a c g \$ a c
a c a a c g \$
a c g \$ a c a
c a a c g \$ a
c g \$ a c a a
g \$ a c a a c

Burrows-Wheeler transform

g \$ a c a a c

c a a c g \$ a

\$ a c a a c g

a a c g \$ a c

a c a a c g \$

a c g \$ a c a

c g \$ a c a a

\$ a c a a c **g**

a a c g \$ a **c**

a c a a c g **\$**

a c g \$ a c **a**

c a a c g \$ **a**

c g \$ a c a **a**

g \$ a c a a **c**

Burrows-Wheeler transform

\$ a c a a c g

a a c g \$ a c

a c a a c g \$

a c g \$ a c a

c a a c g \$ a

c g \$ a c a a

g \$ a c a a c

\$ a c a a c g

a a c g \$ a c

a c a a c g \$

a c g \$ a c a

c a a c g \$ a

c g \$ a c a a

g \$ a c a a c

First-last property

$\$$ ₁	a	c	a	a	c	g
a ₁	a	c	g	$\$$	a	c
a ₂	c	a	a	c	g	$\$$
a ₃	c	g	$\$$	a	c	a
c ₁	a	a	c	g	$\$$	a
c ₂	g	$\$$	a	c	a	a
g ₁	$\$$	a	c	a	a	c

First-last property

$\$$ ₁ a c a a c **g**₁
a₁ a c g \$ a **c**
a₂ c a a c g **\$**₁
a₃ c g \$ a c **a**
c₁ a a c g \$ **a**
c₂ g \$ a c a **a**
g₁ \$ a c a a **c**

First-last property

$\$$ ₁ a c a a c **g**₁

a₁ a c g \$ a **c**

a₂ c a a c g **\$**₁

a₃ c g \$ a c **a**

c₁ a a c g \$ **a**

c₂ g \$ a c a **a**

g₁ \$ a c a a **c**

First-last property

a_3 c g \$ a c **a**
 c_1 a a c g \$ **a**
 c_2 g \$ a c a **a**

$\$1$ a c a a c **g_1**
 a_1 a c g \$ a **c**
 a_2 c a a c g **$\$1$**
 a_3 c g \$ a c **a**
 c_1 a a c g \$ **a**
 c_2 g \$ a c a **a**
 g_1 \$ a c a a **c**

First-last property

a a₃ c g \$ a c
a c₁ a a c g \$
a c₂ g \$ a c a

\$₁ a c a a c **g**₁
a₁ a c g \$ a **c**
a₂ c a a c g \$₁
a₃ c g \$ a c **a**
c₁ a a c g \$ **a**
c₂ g \$ a c a **a**
g₁ \$ a c a a **c**

First-last property

a a₃ c g \$ a c
a c₁ a a c g \$
a c₂ g \$ a c a

\$₁ a c a a c **g**₁
a₁ a c g \$ a **c**
a₂ c a a c g \$₁
a₃ c g \$ a c **a**
c₁ a a c g \$ **a**
c₂ g \$ a c a **a**
g₁ \$ a c a a **c**

First-last property

a₁ **a₃** c g \$ a c
a₂ **c₁** a a c g \$
a₃ **c₂** g \$ a c a

\$₁ a c a a c **g₁**
a₁ a c g \$ a **c**
a₂ c a a c g **\$₁**
a₃ c g \$ a c **a**
c₁ a a c g \$ **a**
c₂ g \$ a c a **a**
g₁ \$ a c a a **c**

First-last property

a_3 c g \$ a c a_1
 c_1 a a c g \$ a_2
 c_2 g \$ a c a a_3

$\$1$ a c a a c g_1
 a_1 a c g \$ a c
 a_2 c a a c g $\$1$
 a_3 c g \$ a c a
 c_1 a a c g \$ a
 c_2 g \$ a c a a
 g_1 \$ a c a a c

First-last property

a_3 c g \$ a c a_1
 c_1 a a c g \$ a_2
 c_2 g \$ a c a a_3

$\$1$ a c a a c g_1
 a_1 a c g \$ a c
 a_2 c a a c g $\$1$
 a_3 c g \$ a c a
 c_1 a a c g \$ a
 c_2 g \$ a c a a
 g_1 \$ a c a a c

First-last property

a_3 c g \$ a c a_1
 c_1 a a c g \$ a_2
 c_2 g \$ a c a a_3

$\$1$ a c a a c g_1
 a_1 a c g \$ a c
 a_2 c a a c g $\$1$
 a_3 c g \$ a c a_1
 c_1 a a c g \$ a_2
 c_2 g \$ a c a a_3
 g_1 \$ a c a a c

First-last property

$\$$ ₁	a	c	a	a	c	g	g ₁		
a	a ₁	c	g	$\$$	a	c	c ₁		
a	a ₂	c	a	a	c	g	$\$$ ₁		
a	a ₃	c	g	$\$$	a	c	a	a ₁	
c	c ₁	a	a	c	g	$\$$	a	a ₂	
c	c ₂	g	$\$$	a	c	a	a	a ₃	
g	g ₁	$\$$	a	c	a	a	a	c	c ₂

Suffix array

6. \$
2. a a c g \$
0. a c a a c g \$
3. a c g \$
1. c a a c g \$
4. c g \$
5. g \$

Suffix array and BWT

6. \$
2. a a c g \$
0. a c a a c g \$
3. a c g \$
1. c a a c g \$
4. c g \$
5. g \$

\$ a c a a c g
a a c g \$ a c
a c a a c g \$
a c g \$ a c a
c a a c g \$ a
c g \$ a c a a
g \$ a c a a c

Suffix array and BWT

$B[i] = \$$ if $S[i] = 0$

$B[i] = X[S[i] - 1]$ otherwise

Suffix array interval

6. \$

2. a a c g \$

0. a c a a c g \$

3. a c g \$

1. c a a c g \$

4. c g \$

5. g \$

$R("a") = (0, 6)$

Suffix array interval

6. \$

$R(\text{"a"}) = (1, 3)$

2. a a c g \$

0. a c a a c g \$

3. a c g \$

1. c a a c g \$

4. c g \$

5. g \$

Suffix array interval

6. \$ $R(\text{"a c"}) = (2, 3)$

2. a a c g \$

0. a c a a c g \$

3. a c g \$

1. c a a c g \$

4. c g \$

5. g \$

Suffix array interval

$$R_L(W) = \min \{k: W \text{ is prefix of } X_{S[k]} \}$$

$$R_H(W) = \max \{k: W \text{ is prefix of } X_{S[k]} \}$$

Suffix array interval

$$C(x) = | \{ 0 \leq j \leq n-2 : X[j] < x \} |$$

a c a a c g \$

$$C(a) = 0, C(c) = 3, C(g) = 5, \dots$$

Suffix array interval

→ 6. \$
→ 2. a a c g \$
0. a c a a c g \$
3. a c g \$
→ 1. c a a c g \$
4. c g \$
→ 5. g \$
→

Suffix array interval

$$O(x, i) = | \{ 0 \leq j \leq i : B[j] = x \} |$$

g c \$ a a a c

$$O(a, 0) = 0, O(a, 1) = 0, O(a, 2) = 0, \\ O(a, 3) = 1, O(a, 4) = 2, \dots$$

Suffix array interval

							a	c	g	t
$\$$ ₁	a	c	a	a	c	g ₁	0	0	1	0
a ₁	a	c	g	$\$$	a	c ₁	0	1	1	0
a ₂	c	a	a	c	g	$\$$ ₁	0	1	1	0
a ₃	c	g	$\$$	a	c	a ₁	1	1	1	0
c ₁	a	a	c	g	$\$$	a ₂	2	1	1	0
c ₂	g	$\$$	a	c	a	a ₃	3	1	1	0
g ₁	$\$$	a	c	a	a	c ₂	3	2	1	0

Suffix array interval

$$R_L(xW) = C(x) + O(x, R_L(W) - 1)$$

$$R_H(xW) = C(x) + O(x, R_H(W)) - 1$$

$$R_L("") = 0$$

$$R_H("") = \text{len}(X) - 1$$

Pattern search

								a	c	g	t	
$C(\$) = 0$	→											
$C(a) = 1$	→	\$ ₁	a	c	a	a	c	g ₁	0	0	1	0
		a ₁	a	c	g	\$	a	c ₁	0	1	1	0
		a ₂	c	a	a	c	g	\$ ₁	0	1	1	0
$C(c) = 4$	→	a ₃	c	g	\$	a	c	a ₁	1	1	1	0
		c ₁	a	a	c	g	\$	a ₂	2	1	1	0
$C(g) = 6$	→	c ₂	g	\$	a	c	a	a ₃	3	1	1	0
$C(t) = 7$	→	g ₁	\$	a	c	a	a	c ₂	3	2	1	0

Pattern search

a a c

							a	c	g	t		
$C(\$) = 0$	→	\$ ₁	a	c	a	a	c	g ₁	0	0	1	0
$C(a) = 1$	→	a ₁	a	c	g	\$	a	c ₁	0	1	1	0
		a ₂	c	a	a	c	g	\$ ₁	0	1	1	0
$C(c) = 4$	→	a ₃	c	g	\$	a	c	a ₁	1	1	1	0
		c ₁	a	a	c	g	\$	a ₂	2	1	1	0
$C(g) = 6$	→	c ₂	g	\$	a	c	a	a ₃	3	1	1	0
$C(t) = 7$	→	g ₁	\$	a	c	a	a	c ₂	3	2	1	0

Pattern search

a a c

							a	c	g	t		
$C(\$) = 0$	→	\$ ₁	a	c	a	a	c	g ₁	0	0	1	0
$C(a) = 1$	→	a ₁	a	c	g	\$	a	c ₁	0	1	1	0
		a ₂	c	a	a	c	g	\$ ₁	0	1	1	0
		a ₃	c	g	\$	a	c	a ₁	1	1	1	0
$C(c) = 4$	→	c ₁	a	a	c	g	\$	a ₂	2	1	1	0
		c ₂	g	\$	a	c	a	a ₃	3	1	1	0
$C(g) = 6$	→	g ₁	\$	a	c	a	a	c ₂	3	2	1	0
$C(t) = 7$	→											

$$R_L(“”) = 0$$

$$R_H(“”) = 6$$

Pattern search

a a c

							a	c	g	t		
$C(\$) = 0$	→	\$ ₁	a	c	a	a	c	g ₁	0	0	1	0
$C(a) = 1$	→	a ₁	a	c	g	\$	a	c ₁	0	1	1	0
		a ₂	c	a	a	c	g	\$ ₁	0	1	1	0
$C(c) = 4$	→	a ₃	c	g	\$	a	c	a ₁	1	1	1	0
		c ₁	a	a	c	g	\$	a ₂	2	1	1	0
$C(g) = 6$	→	c ₂	g	\$	a	c	a	a ₃	3	1	1	0
$C(t) = 7$	→	g ₁	\$	a	c	a	a	c ₂	3	2	1	0

$$R_L(xW) = C(x) + O(x, R_L(W)) - 1$$

$$R_H(xW) = C(x) + O(x, R_H(W)) - 1$$

Pattern search

a a c

$$R_L(c) = C(c) + O(c, R_L("")) - 1$$

$$R_H(c) = C(c) + O(c, R_H("")) - 1$$

							a	c	g	t		
$C(\$) = 0$	→	\$ ₁	a	c	a	a	c	g ₁	0	0	1	0
$C(a) = 1$	→	a ₁	a	c	g	\$	a	c ₁	0	1	1	0
		a ₂	c	a	a	c	g	\$ ₁	0	1	1	0
$C(c) = 4$	→	a ₃	c	g	\$	a	c	a ₁	1	1	1	0
		c ₁	a	a	c	g	\$	a ₂	2	1	1	0
$C(g) = 6$	→	c ₂	g	\$	a	c	a	a ₃	3	1	1	0
$C(t) = 7$	→	g ₁	\$	a	c	a	a	c ₂	3	2	1	0

Pattern search

a a c

							a	c	g	t		
$C(\$) = 0$	→	\$ ₁	a	c	a	a	c	g ₁	0	0	1	0
$C(a) = 1$	→	a ₁	a	c	g	\$	a	c ₁	0	1	1	0
		a ₂	c	a	a	c	g	\$ ₁	0	1	1	0
		a ₃	c	g	\$	a	c	a ₁	1	1	1	0
$C(c) = 4$	→	c ₁	a	a	c	g	\$	a ₂	2	1	1	0
		c ₂	g	\$	a	c	a	a ₃	3	1	1	0
$C(g) = 6$	→	g ₁	\$	a	c	a	a	c ₂	3	2	1	0
$C(t) = 7$	→											

$$R_L(ac) = C(a) + O(a, R_L(c)) - 1$$

$$R_H(ac) = C(a) + O(a, R_H(c)) - 1$$

Pattern search

a a c

							a	c	g	t		
$C(\$) = 0$	→	\$ ₁	a	c	a	a	c	g ₁	0	0	1	0
$C(a) = 1$	→	a ₁	a	c	g	\$	a	c ₁	0	1	1	0
		a ₂	c	a	a	c	g	\$ ₁	0	1	1	0
		a ₃	c	g	\$	a	c	a ₁	1	1	1	0
$C(c) = 4$	→	c ₁	a	a	c	g	\$	a ₂	2	1	1	0
		c ₂	g	\$	a	c	a	a ₃	3	1	1	0
$C(g) = 6$	→	g ₁	\$	a	c	a	a	c ₂	3	2	1	0
$C(t) = 7$	→											

$$R_L(ac) = C(a) + O(a, 3)$$

$$R_H(ac) = C(a) + O(a, 5) - 1$$

Pattern search

a a c

							a	c	g	t		
$C(\$) = 0$	→	\$ ₁	a	c	a	a	c	g ₁	0	0	1	0
$C(a) = 1$	→	a ₁	a	c	g	\$	a	c ₁	0	1	1	0
		a ₂	c	a	a	c	g	\$ ₁	0	1	1	0
		a ₃	c	g	\$	a	c	a ₁	1	1	1	0
$C(c) = 4$	→	c ₁	a	a	c	g	\$	a ₂	2	1	1	0
		c ₂	g	\$	a	c	a	a ₃	3	1	1	0
$C(g) = 6$	→	g ₁	\$	a	c	a	a	c ₂	3	2	1	0
$C(t) = 7$	→											

$$R_L(ac) = C(a) + O(a, 3)$$

$$R_H(ac) = C(a) + O(a, 5) - 1$$

Pattern search

a a c

$$R_L(ac) = 1 + 1 = 2$$

$$R_H(ac) = 1 + 3 - 1 = 3$$

							a	c	g	t		
$C(\$) = 0$	→											
$C(a) = 1$	→	\$₁	a	c	a	a	c	g₁	0	0	1	0
		a₁	a	c	g	\$	a	c₁	0	1	1	0
		a₂	c	a	a	c	g	\$₁	0	1	1	0
$C(c) = 4$	→	a₃	c	g	\$	a	c	a₁	1	1	1	0
		c₁	a	a	c	g	\$	a₂	2	1	1	0
$C(g) = 6$	→	c₂	g	\$	a	c	a	a₃	3	1	1	0
$C(t) = 7$	→	g₁	\$	a	c	a	a	c₂	3	2	1	0

Pattern search

a a c

							a	c	g	t		
$C(\$) = 0$	→	\$ ₁	a	c	a	a	c	g ₁	0	0	1	0
$C(a) = 1$	→	a ₁	a	c	g	\$	a	c ₁	0	1	1	0
		a ₂	c	a	a	c	g	\$ ₁	0	1	1	0
		a ₃	c	g	\$	a	c	a ₁	1	1	1	0
$C(c) = 4$	→	c ₁	a	a	c	g	\$	a ₂	2	1	1	0
		c ₂	g	\$	a	c	a	a ₃	3	1	1	0
$C(g) = 6$	→	g ₁	\$	a	c	a	a	c ₂	3	2	1	0
$C(t) = 7$	→											

$$R_L(aac) = C(a) + O(a, R_L(ac)) - 1$$

$$R_H(aac) = C(a) + O(a, R_H(ac)) - 1$$

Pattern search

a a c

							a	c	g	t		
$C(\$) = 0$	→	\$ ₁	a	c	a	a	c	g ₁	0	0	1	0
$C(a) = 1$	→	a ₁	a	c	g	\$	a	c ₁	0	1	1	0
		a ₂	c	a	a	c	g	\$ ₁	0	1	1	0
		a ₃	c	g	\$	a	c	a ₁	1	1	1	0
$C(c) = 4$	→	c ₁	a	a	c	g	\$	a ₂	2	1	1	0
		c ₂	g	\$	a	c	a	a ₃	3	1	1	0
$C(g) = 6$	→	g ₁	\$	a	c	a	a	c ₂	3	2	1	0
$C(t) = 7$	→											

$$R_L(aac) = C(a) + O(a, 2) - 1$$

$$R_H(aac) = C(a) + O(a, 3) - 1$$

Pattern search

a a c

$$R_L(\text{aac}) = 1 + 0 = 1$$

$$R_H(\text{aac}) = 1 + 1 - 1 = 1$$

								a	c	g	t		
$C(\$) = 0$	→												
$C(a) = 1$	→	\$ ₁	a	c	a	a	c	g ₁	0	0	1	0	
			a ₁	a	c	g	\$	a	c ₁	0	1	1	0
			a ₂	c	a	a	c	g	\$ ₁	0	1	1	0
$C(c) = 4$	→		a ₃	c	g	\$	a	c	a ₁	1	1	1	0
			c ₁	a	a	c	g	\$	a ₂	2	1	1	0
$C(g) = 6$	→		c ₂	g	\$	a	c	a	a ₃	3	1	1	0
$C(t) = 7$	→		g ₁	\$	a	c	a	a	c ₂	3	2	1	0

Pattern search

For how long does it work?

Pattern search

For how long does it work?

$O(m)$

Pattern search with errors

$\$$ ₁ a c a a c **g**₁

a₁ a c g \$ a **c**₁

a₂ c a a c g **\$**₁

a₃ c g \$ a c **a**₁

c₁ a a c g \$ **a**₂

c₂ g \$ a c a **a**₃

g₁ \$ a c a a **c**₂

g c a

Pattern search with errors

$\$_1$	a	c	a	a	c	g_1		g	c	a
a_1	a	c	g	$\$$	a	c_1				
a_2	c	a	a	c	g	$\$_1$				
a_3	c	g	$\$$	a	c	a_1				
c_1	a	a	c	g	$\$$	a_2				
c_2	g	$\$$	a	c	a	a_3				
g_1	$\$$	a	c	a	a	c_2				

Pattern search with errors

$\$$ ₁	a	c	a	a	c	g	g ₁	g	c	a
a ₁	a	c	g	\$	a	c	c ₁			
a ₂	c	a	a	c	g	\$	$\$$ ₁			
a ₃	c	g	\$	a	c	a	a ₁			
c ₁	a	a	c	g	\$	a	a ₂			
c ₂	g	\$	a	c	a	a	a ₃			
g ₁	\$	a	c	a	a	c	c ₂			

Pattern search with errors

$\$$ ₁	a	c	a	a	c	g	g ₁		
a ₁	a	c	g	$\$$	a	c ₁			
a ₂	c	a	a	c	g	$\$$ ₁			
a ₃	c	g	$\$$	a	c	a ₁			
c ₁	a	a	c	g	$\$$	a ₂			
c ₂	g	$\$$	a	c	a	a ₃			
g ₁	$\$$	a	c	a	a	c ₂			

g c a

Pattern search with errors

$\$$ ₁	a	c	a	a	c	g	g ₁		
a ₁	a	c	g	$\$$	a	c ₁			
a ₂	c	a	a	c	g	$\$$ ₁			
a ₃	c	g	$\$$	a	c	a ₁			
c ₁	a	a	c	g	$\$$	a ₂			
c ₂	g	$\$$	a	c	a	a ₃			
g ₁	$\$$	a	c	a	a	c ₂			

g c a

Pattern search with errors

$\$$ ₁	a	c	a	a	c	g	g ₁		
a ₁	a	c	g	$\$$	a	c ₁			
a ₂	c	a	a	c	g	$\$$ ₁			
a ₃	c	g	$\$$	a	c	a ₁			
c ₁	a	a	c	g	$\$$	a ₂			
c ₂	g	$\$$	a	c	a	a ₃			
g ₁	$\$$	a	c	a	a	c ₂			

g c a

Pattern search with errors

$\$$ ₁	a	c	a	a	c	g	g ₁
a ₁	a	c	g	$\$$	a	c ₁	
a ₂	c	a	a	c	g	$\$$ ₁	
a ₃	c	g	$\$$	a	c	a ₁	
c ₁	a	a	c	g	$\$$	a ₂	
c ₂	g	$\$$	a	c	a	a ₃	
g ₁	$\$$	a	c	a	a	c ₂	

g c a

Pattern search with errors

$\$$ ₁	a	c	a	a	c	g	g ₁		
a ₁	a	c	g	$\$$	a	c ₁			
a ₂	c	a	a	c	g	$\$$ ₁			
a ₃	c	g	$\$$	a	c	a ₁			
c ₁	a	a	c	g	$\$$	a ₂			
c ₂	g	$\$$	a	c	a	a ₃			
g ₁	$\$$	a	c	a	a	c ₂			

g c a

Pattern search with errors

$\$$ ₁	a	c	a	a	c	g	g ₁		
a	a	c	g	\$	a	c	c ₁		g c a
a	c	a	a	c	g	\$	$\$$ ₁		
a	c	g	\$	a	c	a	a ₁		
c ₁	a	a	c	g	\$	a	a ₂		
c	g	\$	a	c	a	a	a ₃		
g	\$	a	c	a	a	c	c ₂		

Pattern search with errors

$\$$ ₁ a c a a c **g**₁

a₁ a c g \$ a **c**₁

a₂ c a a c g **\$**₁

a₃ c g \$ a c **a**₁

c₁ a a c g \$ **a**₂

c₂ g \$ a c a **a**₃

g₁ \$ a c a a **c**₂

g c a

Pattern search with errors

$\$$ ₁ a c a a c **g**₁

a₁ a c g \$ a **c**₁

a₂ c a a c g **\$**₁

a₃ c g \$ a c **a**₁

c₁ a a c g \$ **a**₂

c₂ g \$ a c a **a**₃

g₁ \$ a c a a **c**₂

g c a

Pattern search with errors

$\$$ ₁	a	c	a	a	c	g	g ₁
a ₁	a	c	g	$\$$	a	c ₁	
a ₂	c	a	a	c	g	$\$$ ₁	
a ₃	c	g	$\$$	a	c	a ₁	
c ₁	a	a	c	g	$\$$	a ₂	
c ₂	g	$\$$	a	c	a	a ₃	
g ₁	$\$$	a	c	a	a	c ₂	

g c a



Pattern search with errors

$\$$ ₁	a	c	a	a	c	g	g ₁		g	c	a
a ₁	a	c	g	$\$$	a	c ₁					
a ₂	c	a	a	c	g	$\$$ ₁					
a ₃	c	g	$\$$	a	c	a ₁					
c ₁	a	a	c	g	$\$$	a ₂					
c ₂	g	$\$$	a	c	a	a ₃					
g ₁	$\$$	a	c	a	a	c ₂					

Pattern search with errors

$\$$ ₁	a	c	a	a	c	g	g ₁			c	c	a
a ₁	a	c	g	$\$$	a	c ₁						
a ₂	c	a	a	c	g	$\$$ ₁						
a ₃	c	g	$\$$	a	c	a ₁						
c ₁	a	a	c	g	$\$$	a ₂						
c ₂	g	$\$$	a	c	a	a ₃						
g ₁	$\$$	a	c	a	a	c ₂						

Pattern search with errors

$\$$ ₁	a	c	a	a	c	g	g ₁
a ₁	a	c	g	$\$$	a	c ₁	
a ₂	c	a	a	c	g	$\$$ ₁	
a ₃	c	g	$\$$	a	c	a ₁	
c ₁	a	a	c	g	$\$$	<u>a₂</u>	
c ₂	g	$\$$	a	c	a	a ₃	
g ₁	$\$$	a	c	a	a	c ₂	

c c a



Pattern search with errors

$\$$ ₁	a	c	a	a	c	g	g ₁	a	c	a
a ₁	a	c	g	$\$$	a	c ₁				
a ₂	c	a	a	c	g	$\$$ ₁				
a ₃	c	g	$\$$	a	c	a ₁				
c ₁	a	a	c	g	$\$$	a ₂				
c ₂	g	$\$$	a	c	a	a ₃				
g ₁	$\$$	a	c	a	a	c ₂				

Pattern search with errors

$\$$ ₁	a	c	a	a	c	g	g ₁	a	c	a
a ₁	a	c	g	$\$$	a	c ₁				
a ₂	c	a	a	c	g	$\$$ ₁				
a ₃	c	g	$\$$	a	c	a ₁				
c ₁	a	a	c	g	$\$$	a ₂				
c ₂	g	$\$$	a	c	a	a ₃				
g ₁	$\$$	a	c	a	a	c ₂				

Pattern search with errors

$\$$ ₁	a	c	a	a	c	g	g ₁	a	c	a
a ₁	a	c	g	$\$$	a	c ₁				
a ₂	c	a	a	c	g	$\$$ ₁				
a ₃	c	g	$\$$	a	c	a ₁				
c ₁	a	a	c	g	$\$$	a ₂				
c ₂	g	$\$$	a	c	a	a ₃				
g ₁	$\$$	a	c	a	a	c ₂				

Pattern search with errors

$\$$ ₁	a	c	a	a	c	g	g ₁
a ₁	a	c	g	$\$$	a	c ₁	
a ₂	c	a	a	c	g	$\$$ ₁	
a ₃	c	g	$\$$	a	c	a ₁	
c ₁	a	a	c	g	$\$$	a ₂	
c ₂	g	$\$$	a	c	a	a ₃	
g ₁	$\$$	a	c	a	a	c ₂	

a c a



Bowtie

- Backtracking

Bowtie

- Backtracking
 - Quality-aware
 - Limit total number of backtracks
 - Limit quality distance

Bowtie

- Backtracking
 - Quality-aware
 - Limit total number of backtracks
 - Limit quality distance
- Seeds
 - Selected at high-quality end
 - Used to prevent excessive backtracking

Burrows-Wheeler mirror transform

a c a a c g \$ a a c

Burrows-Wheeler mirror transform

\$ g c a a c a a a c

Burrows-Wheeler mirror transform

\$ g c a a c a c a a

Burrows-Wheeler mirror transform

\$ g c a a c a c a a

The diagram illustrates the Burrows-Wheeler mirror transform. It shows two strings of characters. The first string is "\$ g c a a c a", where the characters 'c', 'a', and 'a' are highlighted in blue. The second string is "c a a", where the first character 'c' is bolded. A red arrow points from the top of the 'c' in the second string to the top of the 'c' in the first string, indicating a mapping or transformation between the two strings.

Burrows-Wheeler mirror transform

a c a a c g \$

 a a c

Bowtie

- Backtracking
 - Quality-aware
 - Limit total number of backtracks
 - Limit quality distance
- Seeds
 - Selected at high-quality end
 - Used to prevent excessive backtracking

Seeds

accg...

..cgaa

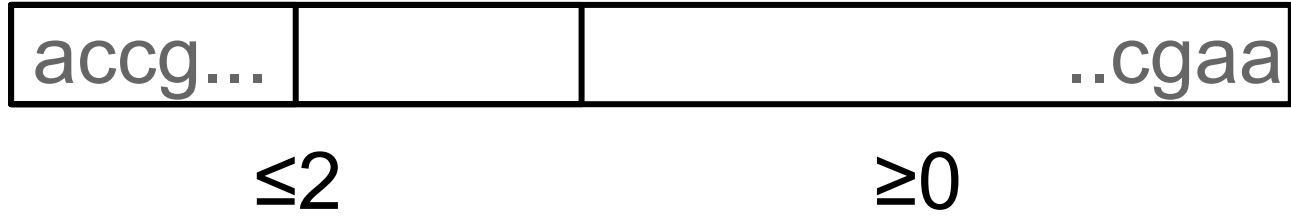
Seeds



≤ 2

≥ 0

Seeds



Seeds

1.

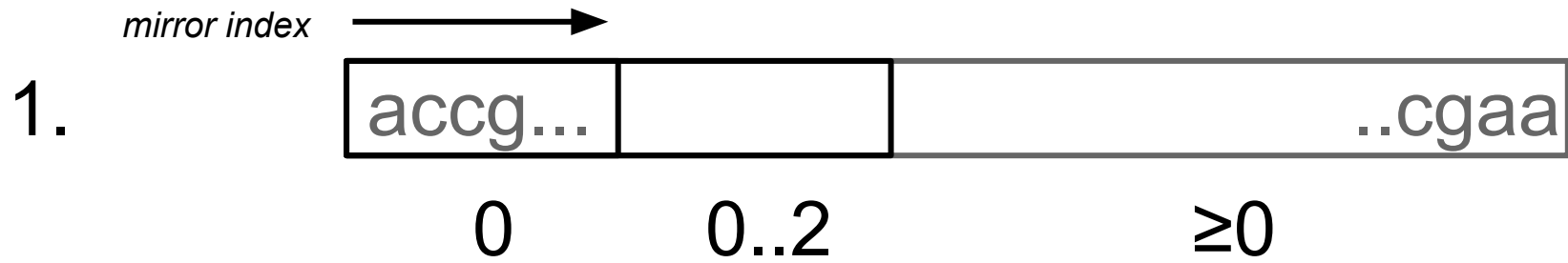


0

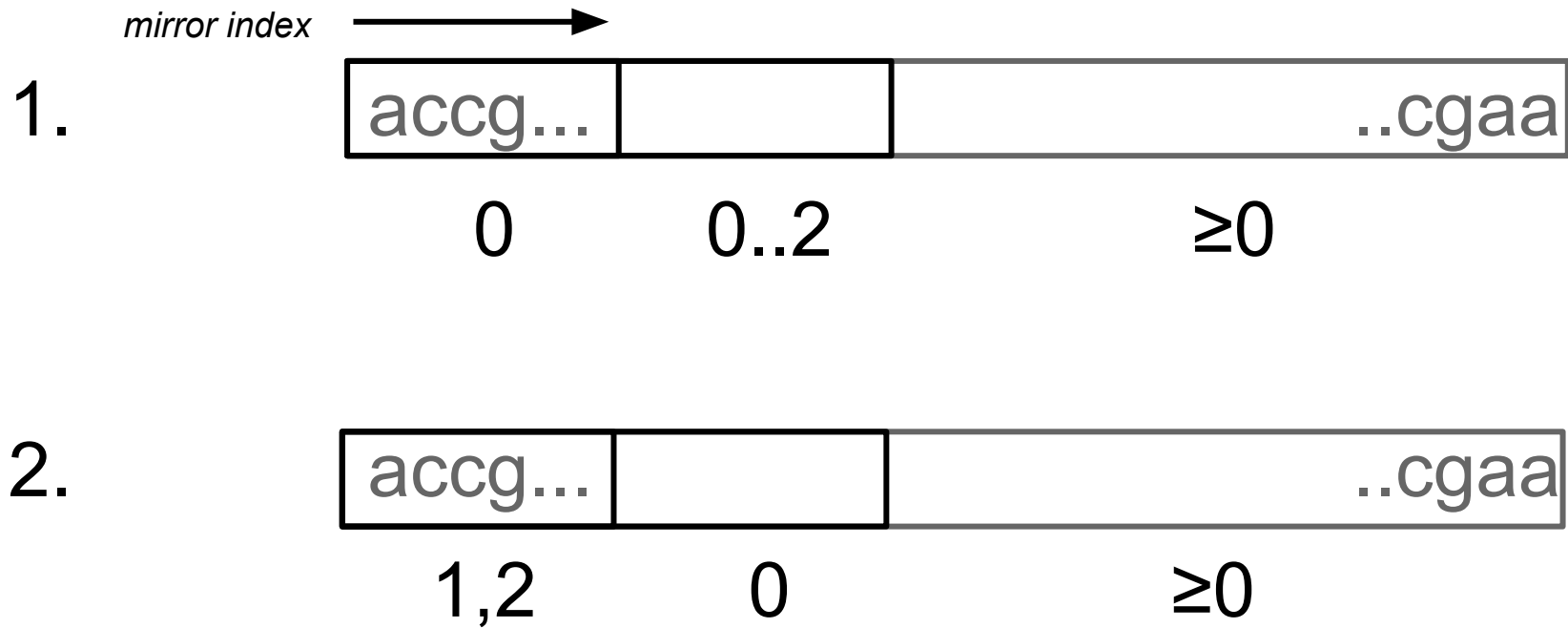
0..2

≥ 0

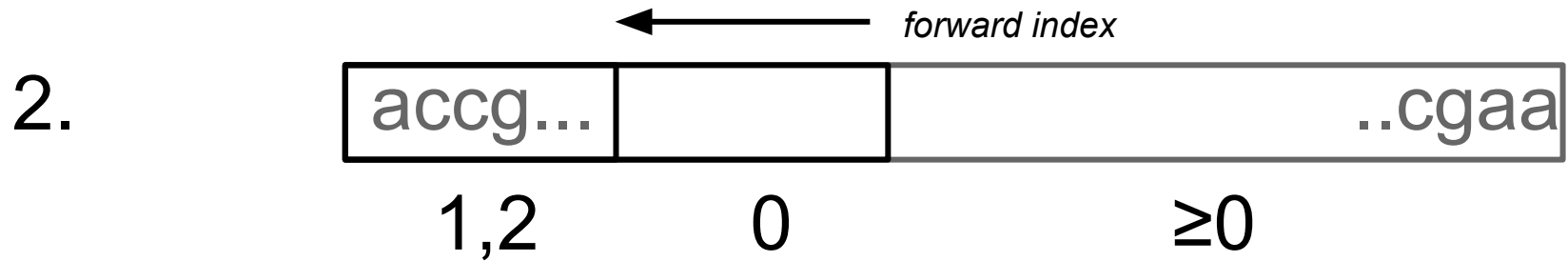
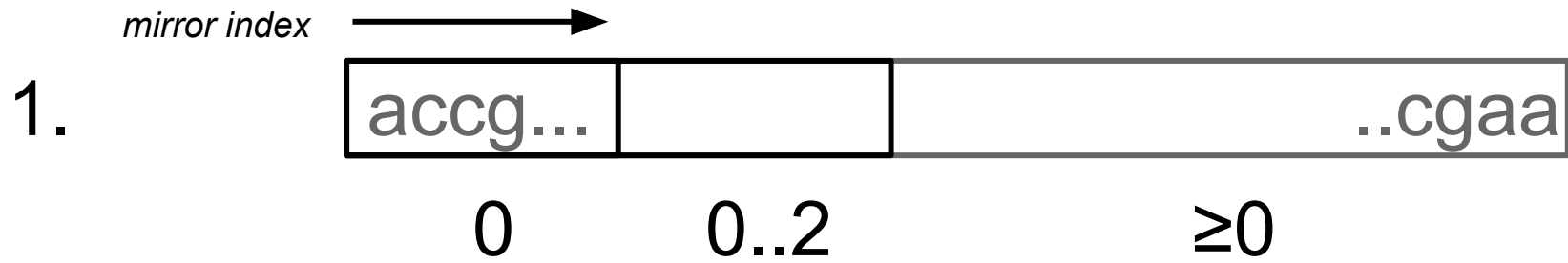
Seeds



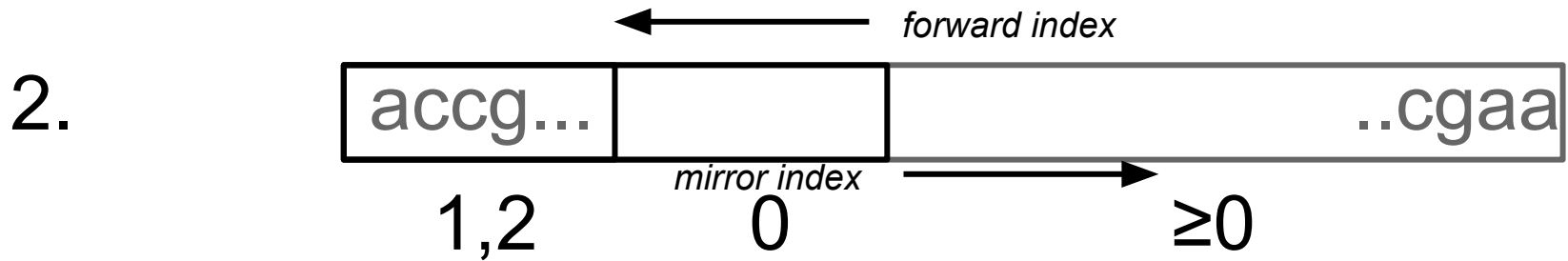
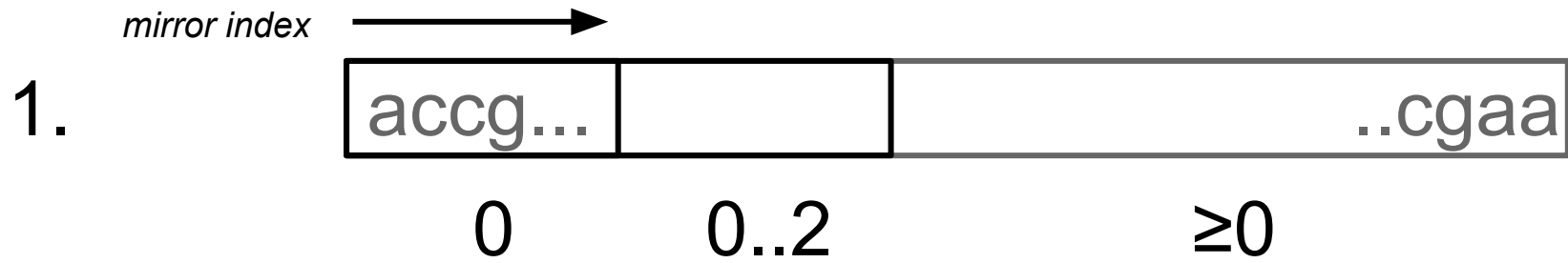
Seeds



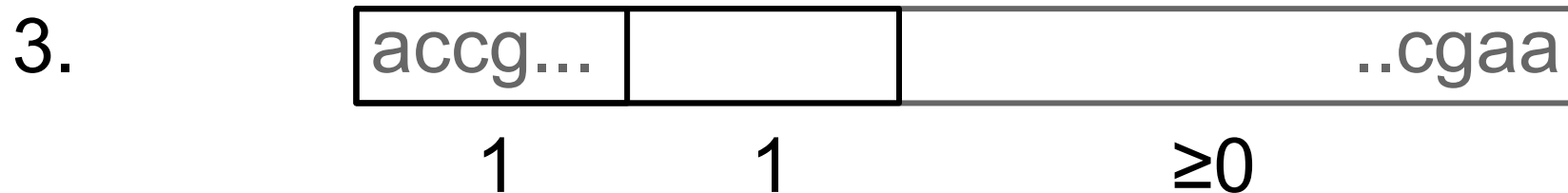
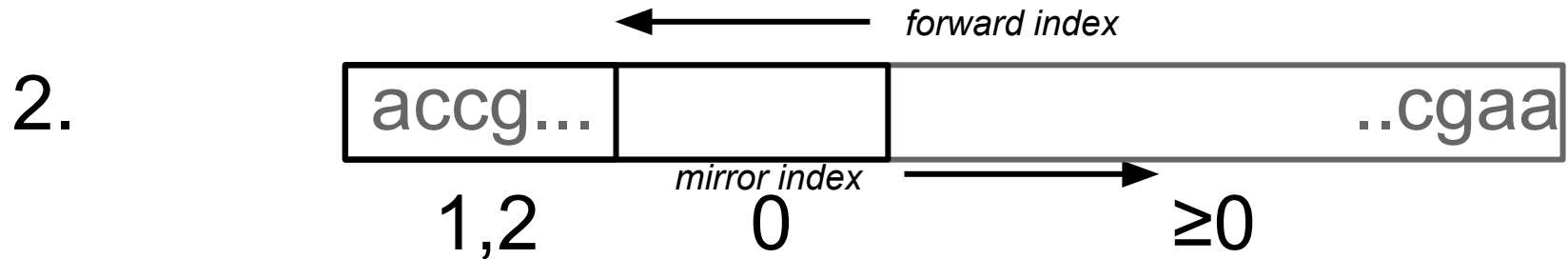
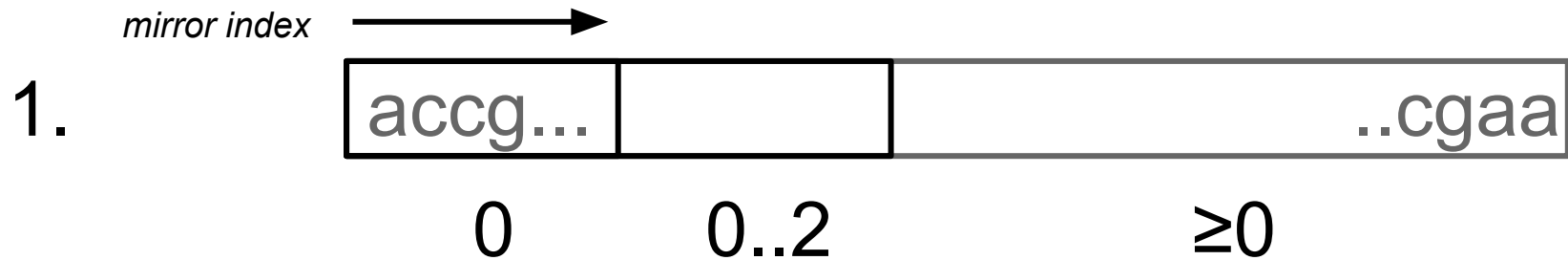
Seeds



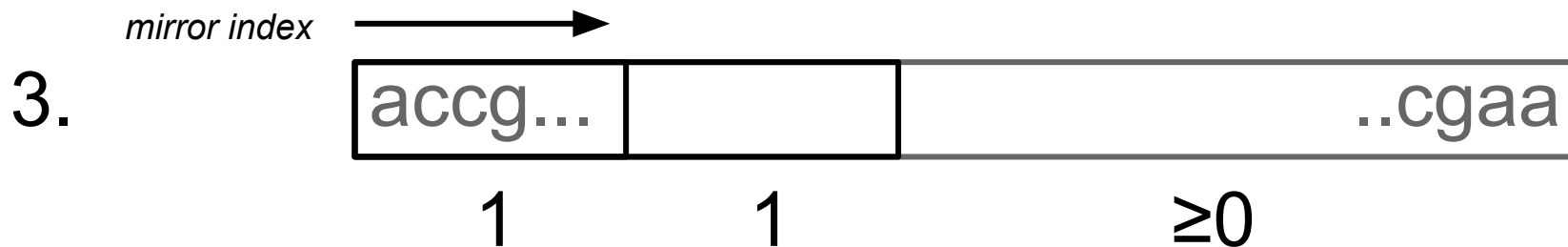
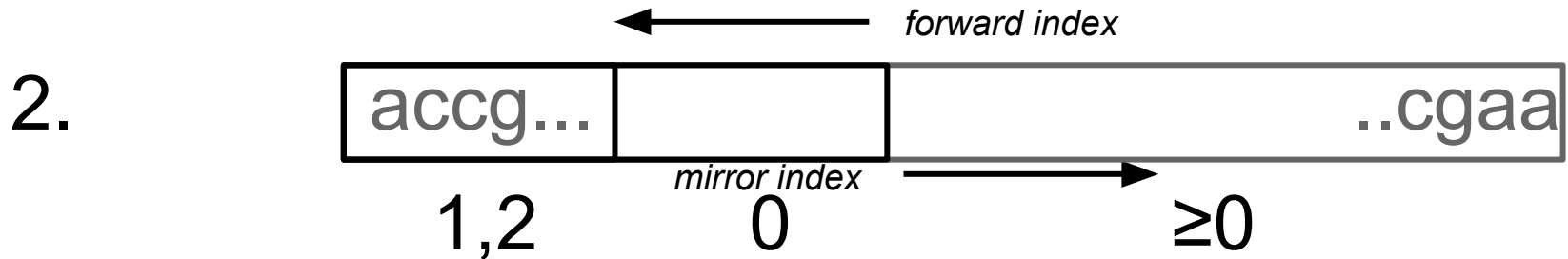
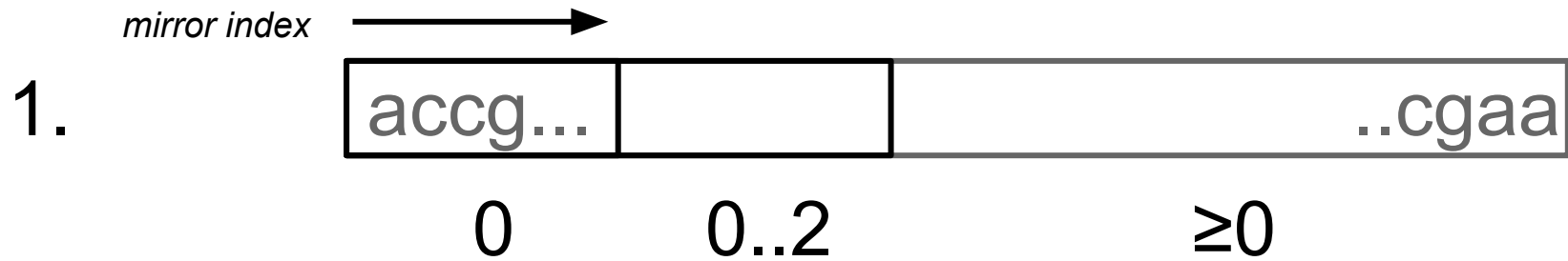
Seeds



Seeds



Seeds



BWA

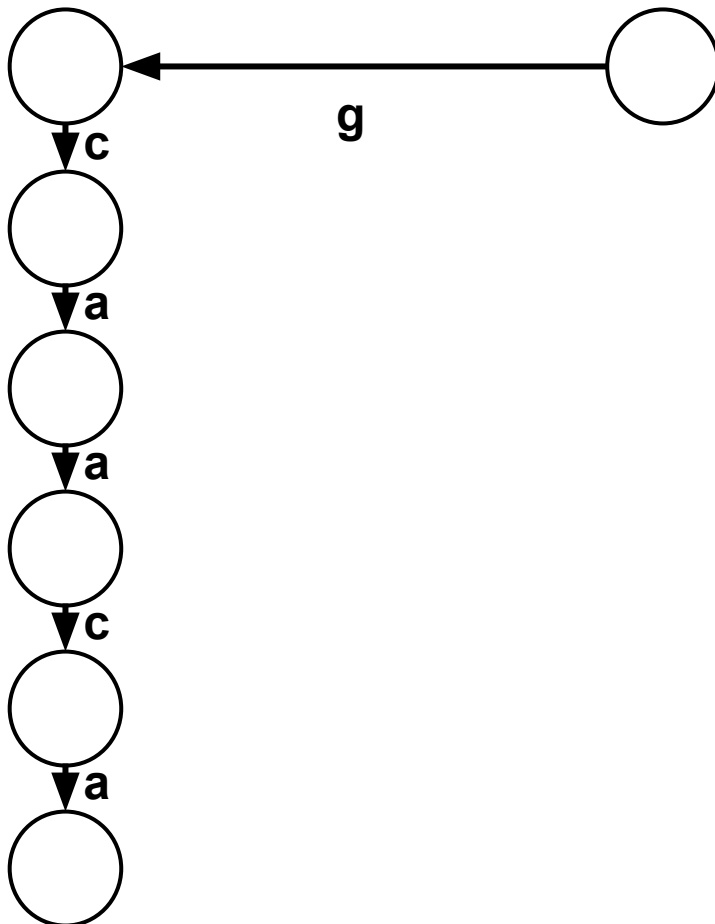
Mapping IonTorrent/454 reads

Prefix trie

a c a a c g

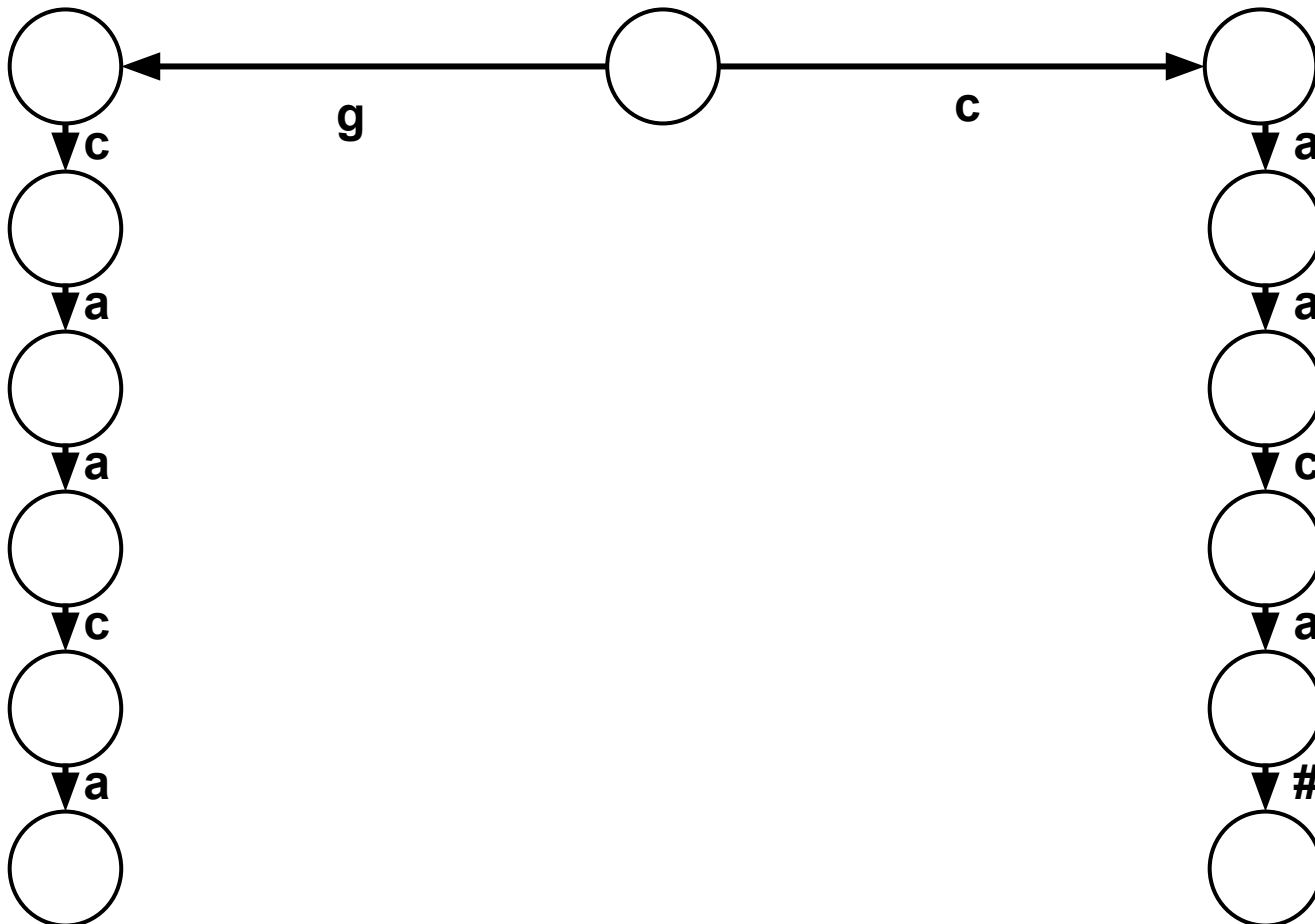
Prefix trie

a c a a c g



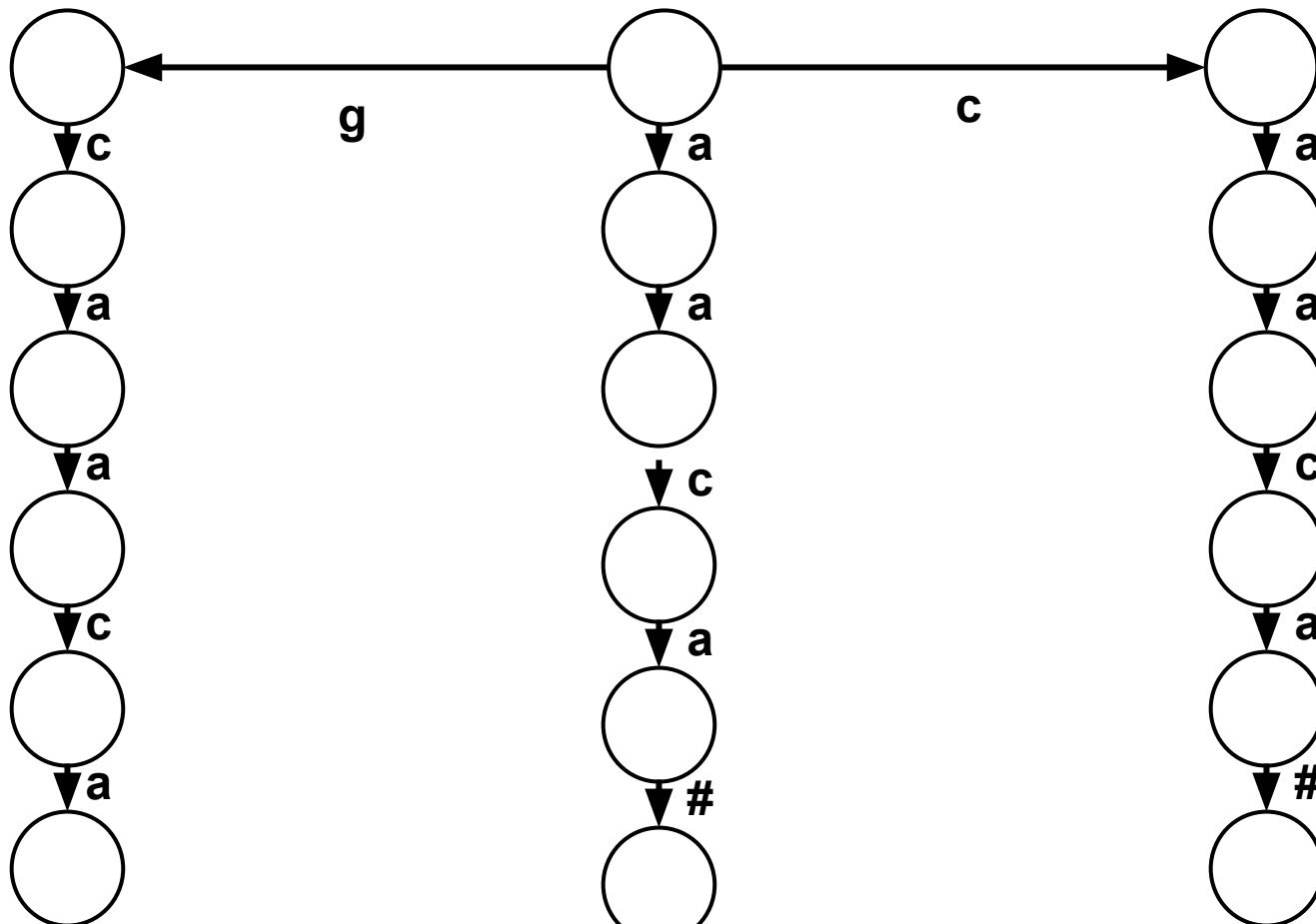
Prefix trie

a c a a c g



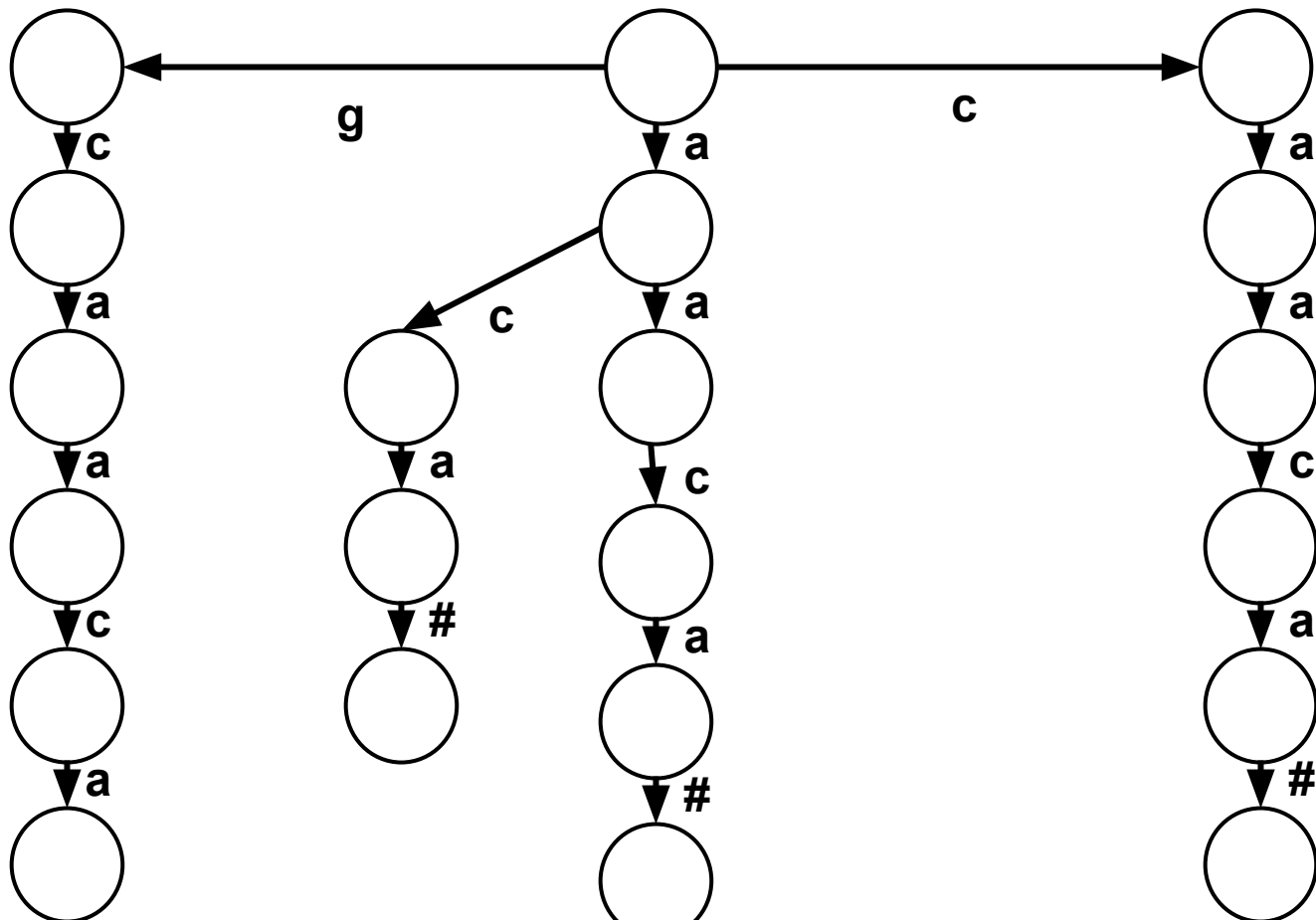
Prefix trie

a c a a c g



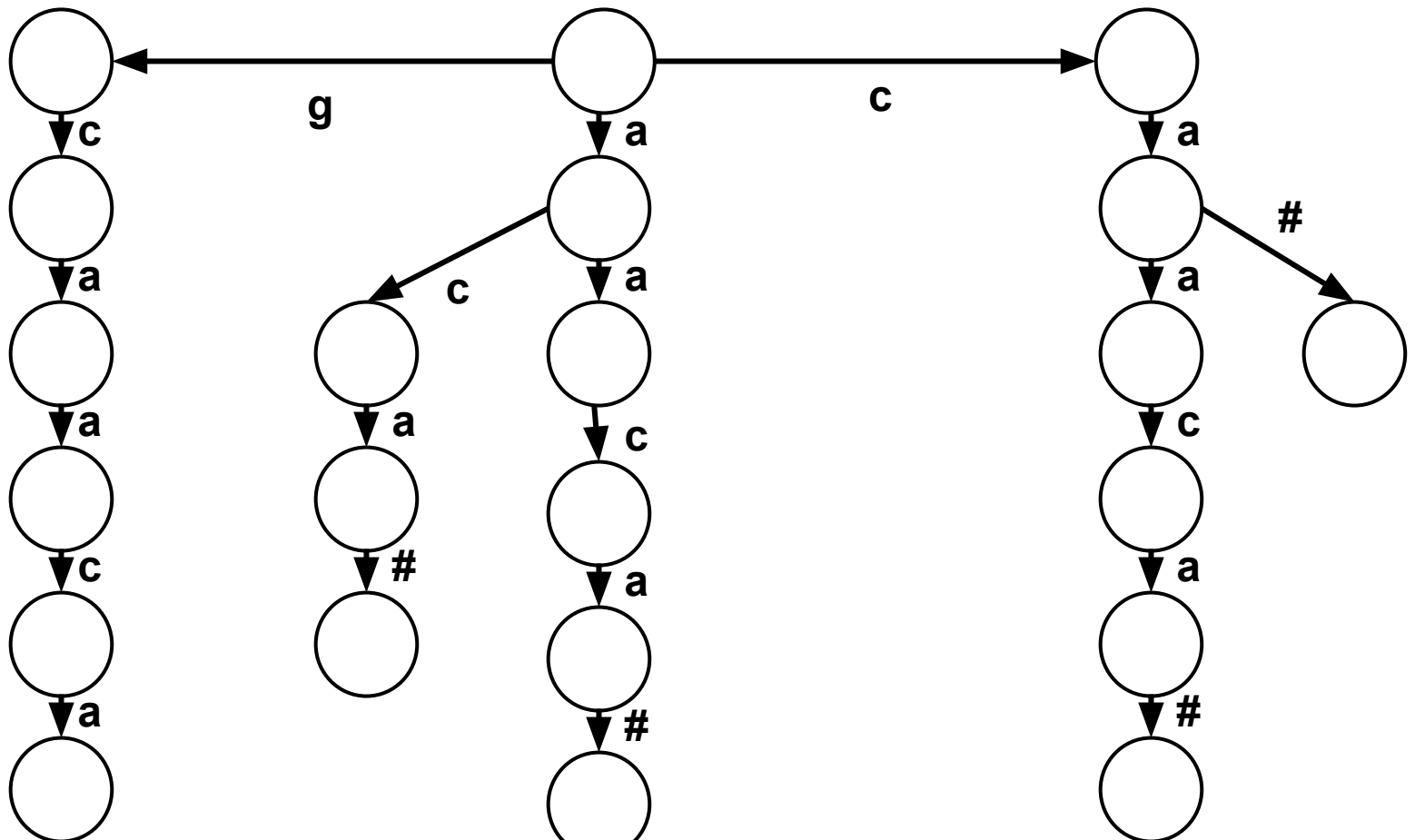
Prefix trie

a c a a c g



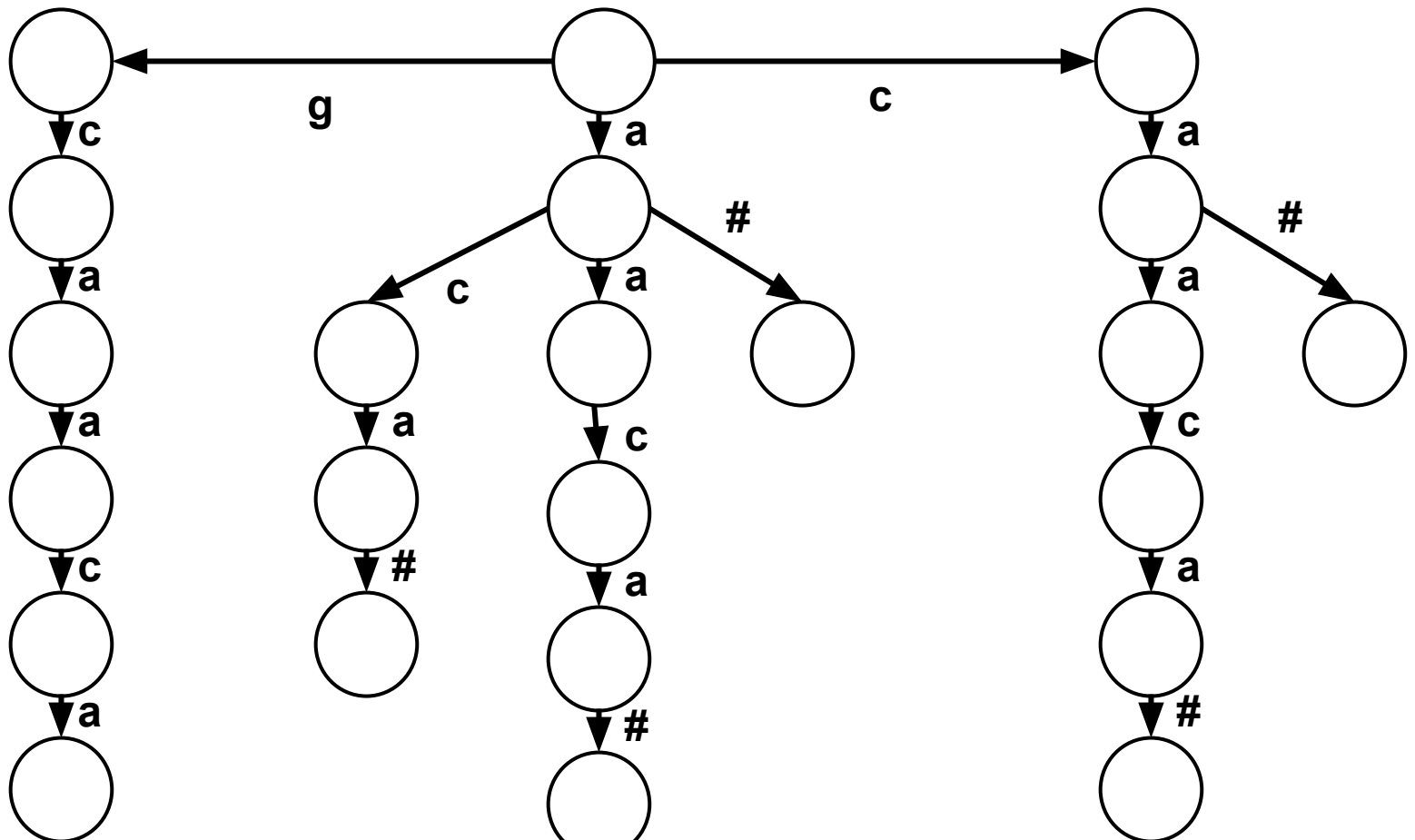
Prefix trie

a c a a c g



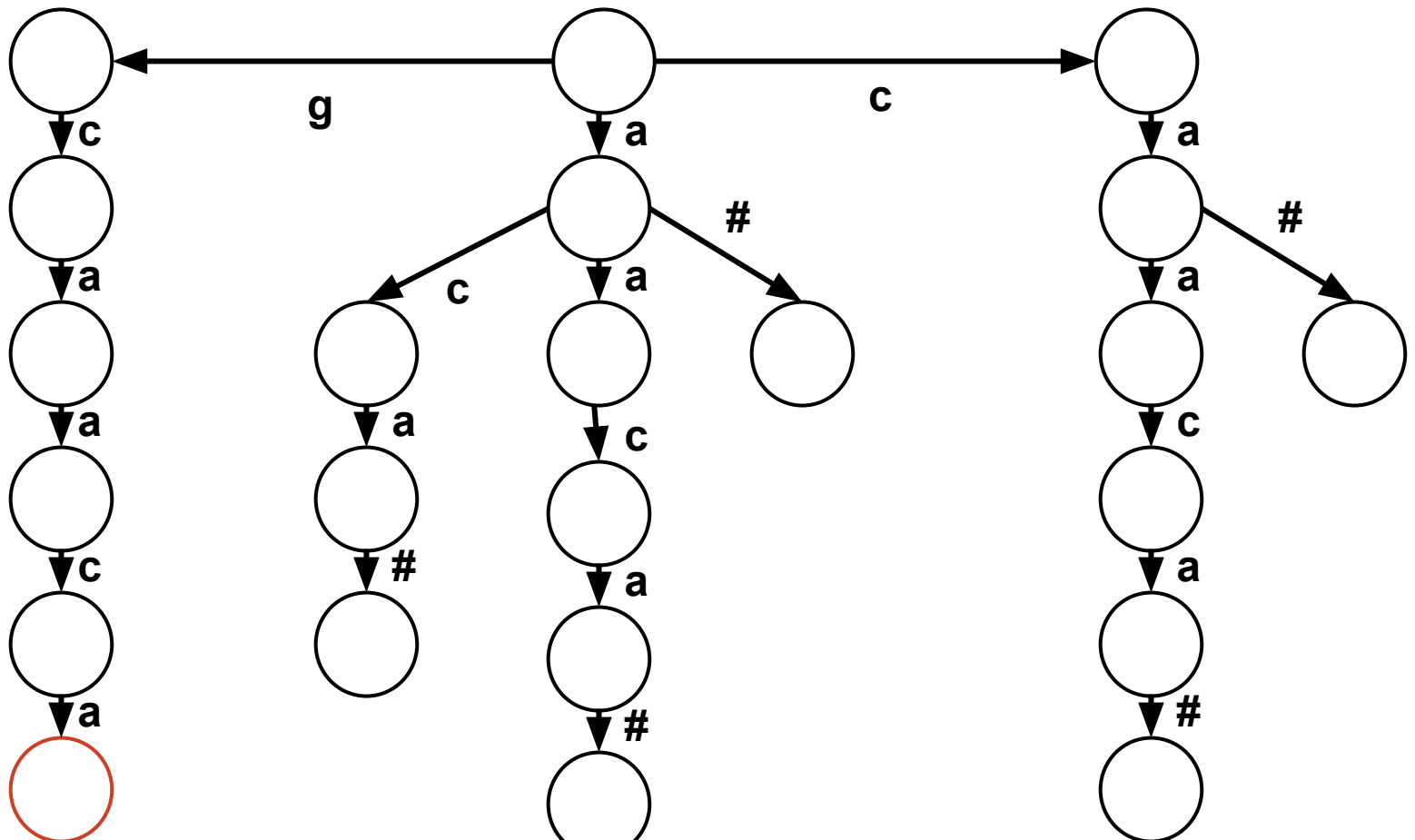
Prefix trie

a c a a c g



Prefix trie

a c a a c g

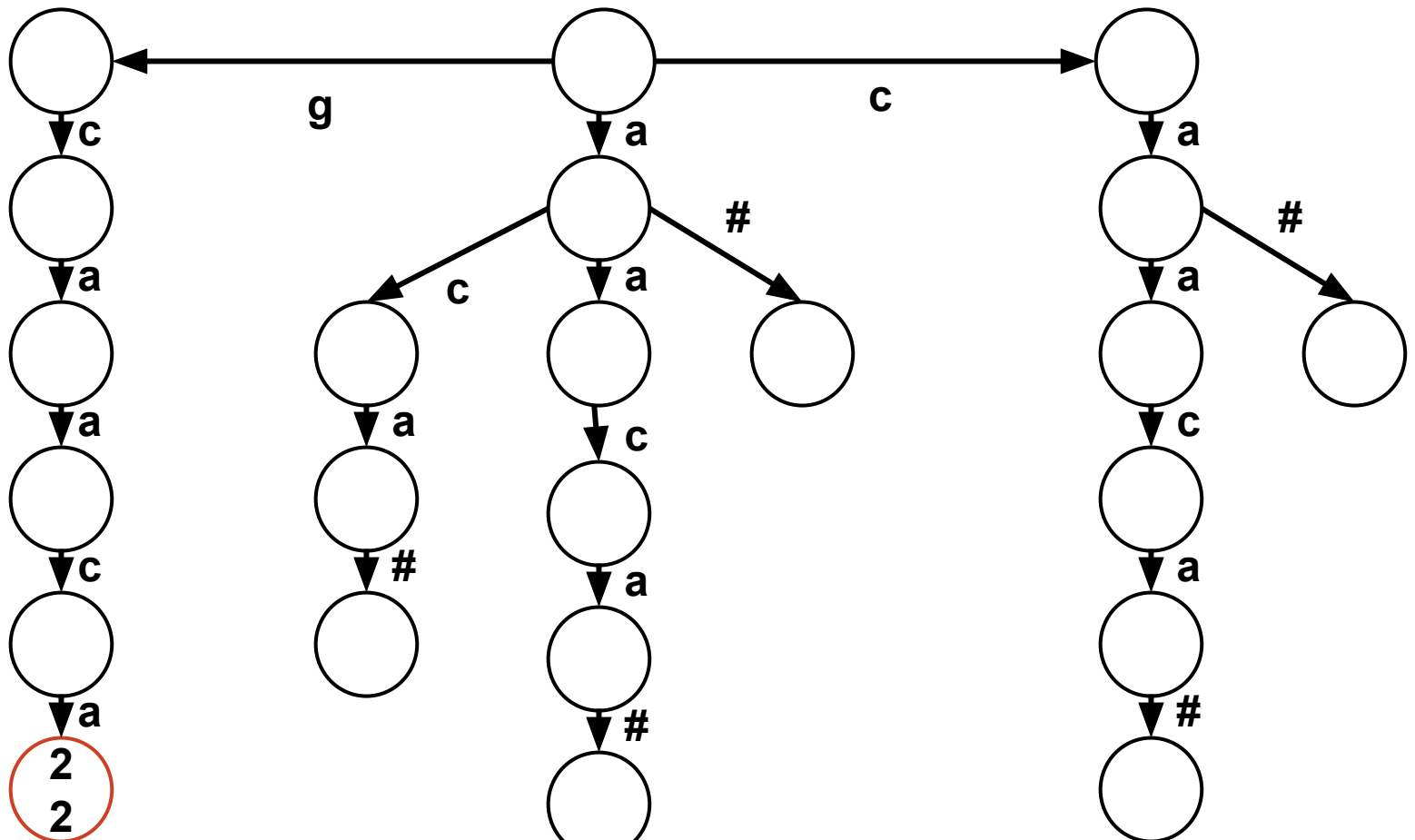


Suffix array

\$
a a c g \$
a c a a c g \$
a c g \$
c a a c g \$
c g \$
g \$

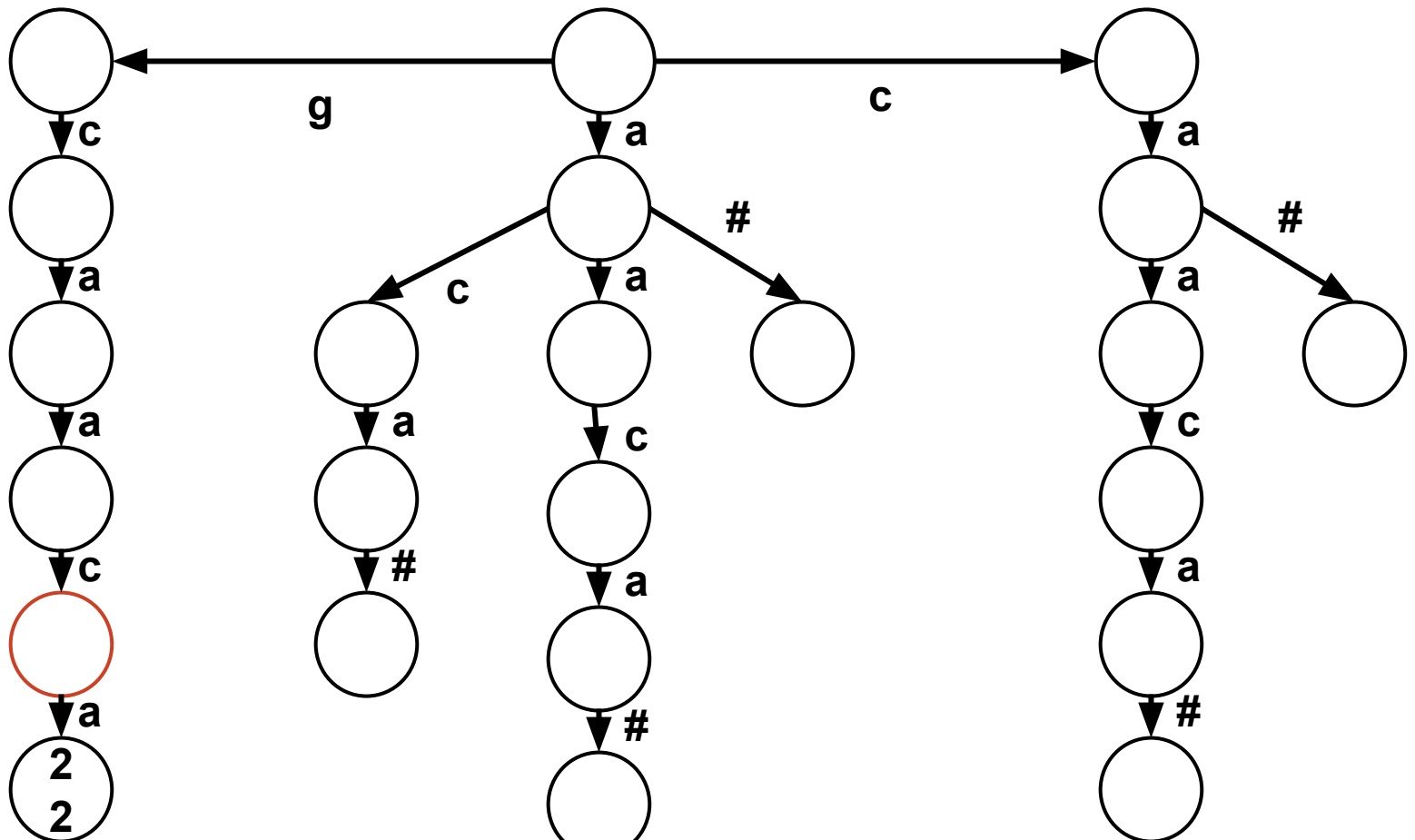
Prefix trie

a c a a c g



Prefix trie

a c a a c g



Suffix array

\$

a a c g \$

a c a a c g \$

a c g \$

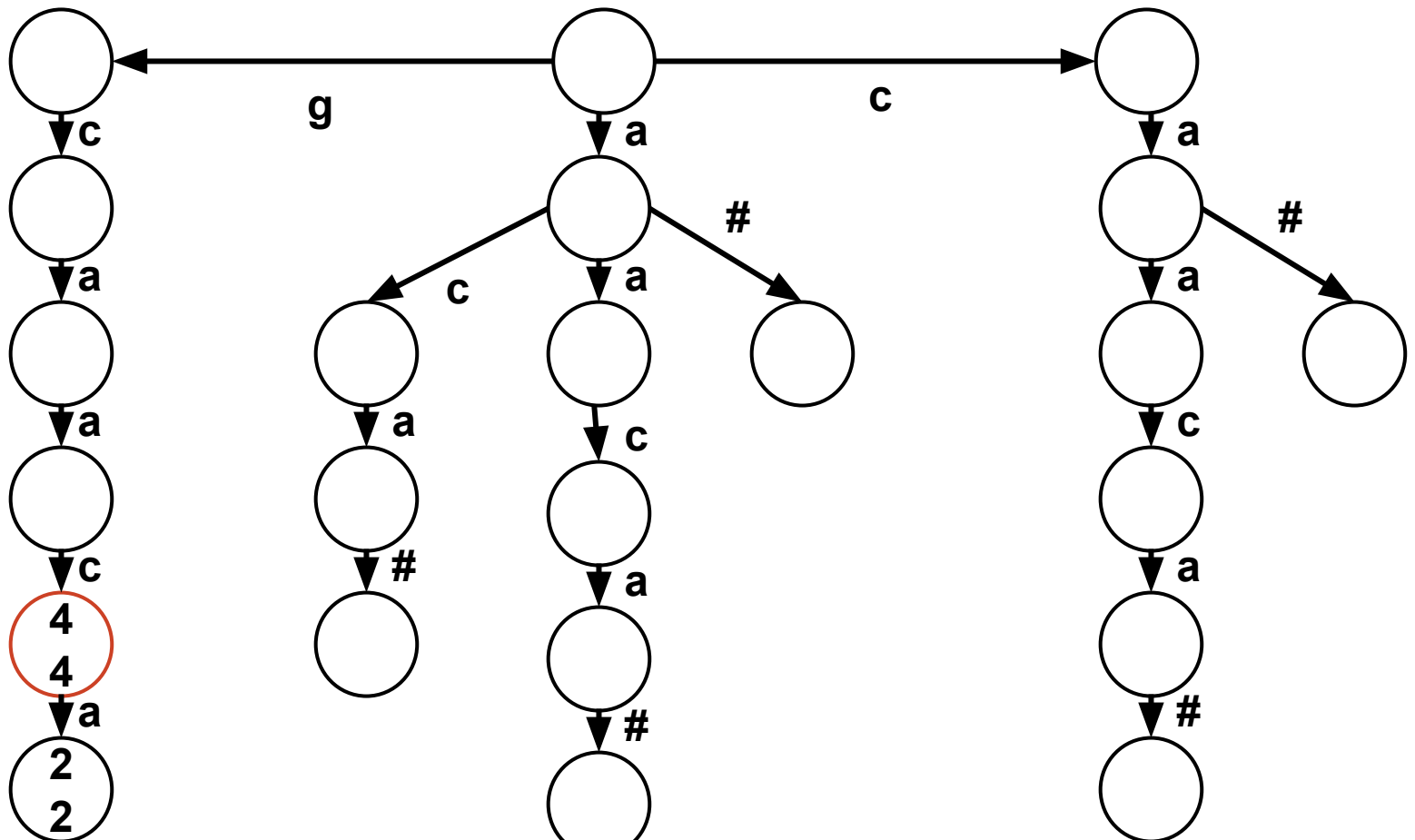
c a a c g \$

c g \$

g \$

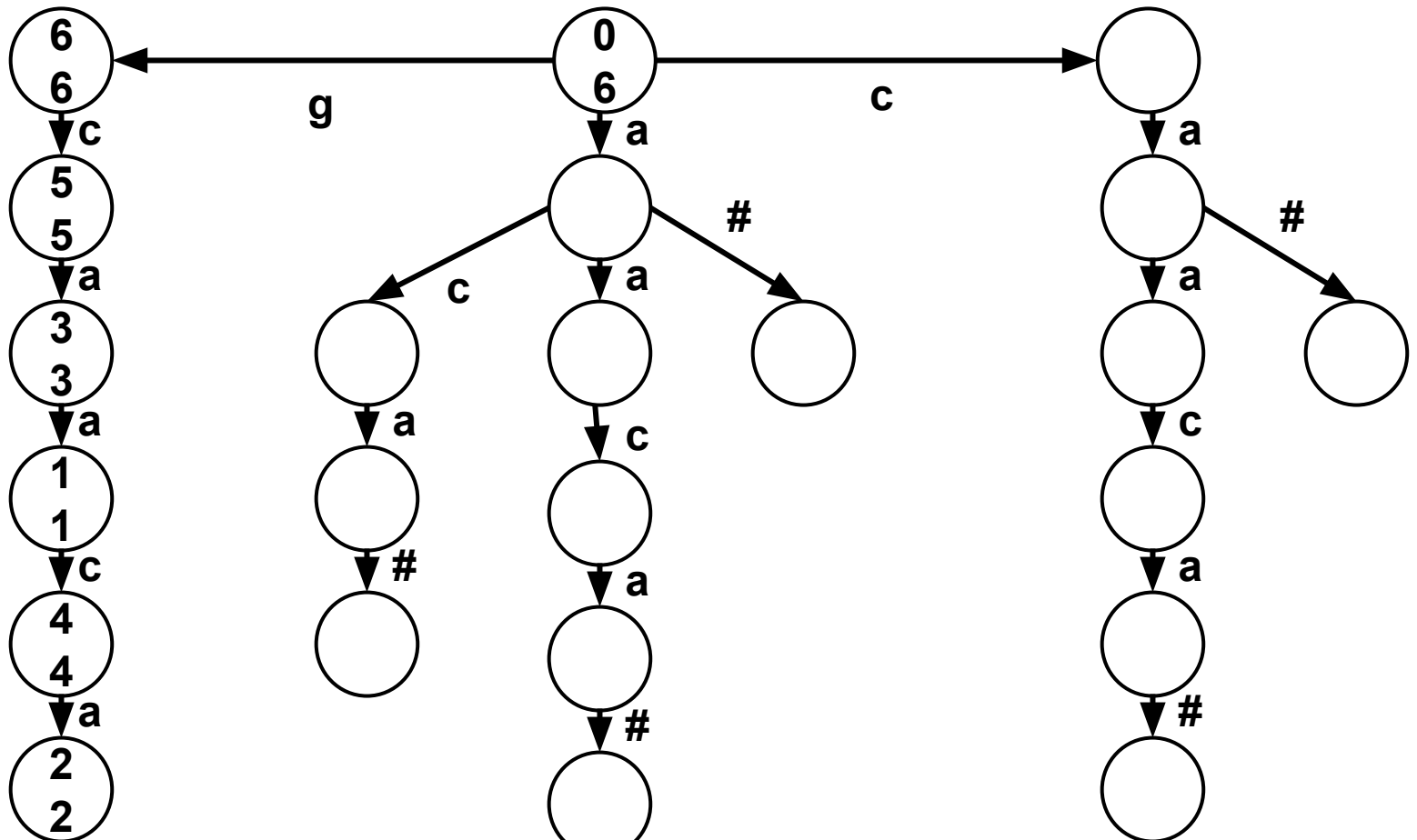
Prefix trie

a c a a c g



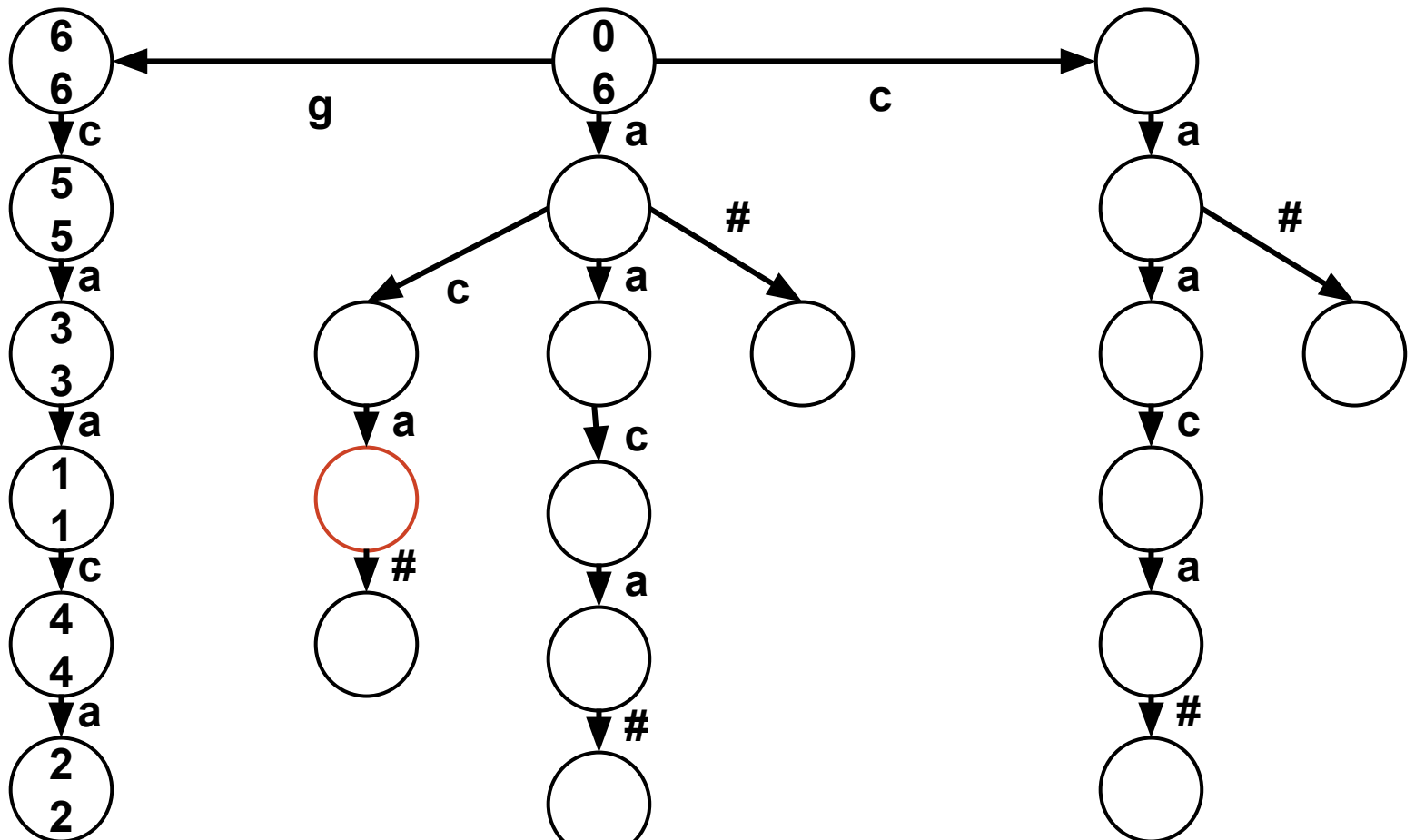
Prefix trie

a c a a c g



Prefix trie

a c a a c g

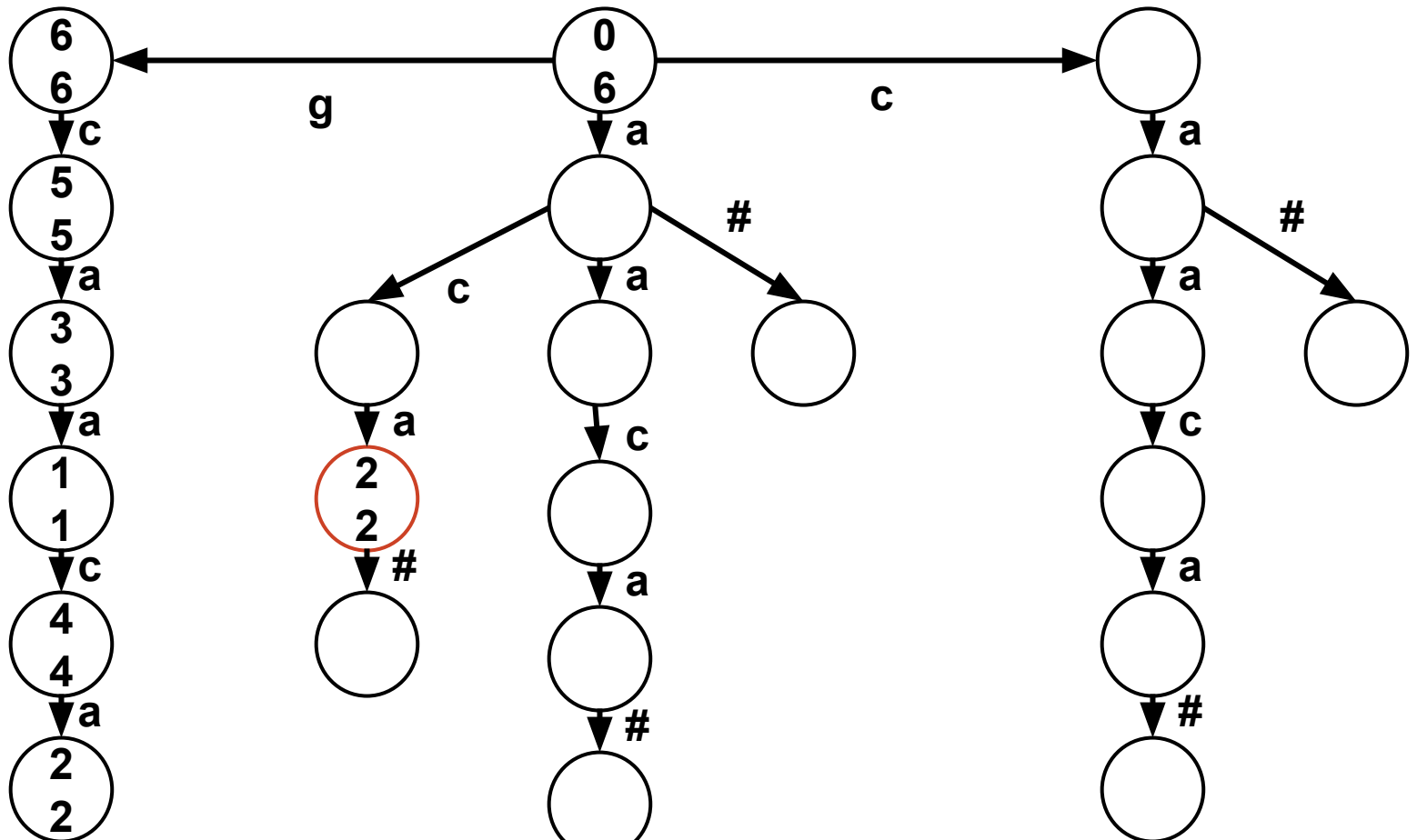


Suffix array

\$
a a c g \$
a c a a c g \$
a c g \$
c a a c g \$
c g \$
g \$

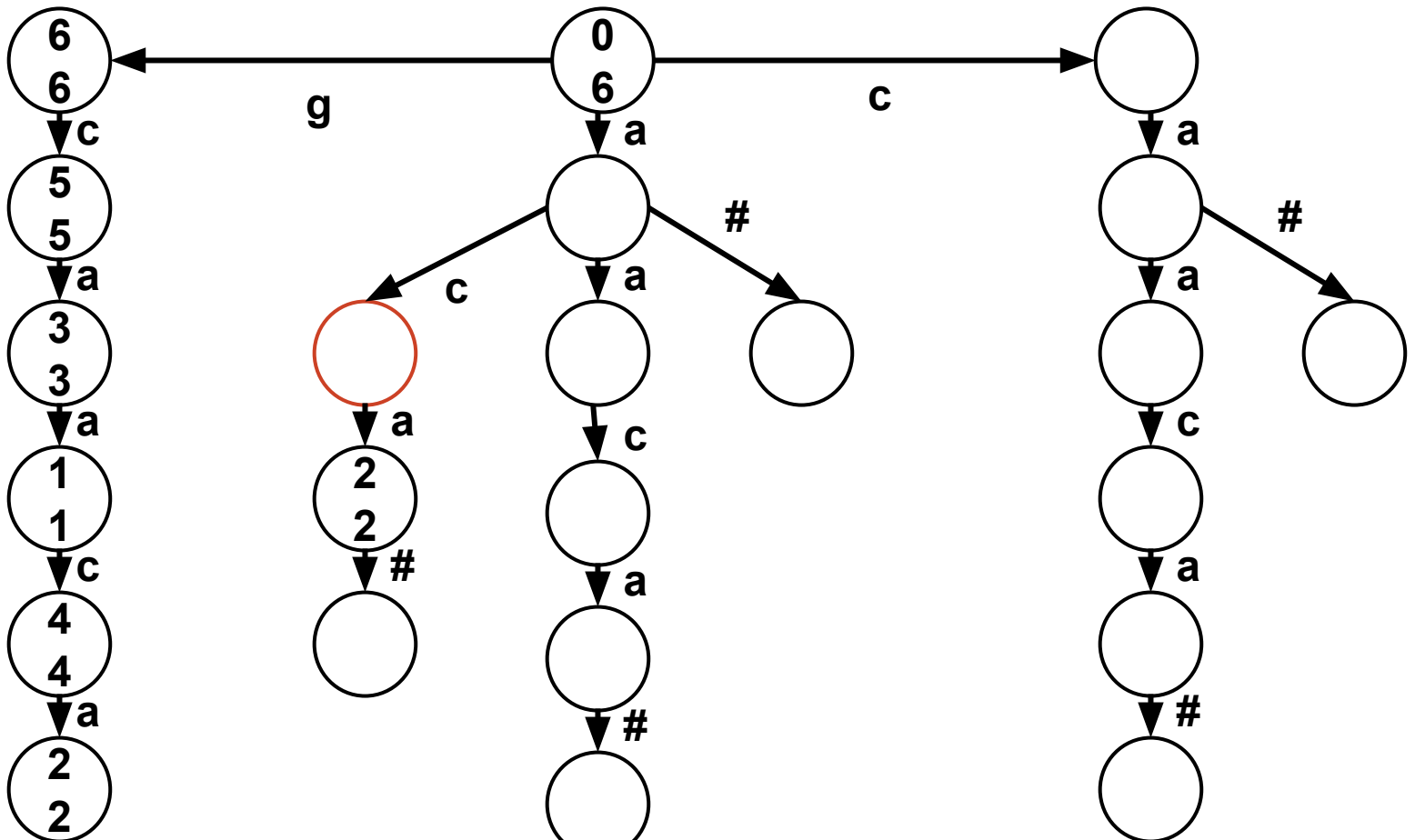
Prefix trie

a c a a c g



Prefix trie

a c a a c g

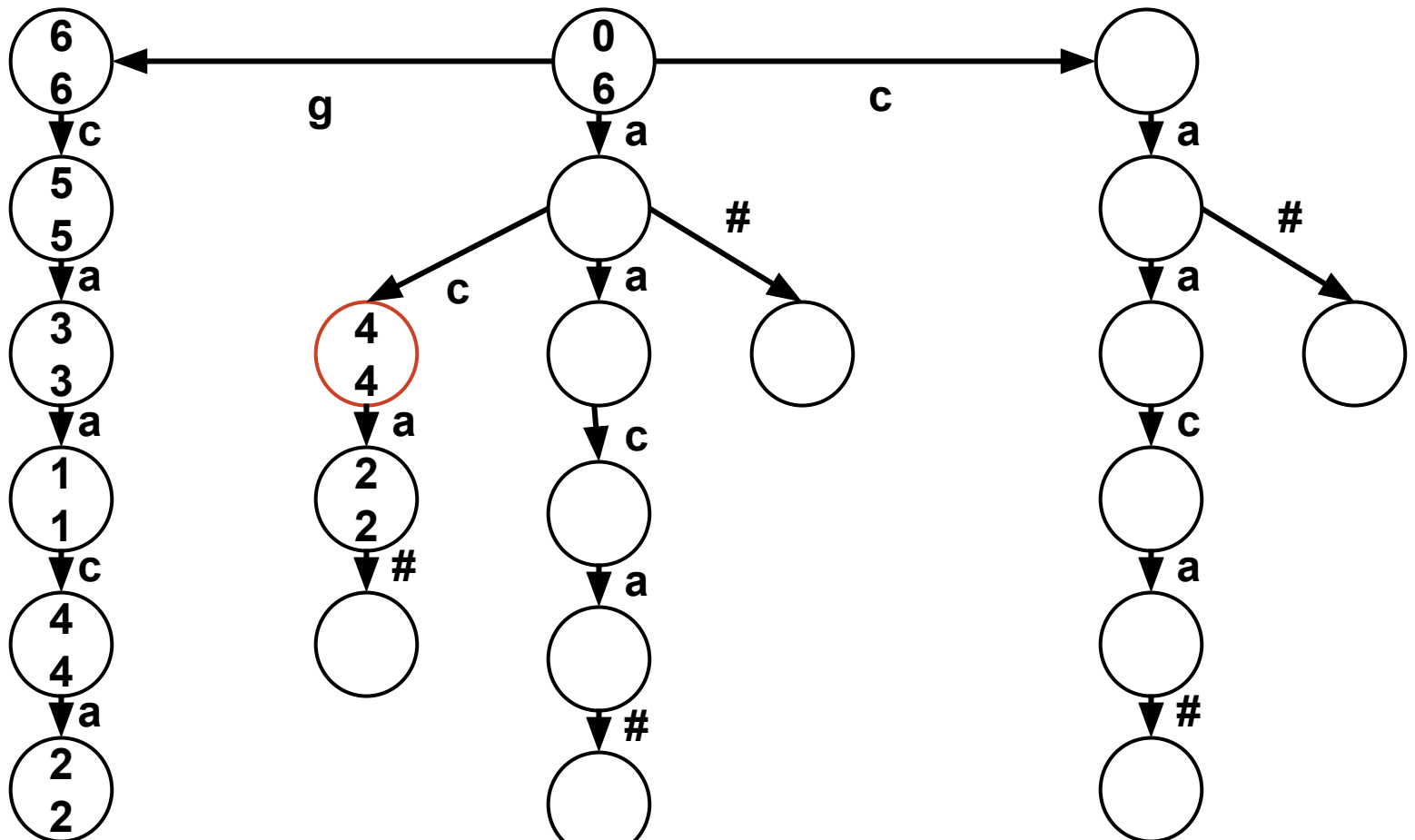


Suffix array

\$
a a c g \$
a c a a c g \$
a c g \$
c a a c g \$
c g \$
g \$

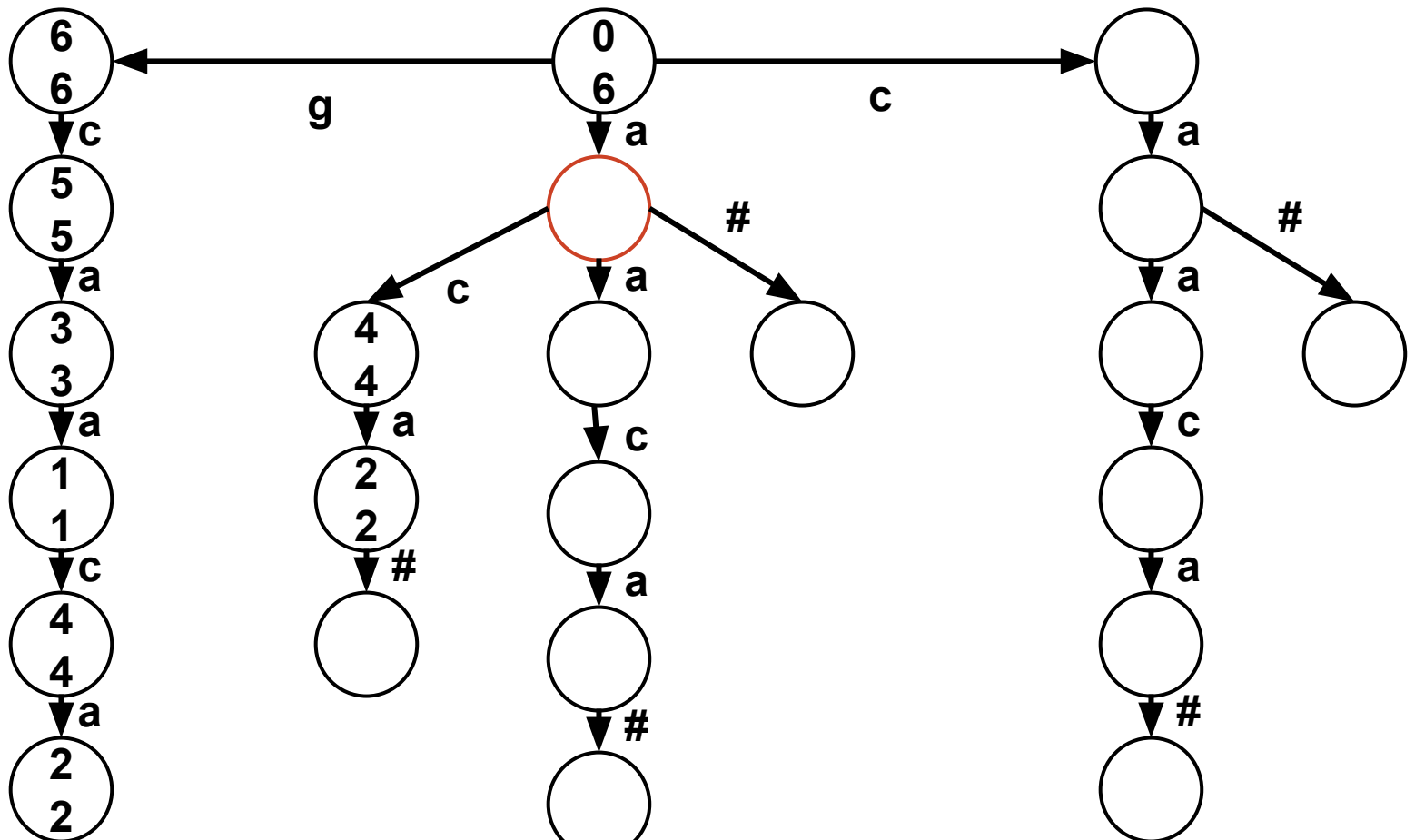
Prefix trie

a c a a c g



Prefix trie

a c a a c g



Suffix array

\$

a a c g \$

a c a a c g \$

a c g \$

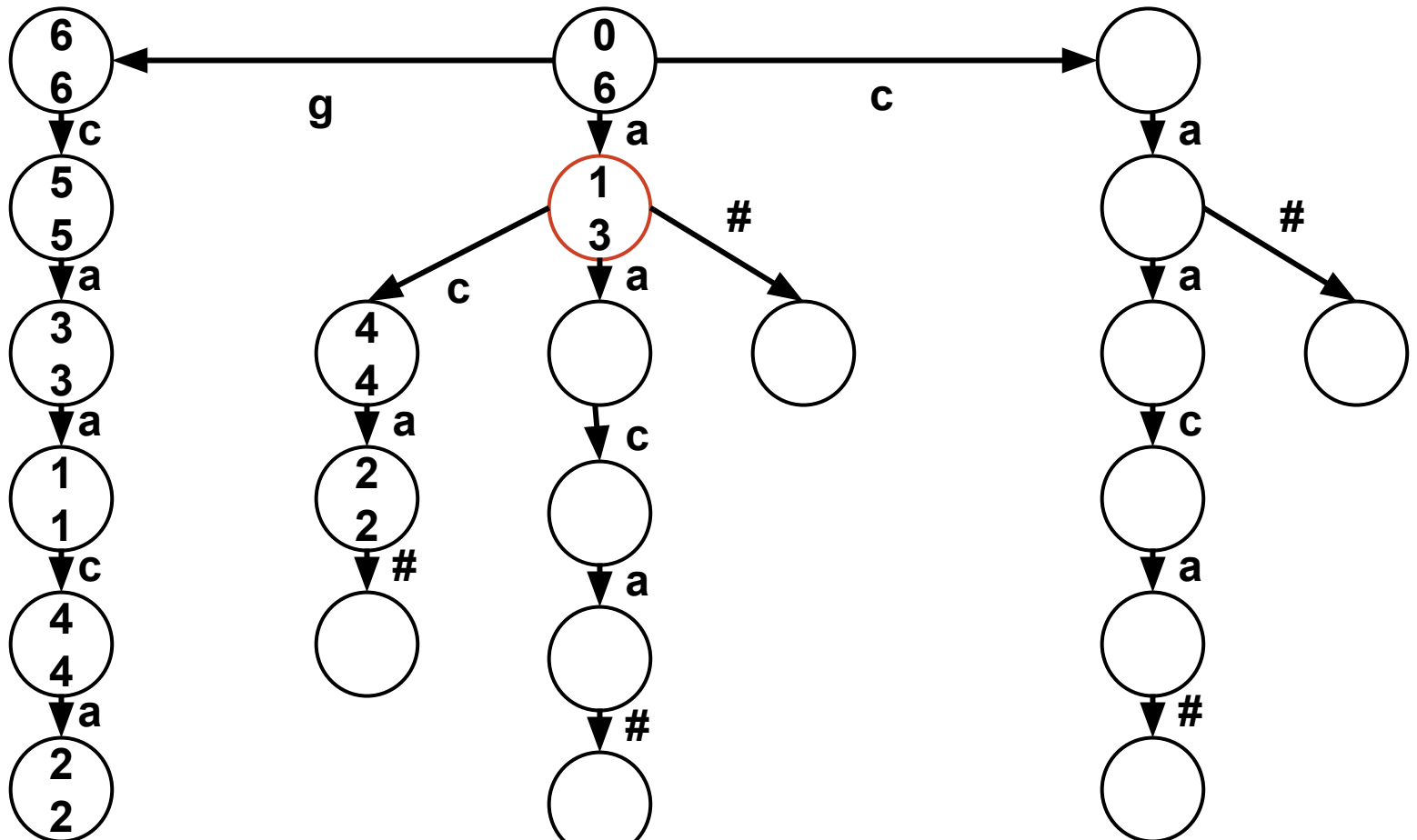
c a a c g \$

c g \$

g \$

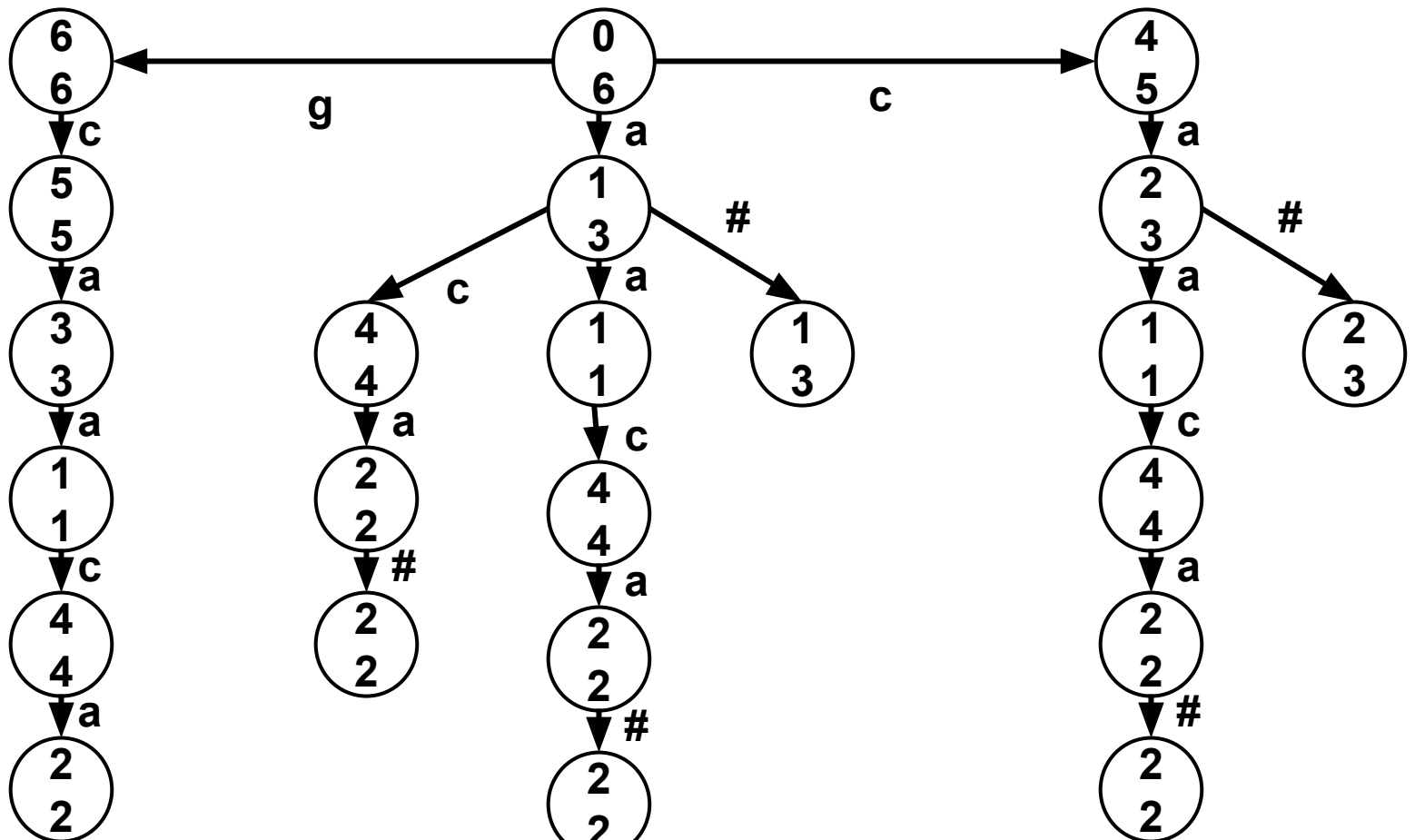
Prefix trie

a c a a c g



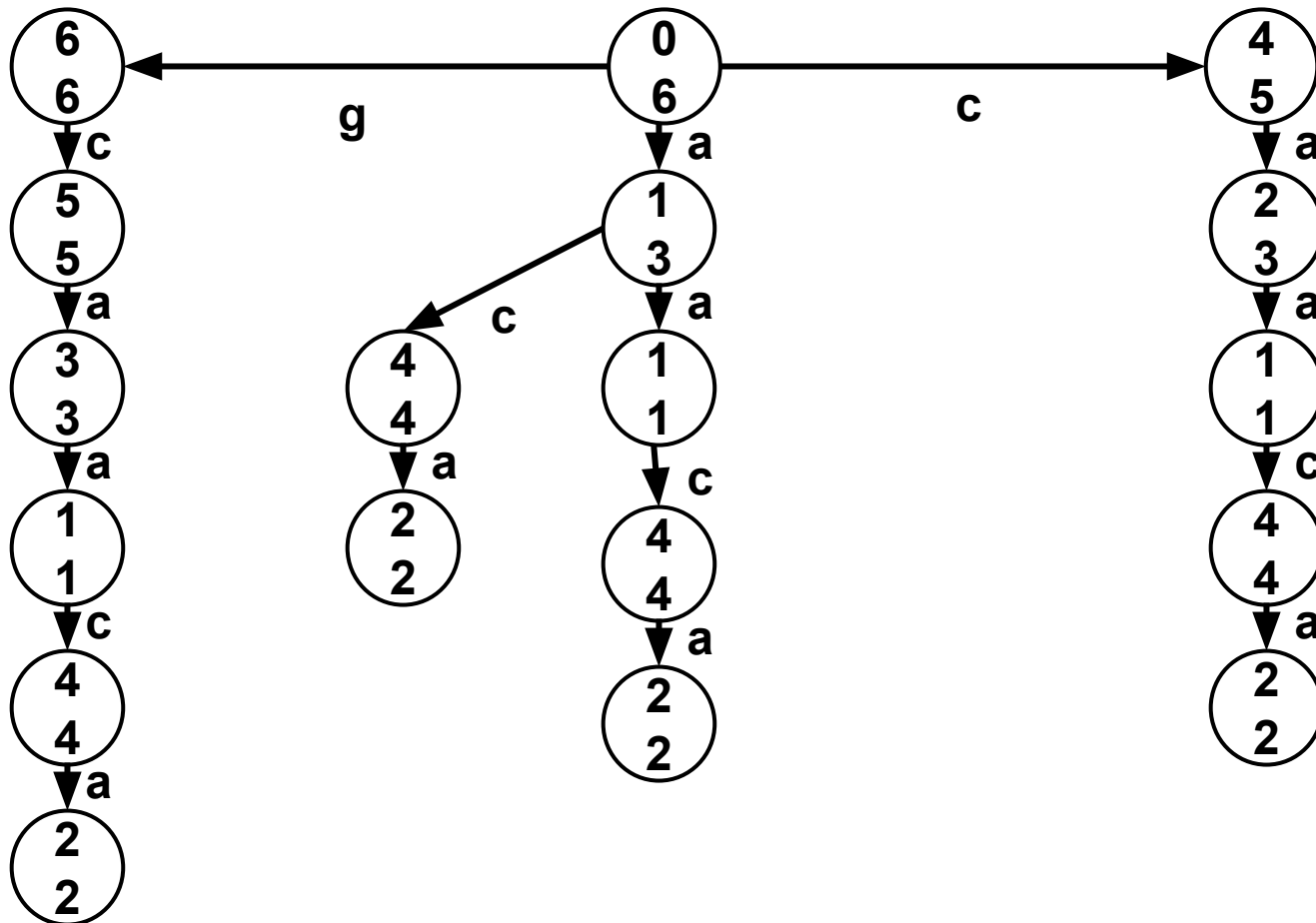
Prefix trie

a c a a c g



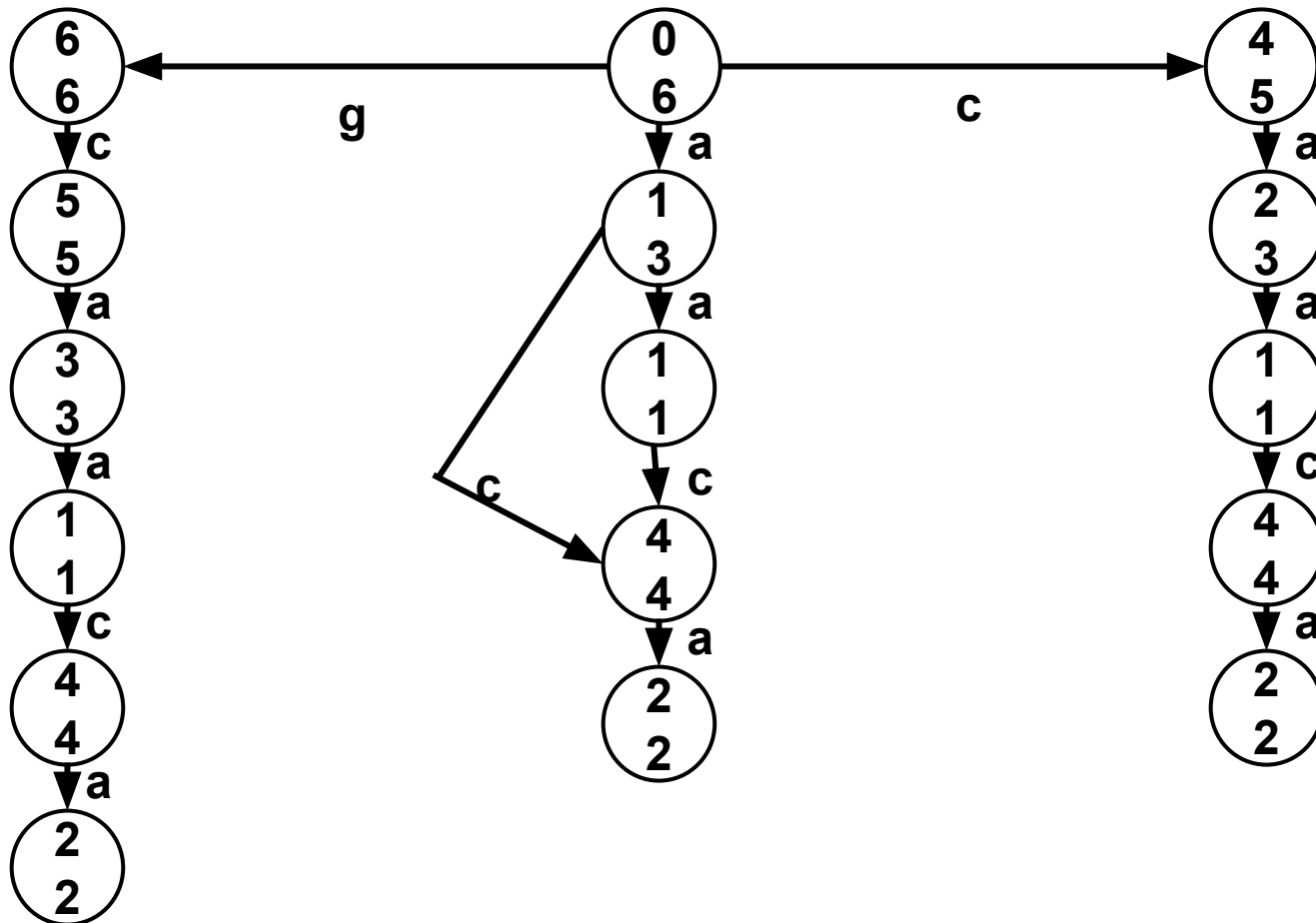
Prefix DAWG

a c a a c g



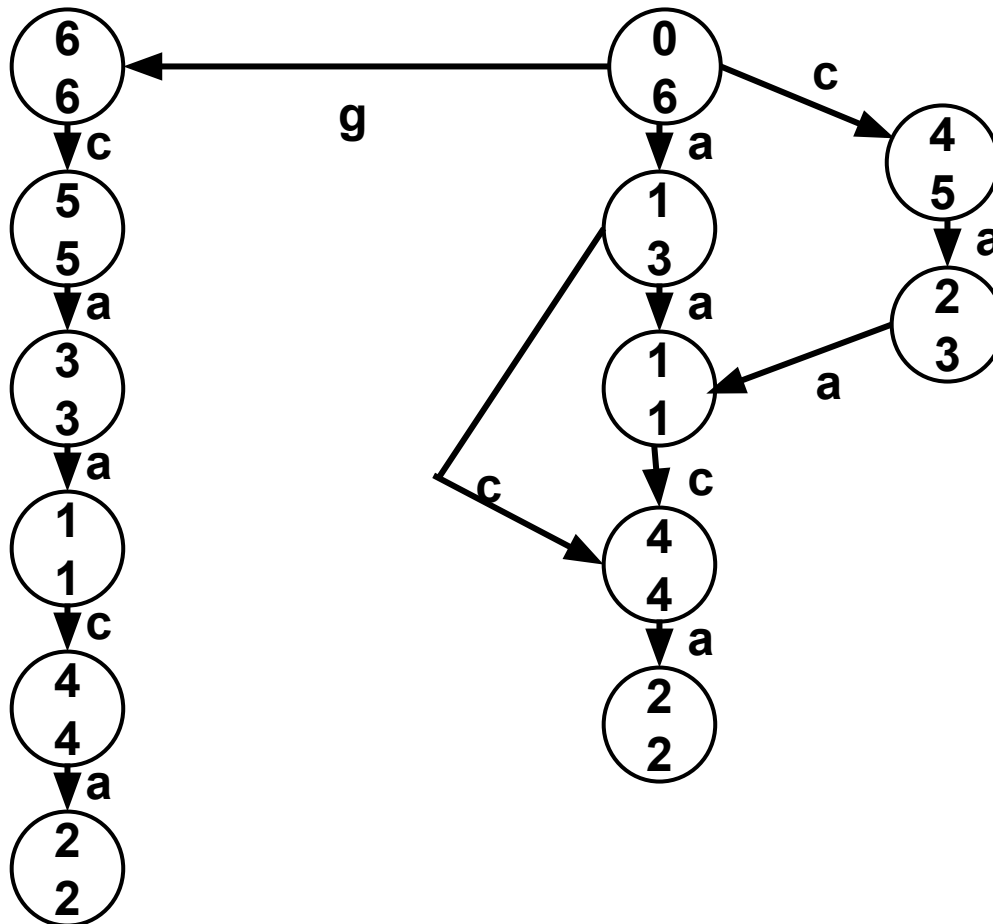
Prefix DAWG

a c a a c g



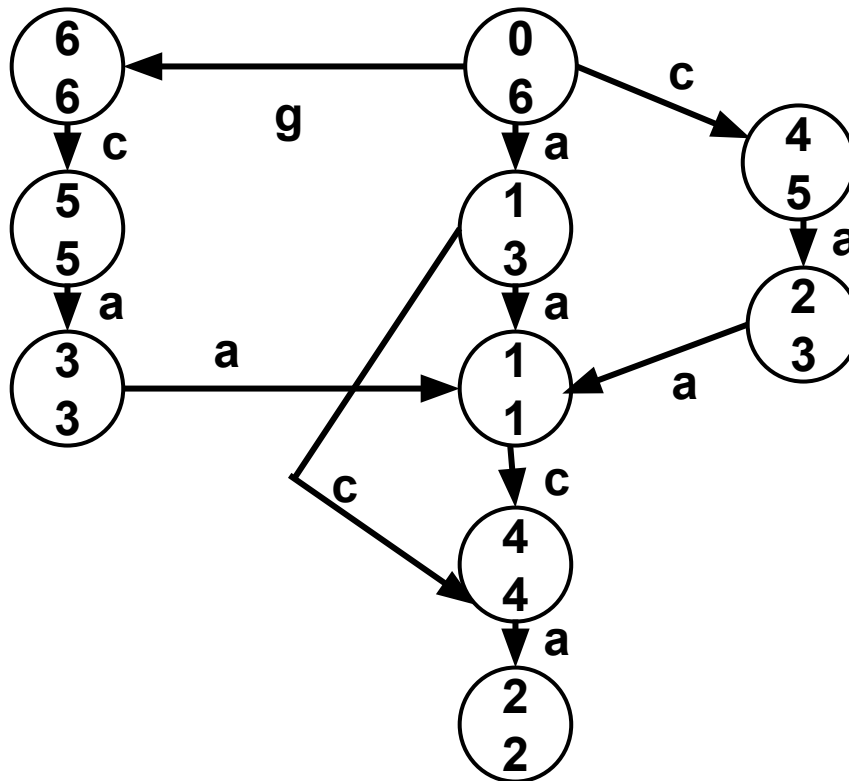
Prefix DAWG

a c a a c g



Prefix DAWG

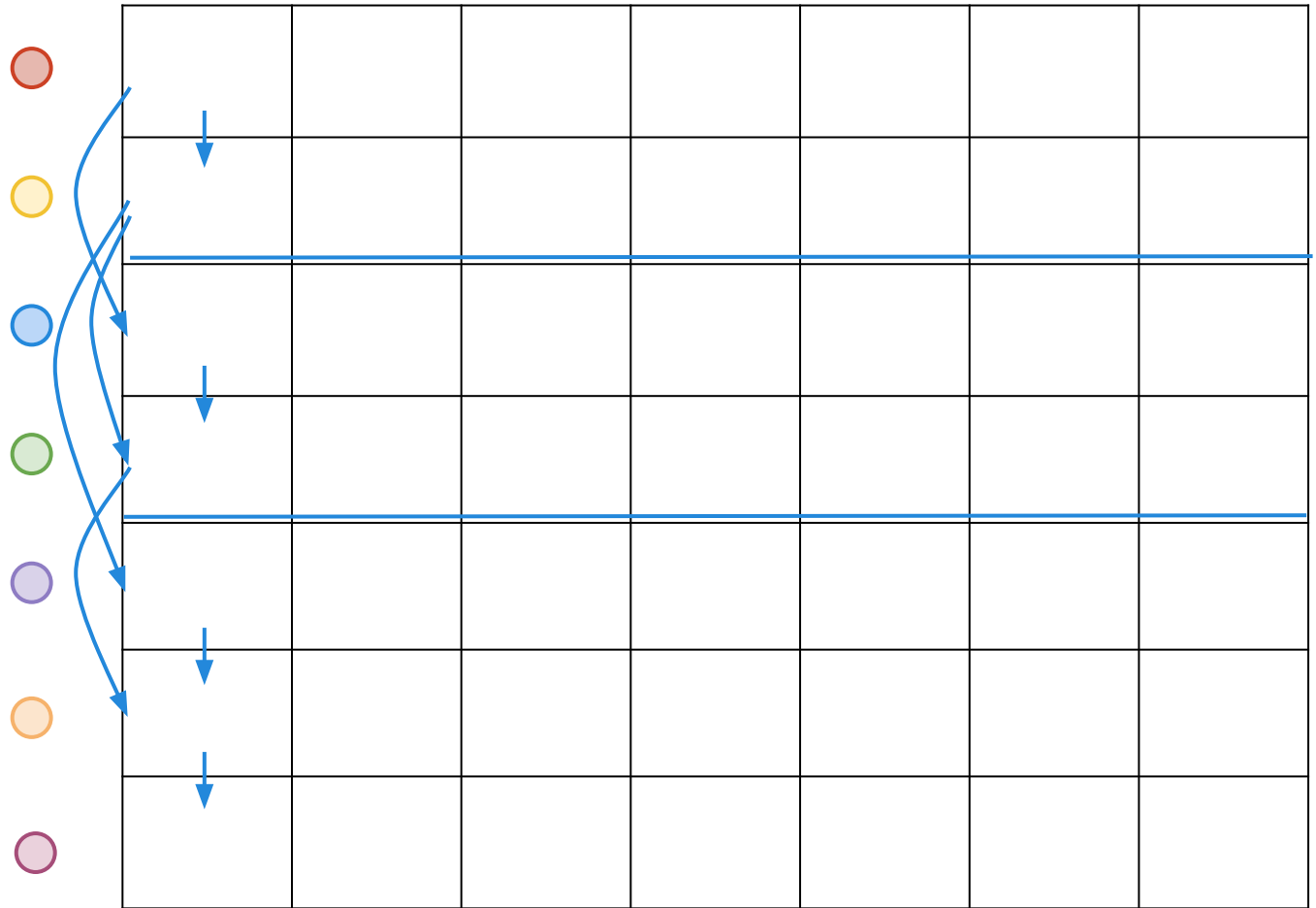
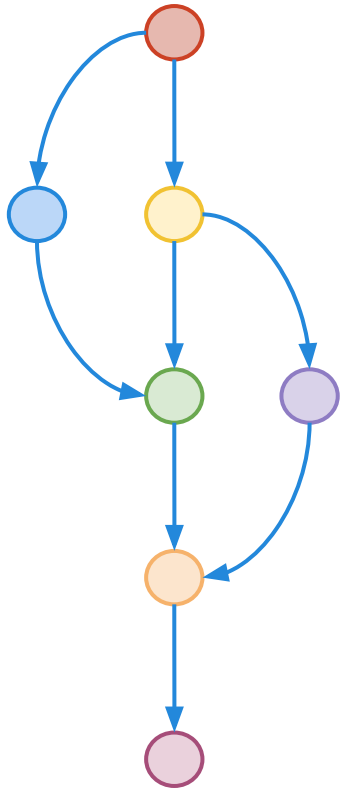
a c a a c g



BWT-SW

- Construct prefix trie $T(X)$ from the reference
- Construct prefix DAWG $\Gamma(W)$ from the query
- Calculate score between $\Gamma(W)$ and $T(X)$
using dynamic programming

BWT-SW



Alignment

- Substitution scoring matrix
 - $S(x, x) = 1$
 - $S(x, y) = -1, x \neq y$

- Affine gap model
 - Gap open score S_{op} ($S_{op} = 1$)
 - Gap extension score S_{ex} ($S_{ex} = 1$)

Dynamic programming

$u \in \Gamma(W)$, u' is a parent of u

$v \in T(X)$, v' is a parent of v

- $I_{uv|u'} = \max \{ I_{u'v} - S_{ex'}, G_{u'v} - S_{op} \}$
- $D_{uv|u'} = \max \{ D_{uv'} - S_{ex'}, G_{uv'} - S_{op} \}$
- $G_{uv|u'} = \max \{ G_{u'v} + S(uu', vv'), I_{uv}, D_{uv} \}$

Dynamic programming

- Select maximal score $G_{uv | u^*}$ (among all parents of u)
- Set $(G_{uv}, l_{uv}, D_{uv}) = (G_{uv | u^*}, l_{uv | u^*}, D_{uv | u^*})$
if $G_{uv | u^*} > 0$
- Set $(G_{uv}, l_{uv}, D_{uv}) = (-\infty, -\infty, -\infty)$
otherwise

Dynamic programming

- Traverse both prefix trie and prefix DAWG from parents to children
- Stop if vertex has non-positive score
- Allows to traverse only a few vertices near the root
- Works much faster than usual Smith-Waterman: $O(n^{0.628} m)$

Acceleration by standard SW

- Once the vertex has
 - High score
 - Small interval
- Traversing becomes excessive
- Switch to usual Smith-Waterman

Z-best strategy

- Keeping all matching vertices for each node in DAWG requires time and memory
- Keep only Z best scoring vertices
- Possibility to miss some seeds (but we rely on other ones)
- Greatly increases performance

BWA tool

- **bwa index** — index construction
- **bwa aln** — short read aligner
- **bwa sw** — long read aligner
- **bwa mem** — new long read and long sequence local aligner

Alignment applications

Alignment applications

- **Quality assessment**
 - Error rate
 - Insert size distribution
 - Chimeric read/read-pairs
 - Genome fraction
- **SNP calling**
- **Comparative analysis**
 - CNVs
- **Transcriptomics**
 - Gene expression
 - Exon/intron detection

Storing alignments

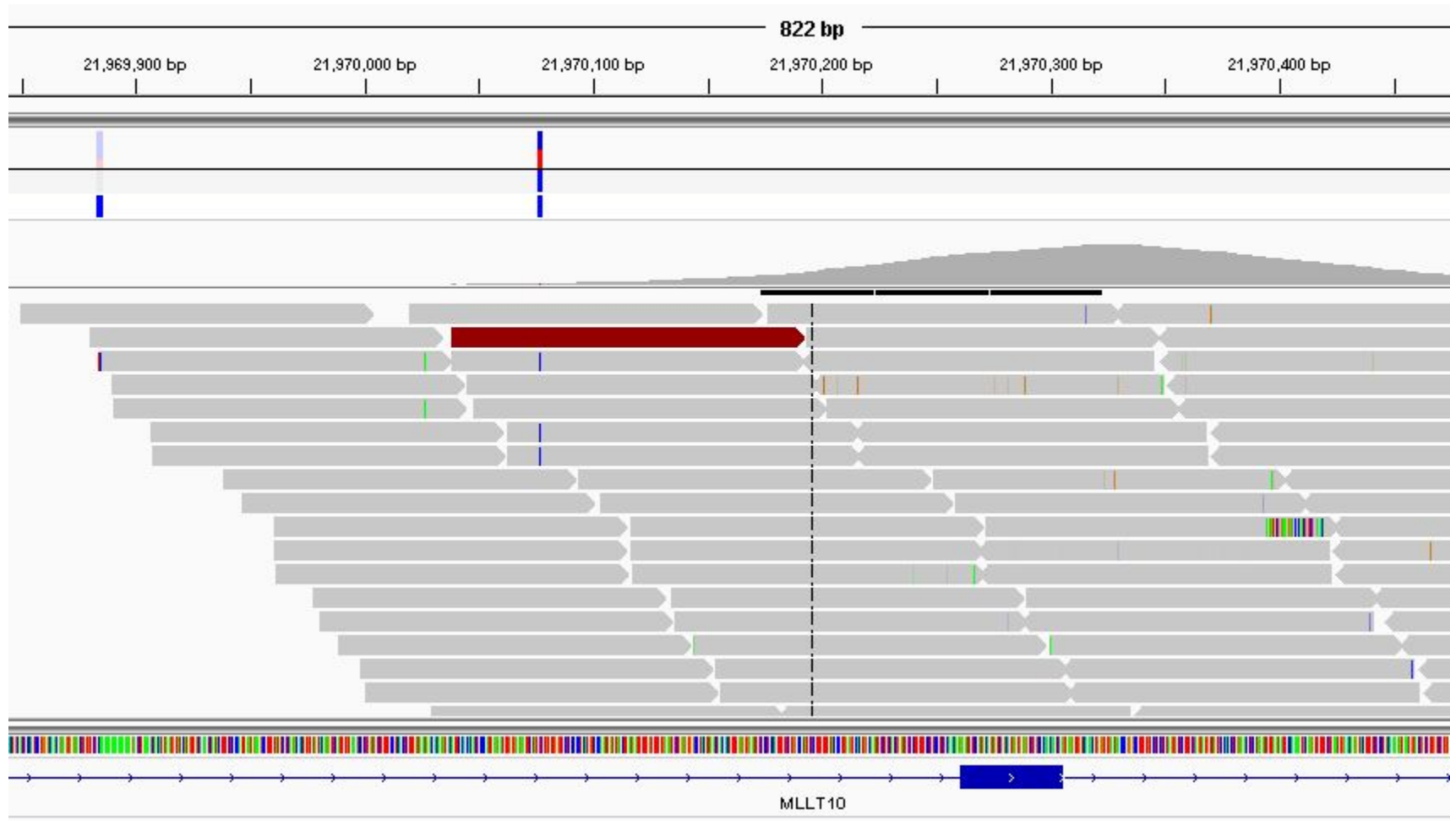
SAM/BAM files

- Read ID (QNAME)
- Reference ID (RNAME)
- Mapping position (POS)
- Mate reference ID (RNEXT)
- Mate position (PNEXT)
- Observed insert length (TLEN)
- Read sequence (SEQ)
- Read quality (QUAL)
- CIGAR string
 - 34M 1I 4M 2D 1X 3M

SAM files

```
@HD VN:1.0 S0:coordinate
@SQ SN:chr20 LN:64444167
@PG ID:TopHat VN:2.0.14 CL:/srv/dna_tools/tophat/tophat -N 3 --read-edit-dist 5 --read-rea
lign-edit-dist 2 -i 50 -I 5000 --max-coverage-intron 5000 -M -o out /data/user446/mapping_tophat/index/chr
20 /data/user446/mapping_tophat/L6_18_GTGAAA_L007_R1_001.fastq
HWI-ST1145:74:C101DACXX:7:1102:4284:73714 16 chr20 190930 3 100M * 0 0
CCGTGTTTAAAGGTGGATGCGGTCACCTCCCAGCTAGGCTTAGGGATTCTTAGTTGGCCTAGGAAATCCAGCTAGTCCTGTCTCTCAGTCCCCCTCT
C BBDDCCDDCCDDDDDCDDDDDDDCDDCCDBC?DDDDDDDDDDDDDDDCDDDDDDDDDDCCCEDDDC?DDDDDDDDDDDDDDDDDDDDDDBDHFFFFDC@@
AS:i:-15 XM:i:3 X0:i:0 XG:i:0 MD:Z:55C20C13A9 NM:i:3 NH:i:2 CC:Z:= CP:i:55352714 HI:i:0
HWI-ST1145:74:C101DACXX:7:1114:2759:41961 16 chr20 193953 50 100M * 0 0
TGCTGGATCATCTGGTTAGTGGCTTCTGACTCAGAGGACCTTCGTCCTGGGGCAGTGGACCTTCCAGTGATTCCCCTGACATAAGGGGCATGGACGA
G DCDDDEDDDDDDDCDDDDDDDCDDDDDEEC>DFFFEJJJJJIGJJJJIHGBHHGJIJJJJJJGJJJIJJJJJIHJJJJJJHHHHHFFFFFFCC
AS:i:-16 XM:i:3 X0:i:0 XG:i:0 MD:Z:60G16T18T3 NM:i:3 NH:i:1
HWI-ST1145:74:C101DACXX:7:1204:14760:4030 16 chr20 270877 50 100M * 0 0
GGCTTTATTGGTAAAAAAGGAATAGCAGATTTAATCAGAAATCCCACCTGGCCCAGCAGCACCAACCAGAAAGAAGGGAAGAAGACAGGAAAAAACCA
C DDDDDDDDDCCDDDDDDDDDEEEEEEEFFFEFFEGHHHFGDJJIHJJJIJJJJIIIGGFJJJIHIIIIJJJJJJIGHHFAHGFHJHFGGHFFFD@BB
AS:i:-11 XM:i:2 X0:i:0 XG:i:0 MD:Z:0A85G13 NM:i:2 NH:i:1
HWI-ST1145:74:C101DACXX:7:1210:11167:8699 0 chr20 271218 50 50M4700N50M * 0
0 GTGGCTCTTCCACAGGAATGTTGAGGATGACATCCATGTCTGGGGTGCACCTGGGTCTCCGAAGCAGAACATCCTCAAATATGACCTCTCG
accepted_hits.sam
```

Alignment visualization with IGV



SNP calling

Process of finding bases in the NGS data that differ from the reference genome

- Typically including an associated statistical confidence score
- Also known as “variant calling”

We need enough coverage to distinguish real variants from sequencing errors

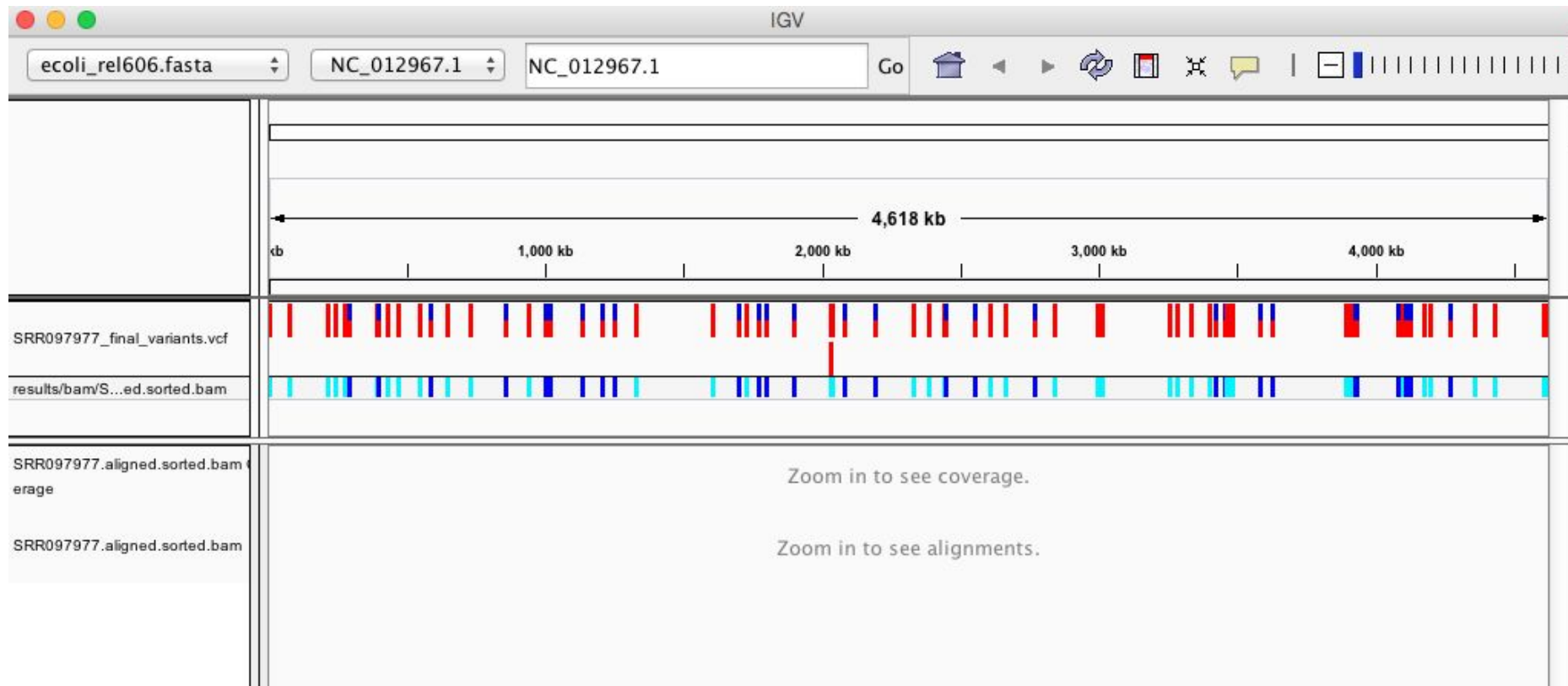
VCF/BCF files

- Chromosome (#CHROM)
- Position (POS)
- Unique identifiers where available (ID)
- Reference base(s) (REF)
- Alternate non-reference alleles (ALT)
- Phred quality score for the variant (QUAL)
- Optional filters (FILTER)
- Additional information (INFO)

VCF files

#CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO
1	10177	.	A	AC	100	PASS	AC=2130;AF=0.425319;AN=5008;NS=2504
1	10235	.	T	TA	100	PASS	AC=6;AF=0.00119808;AN=5008;NS=2504
1	10352	rs145072688	T	TA	100	PASS	AC=2191;AF=0.4375;AN=5008;NS=2504
1	10505	.	A	T	100	PASS	AC=1;AF=0.000199681;AN=5008;NS=2504
1	10506	.	C	G	100	PASS	AC=1;AF=0.000199681;AN=5008;NS=2504
1	10511	.	G	A	100	PASS	AC=1;AF=0.000199681;AN=5008;NS=2504
1	10539	.	C	A	100	PASS	AC=3;AF=0.000599042;AN=5008;NS=2504
1	10542	.	C	T	100	PASS	AC=1;AF=0.000199681;AN=5008;NS=2504
1	10579	.	C	A	100	PASS	AC=1;AF=0.000199681;AN=5008;NS=2504
1	10616	rs376342519	CCGCCGTTGCAAAGGCGGCCG	C	100	PASS	AC=4973;AF=0.993011;AN=5008;NS=2504
1	10642	.	G	A	100	PASS	AC=21;AF=0.00419329;AN=5008;NS=2504
1	11008	.	C	G	100	PASS	AC=441;AF=0.0880591;AN=5008;NS=2504
1	11012	.	C	G	100	PASS	AC=441;AF=0.0880591;AN=5008;NS=2504
1	11063	.	T	G	100	PASS	AC=15;AF=0.00299521;AN=5008;NS=2504
1	13011	.	T	G	100	PASS	AC=3;AF=0.000599042;AN=5008;NS=2504
1	13110	.	G	A	100	PASS	AC=134;AF=0.0267572;AN=5008;NS=2504

Alignment visualization with IGV: one sample



Tools

- Alignment and data processing
 - samtools
 - bcftools
- SNP calling and annotation
 - VarScan
 - SnpEff
- Visualization
 - Tablet
 - IGV
- Pipelines
 - GATK

Thank you!

Questions?