



# Short read alignment

**Andrey Prjibelski**  
**Center for algorithmic biotechnology and**  
**bioinformatics**  
**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

$a_k c_i a_m a_n c_j g_x \$$

# First-last property

$\$$ <sub>1</sub>	a	c	a	a	c	g
a <sub>1</sub>	a	c	g	\$	a	c
a <sub>2</sub>	c	a	a	c	g	\$
a <sub>3</sub>	c	g	\$	a	c	a
c <sub>1</sub>	a	a	c	g	\$	a
c <sub>2</sub>	g	\$	a	c	a	a
g <sub>1</sub>	\$	a	c	a	a	c

# First-last property

$\$_1$  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

$\$_1$  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**  
 $c_1$  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**  $a_3$  c g \$ a c  
**a**  $c_1$  a a c g \$  
**a**  $c_2$  g \$ a c 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**  $a_3$  c g \$ a c  
**a**  $c_1$  a a c g \$  
**a**  $c_2$  g \$ a c 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

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

$a_3$  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	$g_1$
$a_1$	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	$c_2$

# First-last property

$a_2 \ c_1 \ a_1 \ a_3 \ c_2 \ g_1 \ \$_1$

# First-last property

$\$1$	$a_2$	$c_1$	$a_1$	$a_3$	$c_2$	$g_1$
$a_1$	$a_3$	$c_2$	$g_1$	$\$1$	$a_2$	$c_1$
$a_2$	$c_1$	$a_1$	$a_3$	$c_2$	$g_1$	$\$1$
$a_3$	$c_2$	$g_1$	$\$1$	$a_2$	$c_1$	$a_1$
$c_1$	$a_1$	$a_3$	$c_2$	$g_1$	$\$1$	$a_2$
$c_2$	$g_1$	$\$1$	$a_2$	$c_1$	$a_1$	$a_3$
$g_1$	$\$1$	$a_2$	$c_1$	$a_1$	$a_3$	$c_2$

# First-last property

$\$1$	$a_2$	$c_1$	$a_1$	$a_3$	$c_2$	$g_1$
$a_1$	$a_3$	$c_2$	$g_1$	$\$1$	$a_2$	$c_1$
$a_2$	$c_1$	$a_1$	$a_3$	$c_2$	$g_1$	$\$1$
$a_3$	$c_2$	$g_1$	$\$1$	$a_2$	$c_1$	$a_1$
$c_1$	$a_1$	$a_3$	$c_2$	$g_1$	$\$1$	$a_2$
$c_2$	$g_1$	$\$1$	$a_2$	$c_1$	$a_1$	$a_3$
$g_1$	$\$1$	$a_2$	$c_1$	$a_1$	$a_3$	$c_2$

# Suffix array

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

# Suffix array

0. a c a a c g \$

1. c a a c g \$

2. a a c g \$

3. a c g \$

4. c g \$

5. g \$

6. \$

# 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

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



# Suffix array and BWT

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

# Suffix array and BWT

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

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

# BLASR

$k = 4$

AGGCAGGGGGCAGGTCTGCCACCGACCTCT

CCAGGCTGGGGGAGGTGGCACCGTTTCCTCTCT

# 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 \$

$R('') = (0, 6)$

# 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
$\$_1$	a	c	a	a	c	$g_1$	0	0
$a_1$	a	c	g	$\$$	a	$c_1$	0	0
$a_2$	c	a	a	c	g	$\$_1$	0	0
$a_3$	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	$c_2$		

# Suffix array interval

								a
$\$_1$	a	c	a	a	c	$g_1$	0	
$a_1$	a	c	g	$\$$	a	$c_1$	0	
$a_2$	c	a	a	c	g	$\$_1$	0	
$a_3$	c	g	$\$$	a	c	$a_1$	1	
$c_1$	a	a	c	g	$\$$	$a_2$		
$c_2$	g	$\$$	a	c	a	$a_3$		
$g_1$	$\$$	a	c	a	a	$c_2$		

# Suffix array interval

	a							
$\$_1$	a	c	a	a	c	$g_1$	0	
$a_1$	a	c	g	$\$$	a	$c_1$	0	
$a_2$	c	a	a	c	g	$\$_1$	0	
$a_3$	c	g	$\$$	a	c	$a_1$	1	
$c_1$	a	a	c	g	$\$$	$a_2$	2	
$c_2$	g	$\$$	a	c	a	$a_3$		
$g_1$	$\$$	a	c	a	a	$c_2$		

# Suffix array interval

	a							
$\$ _1$	a	c	a	a	c	$g _1$	0	
$a _1$	a	c	g	$\$$	a	$c _1$	0	
$a _2$	c	a	a	c	g	$\$ _1$	0	
$a _3$	c	g	$\$$	a	c	$a _1$	1	
$c _1$	a	a	c	g	$\$$	$a _2$	2	
$c _2$	g	$\$$	a	c	a	$a _3$	3	
$g _1$	$\$$	a	c	a	a	$c _2$		

# Suffix array interval

	a							
$\$$ <sub>1</sub>	a	c	a	a	c	$g$ <sub>1</sub>	0	
$a$ <sub>1</sub>	a	c	g	$\$$	a	$c$ <sub>1</sub>	0	
$a$ <sub>2</sub>	c	a	a	c	g	$\$$ <sub>1</sub>	0	
$a$ <sub>3</sub>	c	g	$\$$	a	c	$a$ <sub>1</sub>	1	
$c$ <sub>1</sub>	a	a	c	g	$\$$	$a$ <sub>2</sub>	2	
$c$ <sub>2</sub>	g	$\$$	a	c	a	$a$ <sub>3</sub>	3	
$g$ <sub>1</sub>	$\$$	a	c	a	a	$c$ <sub>2</sub>	3	

# Suffix array interval

								a	
$\$_1$	a	c	a	a	c	$g_1$	0		
$a_1$	a	c	g	$\$$	a	$c_1$	0	$O(a, 1) = 0$	
$a_2$	c	a	a	c	g	$\$_1$	0		
$a_3$	c	g	$\$$	a	c	$a_1$	1		
$c_1$	a	a	c	g	$\$$	$a_2$	2	$O(a, 4) = 2$	
$c_2$	g	$\$$	a	c	a	$a_3$	3		
$g_1$	$\$$	a	c	a	a	$c_2$	3		



# Suffix array interval

									a	
\$ <sub>1</sub>	a	c	a	a	c	g <sub>1</sub>	0			
a <sub>1</sub>	a	c	g	\$	a	c <sub>1</sub>	0	O(a, 1) = 0		
a <sub>2</sub>	c	a	a	c	g	\$ <sub>1</sub>	0	[a <sub>1</sub> , a <sub>2</sub> ]		
a <sub>3</sub>	c	g	\$	a	c	a <sub>1</sub>	1			
c <sub>1</sub>	a	a	c	g	\$	a <sub>2</sub>	2	O(a, 4) = 2		
c <sub>2</sub>	g	\$	a	c	a	a <sub>3</sub>	3			
g <sub>1</sub>	\$	a	c	a	a	c <sub>2</sub>	3			

# Suffix array interval

		a	c	g	t
$\$$ <sub>1</sub>	a c a a c <b>g</b> <sub>1</sub>	0	0	1	0
a <sub>1</sub>	a c g \$ a <b>c</b> <sub>1</sub>	0	1	1	0
a <sub>2</sub>	c a a c g <b>\$</b> <sub>1</sub>	0	1	1	0
a <sub>3</sub>	c g \$ a c <b>a</b> <sub>1</sub>	1	1	1	0
c <sub>1</sub>	a a c g \$ <b>a</b> <sub>2</sub>	2	1	1	0
c <sub>2</sub>	g \$ a c a <b>a</b> <sub>3</sub>	3	1	1	0
g <sub>1</sub>	\$ a c a a <b>c</b> <sub>2</sub>	3	2	1	0

# Suffix array interval

	a	c	g	t	$O(c, -1) = 0$
$\$$ <sub>1</sub> a c a a c <b>g</b> <sub>1</sub>	0	0	1	0	
<b>a</b> <sub>1</sub> a c g \$ a <b>c</b> <sub>1</sub>	0	1	1	0	<b>[c<sub>1</sub>]</b>
<b>a</b> <sub>2</sub> c a a c g <b>\$</b> <sub>1</sub>	0	1	1	0	
<b>a</b> <sub>3</sub> c g \$ a c <b>a</b> <sub>1</sub>	1	1	1	0	
<b>c</b> <sub>1</sub> a a c g \$ <b>a</b> <sub>2</sub>	2	1	1	0	$O(c, 4) = 1$
<b>c</b> <sub>2</sub> g \$ a c a <b>a</b> <sub>3</sub>	3	1	1	0	
<b>g</b> <sub>1</sub> \$ a c a a <b>c</b> <sub>2</sub>	3	2	1	0	

$$O(c, 4) = 1$$

# First-last property

$\$1$	a	c	a	a	c	$g_1$
$a_1$	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	$c_2$

# Suffix array interval

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

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

$$R_L("") = 0$$

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

# Intuition behind pattern search

**a**

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

# Intuition behind pattern search

**a**

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

# Intuition behind pattern search

**c a**

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



# Intuition behind pattern search

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

$\$ _1$  a c a a c **g<sub>1</sub>**

**a<sub>1</sub>** a c g \$ a **c<sub>1</sub>**

**a<sub>2</sub>** c a a c g **\$<sub>1</sub>**

**a<sub>3</sub>** c g \$ a c **a<sub>1</sub>**

**c<sub>1</sub>** a a c g \$ **a<sub>2</sub>**

**c<sub>2</sub>** g \$ a c a **a<sub>3</sub>**

**g<sub>1</sub>** \$ a c a a **c<sub>2</sub>**

# Intuition behind pattern search

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

$\$_1$  a c a a c  **$g_1$**

**$a_1$**  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  **$c_2$**

# Intuition behind pattern search

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

<b>\$<sub>1</sub></b>	a	c	a	a	c	<b>g<sub>1</sub></b>
<b>a<sub>1</sub></b>	a	c	g	\$	a	<b>c<sub>1</sub></b>
<b>a<sub>2</sub></b>	c	a	a	c	g	<b>\$<sub>1</sub></b>
<b>a<sub>3</sub></b>	c	g	\$	a	c	<b>a<sub>1</sub></b>
<b>c<sub>1</sub></b>	a	a	c	g	\$	<b>a<sub>2</sub></b>
<b>c<sub>2</sub></b>	g	\$	a	c	a	<b>a<sub>3</sub></b>
<b>g<sub>1</sub></b>	\$	a	c	a	a	<b>c<sub>2</sub></b>

# Intuition behind pattern search

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

$\$ _1$  a c a a c **g<sub>1</sub>**

**a<sub>1</sub>** a c g \$ a **c<sub>1</sub>**

**a<sub>2</sub>** c a a c g **\$<sub>1</sub>**

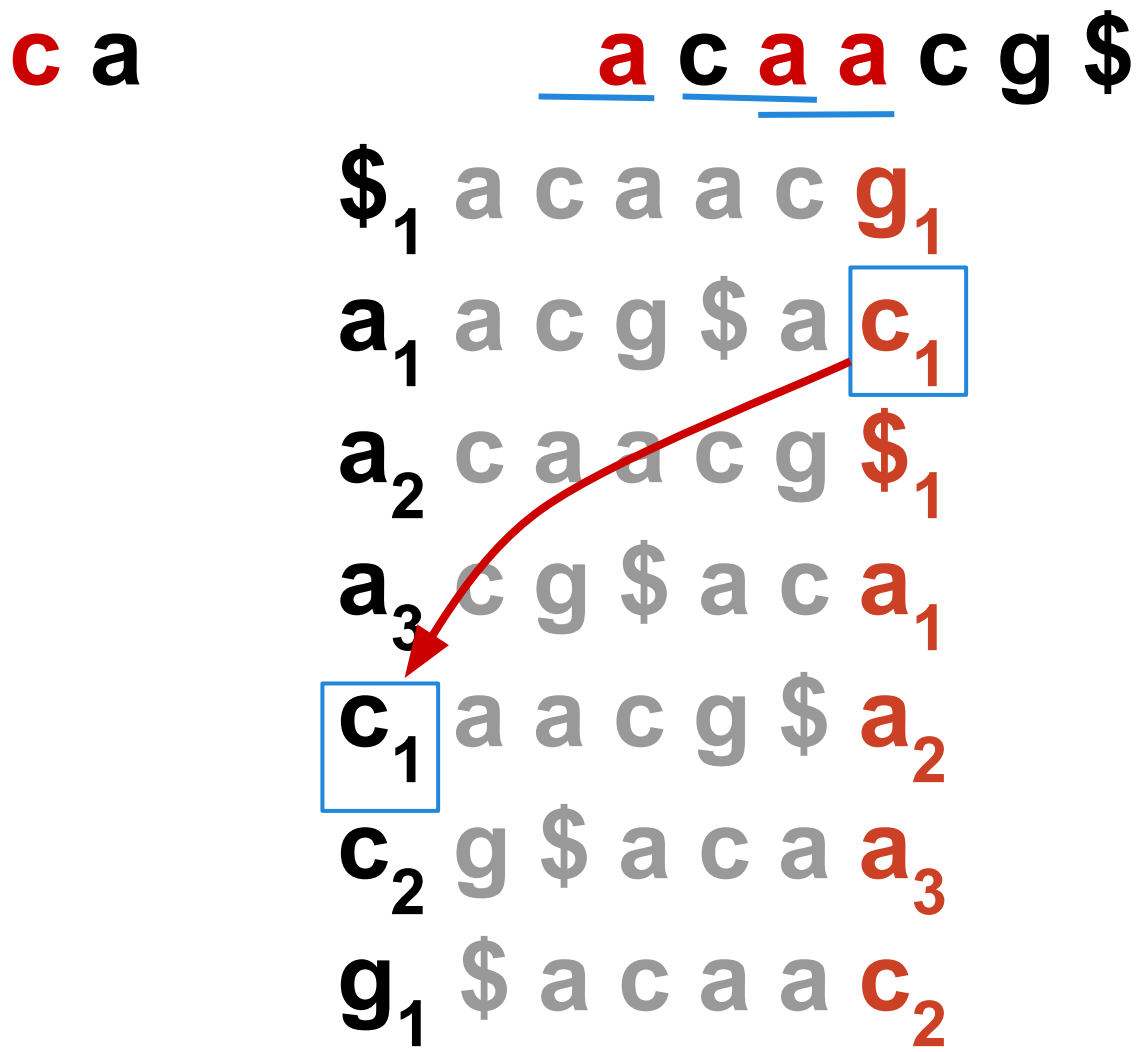
**a<sub>3</sub>** c g \$ a c **a<sub>1</sub>**

**c<sub>1</sub>** a a c g \$ **a<sub>2</sub>**

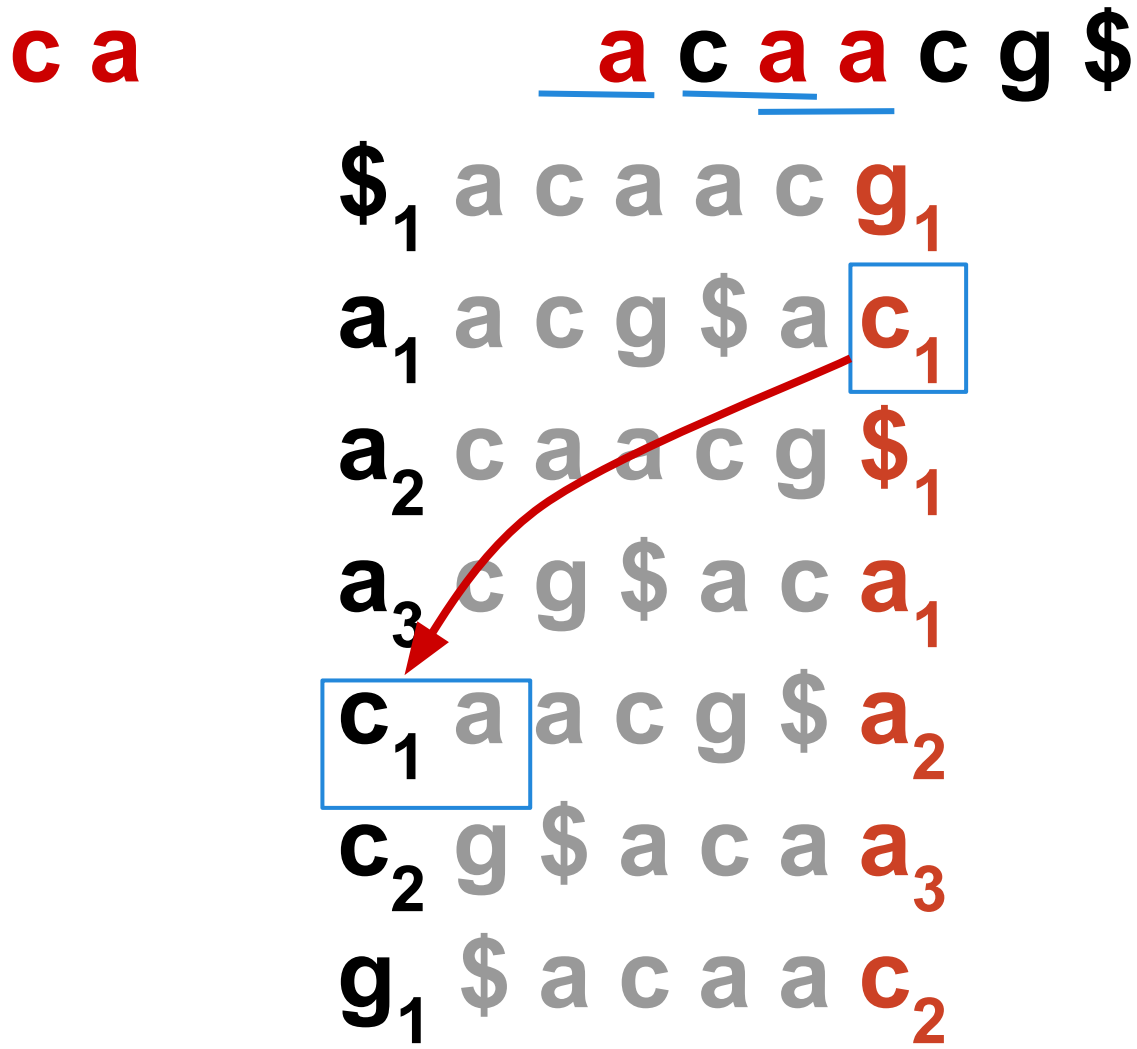
**c<sub>2</sub>** g \$ a c a **a<sub>3</sub>**

**g<sub>1</sub>** \$ a c a a **c<sub>2</sub>**

# Intuition behind pattern search



# Intuition behind pattern search



# 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 a c**

									a	c	g	t
$C(\$) = 0$	→	<b>\$</b> <sub>1</sub>	a	c	a	a	c	<b>g</b> <sub>1</sub>	0	0	1	0
$C(a) = 1$	→	<b>a</b> <sub>1</sub>	a	c	g	\$	a	<b>c</b> <sub>1</sub>	0	1	1	0
		<b>a</b> <sub>2</sub>	c	a	a	c	g	<b>\$</b> <sub>1</sub>	0	1	1	0
$C(c) = 4$	→	<b>a</b> <sub>3</sub>	c	g	\$	a	c	<b>a</b> <sub>1</sub>	1	1	1	0
		<b>c</b> <sub>1</sub>	a	a	c	g	\$	<b>a</b> <sub>2</sub>	2	1	1	0
$C(g) = 6$	→	<b>c</b> <sub>2</sub>	g	\$	a	c	a	<b>a</b> <sub>3</sub>	3	1	1	0
$C(t) = 7$	→	<b>g</b> <sub>1</sub>	\$	a	c	a	a	<b>c</b> <sub>2</sub>	3	2	1	0

# Pattern search

**a a c**

$$R_L("") = 0$$

$$R_H("") = 6$$

									a	c	g	t
$C(\$) = 0$	→	<b>\$</b> <sub>1</sub>	a	c	a	a	c	<b>g</b> <sub>1</sub>	0	0	1	0
$C(a) = 1$	→	<b>a</b> <sub>1</sub>	a	c	g	\$	a	<b>c</b> <sub>1</sub>	0	1	1	0
		<b>a</b> <sub>2</sub>	c	a	a	c	g	<b>\$</b> <sub>1</sub>	0	1	1	0
$C(c) = 4$	→	<b>a</b> <sub>3</sub>	c	g	\$	a	c	<b>a</b> <sub>1</sub>	1	1	1	0
		<b>c</b> <sub>1</sub>	a	a	c	g	\$	<b>a</b> <sub>2</sub>	2	1	1	0
$C(g) = 6$	→	<b>c</b> <sub>2</sub>	g	\$	a	c	a	<b>a</b> <sub>3</sub>	3	1	1	0
$C(t) = 7$	→	<b>g</b> <sub>1</sub>	\$	a	c	a	a	<b>c</b> <sub>2</sub>	3	2	1	0

# Pattern search

**a a c**

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

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

									a	c	g	t
$C(\$) = 0$	→	<b>\$</b> <sub>1</sub>	a	c	a	a	c	<b>g</b> <sub>1</sub>	0	0	1	0
$C(a) = 1$	→	<b>a</b> <sub>1</sub>	a	c	g	\$	a	<b>c</b> <sub>1</sub>	0	1	1	0
		<b>a</b> <sub>2</sub>	c	a	a	c	g	<b>\$</b> <sub>1</sub>	0	1	1	0
$C(c) = 4$	→	<b>a</b> <sub>3</sub>	c	g	\$	a	c	<b>a</b> <sub>1</sub>	1	1	1	0
		<b>c</b> <sub>1</sub>	a	a	c	g	\$	<b>a</b> <sub>2</sub>	2	1	1	0
$C(g) = 6$	→	<b>c</b> <sub>2</sub>	g	\$	a	c	a	<b>a</b> <sub>3</sub>	3	1	1	0
$C(t) = 7$	→	<b>g</b> <sub>1</sub>	\$	a	c	a	a	<b>c</b> <sub>2</sub>	3	2	1	0

# 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
$C(\$) = 0$	→	<b>\$</b> <sub>1</sub>	a	c	a	a	c	<b>g</b> <sub>1</sub>	0	0	1	0
$C(a) = 1$	→	<b>a</b> <sub>1</sub>	a	c	g	\$	a	<b>c</b> <sub>1</sub>	0	1	1	0
		<b>a</b> <sub>2</sub>	c	a	a	c	g	<b>\$</b> <sub>1</sub>	0	1	1	0
$C(c) = 4$	→	<b>a</b> <sub>3</sub>	c	g	\$	a	c	<b>a</b> <sub>1</sub>	1	1	1	0
		<b>c</b> <sub>1</sub>	a	a	c	g	\$	<b>a</b> <sub>2</sub>	2	1	1	0
$C(g) = 6$	→	<b>c</b> <sub>2</sub>	g	\$	a	c	a	<b>a</b> <sub>3</sub>	3	1	1	0
$C(t) = 7$	→	<b>g</b> <sub>1</sub>	\$	a	c	a	a	<b>c</b> <sub>2</sub>	3	2	1	0

# Pattern search

a a c

$$R_L(c) = C(c) + \mathbf{O}(c, -1)$$

$$R_H(c) = C(c) + \mathbf{O}(c, 6) - 1$$

[illegible]

# Pattern search

a a c

$$R_L(c) = \mathbf{C}(c) + 0$$

$$R_H(c) = \mathbf{C}(c) + 2 - 1$$

[illegible]

# Pattern search

**a a c**

$$R_L(c) = 4 + 0 = 4$$

$$R_H(c) = 4 + 2 - 1 = 5$$

									<b>a c g t</b>
$C(\$) = 0$	→	<b>\$</b> <sub>1</sub>	a	c	a	a	c	<b>g</b> <sub>1</sub>	<b>0 0 1 0</b>
$C(a) = 1$	→	<b>a</b> <sub>1</sub>	a	c	g	\$	a	<b>c</b> <sub>1</sub>	<b>0 1 1 0</b>
		<b>a</b> <sub>2</sub>	c	a	a	c	g	<b>\$</b> <sub>1</sub>	<b>0 1 1 0</b>
		<b>a</b> <sub>3</sub>	c	g	\$	a	c	<b>a</b> <sub>1</sub>	<b>1 1 1 0</b>
$C(c) = 4$	→	<b>c</b> <sub>1</sub>	a	a	c	g	\$	<b>a</b> <sub>2</sub>	<b>2 1 1 0</b>
		<b>c</b> <sub>2</sub>	g	\$	a	c	a	<b>a</b> <sub>3</sub>	<b>3 1 1 0</b>
$C(g) = 6$	→	<b>g</b> <sub>1</sub>	\$	a	c	a	a	<b>c</b> <sub>2</sub>	<b>3 2 1 0</b>
$C(t) = 7$	→								

# Pattern search

a a c

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

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

[illegible]



# Pattern search

a a c

$$R_l(ac) = C(a) + \mathbf{O}(a, 3)$$

$$R_H(ac) = C(a) + \mathbf{O(a, 5)} - 1$$

The diagram illustrates the construction of a suffix array for the string "acgt" using a bucketing approach. It shows the iterative insertion of characters into buckets, with the final suffix array being [1, 2, 3, 0].

**Initial State:** The buckets are labeled with the characters 'a', 'c', 'g', and 't'. The initial suffix array is empty.

**Step 1: Inserting '\$' (C(\$)=0)**

- Character: **\$**<sub>1</sub>
- Bucket: **a**
- Current Suffix Array: [0]

**Step 2: Inserting 'a' (C(a)=1)**

- Character: **a**<sub>1</sub>
- Bucket: **a**
- Current Suffix Array: [0, 1]

**Step 3: Inserting 'a' (C(a)=2)**

- Character: **a**<sub>2</sub>
- Bucket: **a**
- Current Suffix Array: [0, 1, 2]

**Step 4: Inserting 'c' (C(c)=4)**

- Character: **c**<sub>1</sub>
- Bucket: **c**
- Current Suffix Array: [0, 1, 2, 3]

**Step 5: Inserting 'c' (C(c)=5)**

- Character: **c**<sub>2</sub>
- Bucket: **c**
- Current Suffix Array: [0, 1, 2, 3, 4]

**Step 6: Inserting 'g' (C(g)=6)**

- Character: **g**<sub>1</sub>
- Bucket: **g**
- Current Suffix Array: [0, 1, 2, 3, 4, 5]

**Step 7: Inserting 't' (C(t)=7)**

- Character: **t**<sub>1</sub>
- Bucket: **t**
- Current Suffix Array: [0, 1, 2, 3, 4, 5, 6]

**Final Suffix Array:** [0, 1, 2, 3, 4, 5, 6]

# Pattern search

a a c

$$R_l(ac) = C(a) + \mathbf{O}(a, 3)$$

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

		a	c	g	t
$C(\$) = 0$	$\$$	0	0	1	0
$C(a) = 1$	$a$	0	1	1	0
	$a_2$	0	1	1	0
$C(c) = 4$	$a_3$	1	1	1	0
	$c_1$	2	1	1	0
$C(g) = 6$	$c_2$	3	1	1	0
$C(t) = 7$	$g_1$	3	2	1	0

# Pattern search

a a c

$$R_L(ac) = C(a) + 1$$

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

		a	c	g	t
$C(\$) = 0$	$\$$	0	0	1	0
$C(a) = 1$	$a$	0	1	1	0
	$a_2$	0	1	1	0
$C(c) = 4$	$a_3$	1	1	1	0
	$c_1$	2	1	1	0
$C(g) = 6$	$c_2$	3	1	1	0
$C(t) = 7$	$g_1$	3	2	1	0

# Pattern search

a a c

$$R_L(ac) = \mathbf{C(a)} + 1$$

$$R_H(ac) = \mathbf{C(a)} + 3 - 1$$

Diagram illustrating the construction of a suffix array for the string "acgt" using a bucketing approach. The suffixes are sorted iteratively based on their first character, then their first two characters, and so on.

**Initial Suffixes and Buckets:**

- $C(\$) = 0$ : Buckets for the first character (\$). Suffixes:  $\$_1$  (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).

**Intermediate Steps:**

- $C(a) = 1$ : Buckets for the first character (a). Suffixes:  $a_1$  (a c g \$ a c),  $a_2$  (c a a c g \$),  $a_3$  (c g \$ a c a),  $a_2$  (a a c g \$ a),  $a_3$  (g \$ a c a a),  $a_1$  (\$ a c a a c).
- $C(c) = 4$ : Buckets for the first character (c). Suffixes:  $c_1$  (a a c g \$ a),  $c_2$  (g \$ a c a a),  $c_1$  (a c g \$ a c),  $c_2$  (c a a c g \$),  $c_1$  (a g \$ a c a),  $c_2$  (\$ a c a a c).
- $C(g) = 6$ : Buckets for the first character (g). Suffixes:  $g_1$  (\$ a c a a c),  $g_1$  (a a c g \$ a),  $g_1$  (c a a c g \$),  $g_1$  (g \$ a c a a),  $g_1$  (a c g \$ a c),  $g_1$  (c g \$ a c a).
- $C(t) = 7$ : Buckets for the first character (t). Suffixes:  $t_1$  (a c g \$ a c),  $t_1$  (c a a c g \$),  $t_1$  (g \$ a c a a),  $t_1$  (a a c g \$ a),  $t_1$  (c g \$ a c a),  $t_1$  (\$ a c a a c).

**Final Sorted Order (Suffix Array):**

	a	c	g	t
0	0	0	1	0
1	0	1	1	0
2	0	1	1	0
3	1	1	1	0
4	2	1	1	0
5	3	1	1	0
6	3	2	1	0

# Pattern search

a a c

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

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

[illegible]

# Pattern search

a a c

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

**a c g t**

0 0 1 0

0 1 1 0

0 1 1 0

1 1 1 0

**2 1 1 0**

**3 1 1 0**

# 3 2 1 0

**\$<sub>1</sub>**

$$a_1$$
 $a_2$  $a_3$ 

**C<sub>1</sub>**

**C<sub>2</sub>**

**g<sub>1</sub>**

a c a a c g<sub>1</sub>

a c g \$ a c<sub>1</sub>

c a a c g \$<sub>1</sub>

c g \$ a c a<sub>1</sub>

a a c g \$ a<sub>2</sub>

g \$ a c a a<sub>3</sub>

\$ a c a a c<sub>2</sub>

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


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

# Pattern search

a a c

$$R_1(aac) = C(a) + O(a, 2 - 1)$$

$$R_H(aac) = C(a) + \mathbf{O(a, 3)} - 1$$

		a	c	g	t							
$C(\$) = 0$		<b>\$<sub>1</sub></b>	a	c	a	a	c	<b>g<sub>1</sub></b>	0	0	1	0
$C(a) = 1$		<b>a<sub>1</sub></b>	a	c	g	\$	a	<b>c<sub>1</sub></b>	0	1	1	0
		<b>a<sub>2</sub></b>	c	a	a	c	g	<b>\$<sub>1</sub></b>	0	1	1	0
$C(c) = 4$		<b>a<sub>3</sub></b>	c	g	\$	a	c	<b>a<sub>1</sub></b>	1	1	1	0
		<b>c<sub>1</sub></b>	a	a	c	g	\$	<b>a<sub>2</sub></b>	2	1	1	0
$C(g) = 6$		<b>c<sub>2</sub></b>	g	\$	a	c	a	<b>a<sub>3</sub></b>	3	1	1	0
$C(t) = 7$		<b>g<sub>1</sub></b>	\$	a	c	a	a	<b>c<sub>2</sub></b>	3	2	1	0

# Pattern search

a a c

$$R_1(\text{aac}) = \mathbf{C}(\mathbf{a}) + 0$$

$$R_H(aac) = \mathbf{C(a)} + 1 - 1$$

Diagram illustrating the construction of a suffix array for the string "acgt" using a bucketing approach. The suffixes are sorted iteratively based on their first character, then their first two characters, and so on.

**Initial Suffixes and Buckets:**

- $C(\$) = 0$ : Buckets for '\$' (rank 0) and 'a' (rank 1).
- $C(a) = 1$ : Buckets for 'a' (rank 1) and 'c' (rank 2).
- $C(c) = 4$ : Buckets for 'c' (rank 4) and 'g' (rank 6).
- $C(g) = 6$ : Buckets for 'g' (rank 6) and 't' (rank 7).
- $C(t) = 7$ : Buckets for 't' (rank 7).

**Final Suffix Array (SA):**

SA	0	1	2	3	4	5	6	7
Suffix	\$	a <sub>1</sub>	a <sub>2</sub>	a <sub>3</sub>	c <sub>1</sub>	c <sub>2</sub>	g <sub>1</sub>	t



# Pattern search

a a c

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

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

Diagram illustrating the construction of a suffix array for the string "acgt" using a suffix tree. The suffix tree has root 'a' and children 'c', 'g', 't'. The suffixes are listed in a table with their corresponding suffix array indices. The suffixes are:  $\$_1 a c a a c g_1$ ,  $a_1 a c g \$ a c_1$ ,  $a_2 c a a c g \$_1$ ,  $a_3 c g \$ a c a_1$ ,  $c_1 a a c g \$ a_2$ ,  $c_2 g \$ a c a a_3$ , and  $g_1 \$ a c a a c_2$ . The suffix array indices are: 0 0 1 0, 0 1 1 0, 0 1 1 0, 1 1 1 0, 2 1 1 0, 3 1 1 0, and 3 2 1 0. The suffixes are sorted lexicographically, and the suffix array indices are the corresponding positions in the original string.

# Pattern search

**a a c**

$$R_L(aac) = 1$$

$$R_H(aac) = 1$$



									<b>a c g t</b>
$C(\$) = 0$	→								
$C(a) = 1$	→	<b>\$<sub>1</sub></b>	a	c	a	a	c	<b>g<sub>1</sub></b>	<b>0 0 1 0</b>
	→	<b>a<sub>1</sub></b>	<b>a</b>	<b>c</b>	g	\$	a	<b>c<sub>1</sub></b>	<b>0 1 1 0</b>
		<b>a<sub>2</sub></b>	c	a	a	c	g	<b>\$<sub>1</sub></b>	<b>0 1 1 0</b>
		<b>a<sub>3</sub></b>	c	g	\$	a	c	<b>a<sub>1</sub></b>	<b>1 1 1 0</b>
$C(c) = 4$	→	<b>c<sub>1</sub></b>	a	a	c	g	\$	<b>a<sub>2</sub></b>	<b>2 1 1 0</b>
		<b>c<sub>2</sub></b>	g	\$	a	c	a	<b>a<sub>3</sub></b>	<b>3 1 1 0</b>
$C(g) = 6$	→	<b>g<sub>1</sub></b>	\$	a	c	a	a	<b>c<sub>2</sub></b>	<b>3 2 1 0</b>
$C(t) = 7$	→								

# 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  $g_1$

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	$g_1$	$g$	$c$	$a$
$a_1$	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	$g_1$	$g$	$c$	$a$
$a_1$	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

$\$$ <sub>1</sub>	a	c	a	a	c	$g$ <sub>1</sub>	$g$	$c$	$a$
$a$ <sub>1</sub>	a	c	g	$\$$	a	$c$ <sub>1</sub>			
$a$ <sub>2</sub>	c	a	a	c	g	$\$$ <sub>1</sub>			
$a$ <sub>3</sub>	c	g	$\$$	a	c	$a$ <sub>1</sub>			
$c$ <sub>1</sub>	a	a	c	g	$\$$	$a$ <sub>2</sub>			
$c$ <sub>2</sub>	g	$\$$	a	c	a	$a$ <sub>3</sub>			
$g$ <sub>1</sub>	$\$$	a	c	a	a	$c$ <sub>2</sub>			



# Pattern search with errors

$\$_1$	a	c	a	a	c	$g_1$	
$a_1$	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

$\$$ <sub>1</sub> a c a a c **g**<sub>1</sub>

**a**<sub>1</sub> a c g \$ a **c**<sub>1</sub>

**a**<sub>2</sub> c a a c g **\$**<sub>1</sub>

**a**<sub>3</sub> c g \$ a c **a**<sub>1</sub>

**c**<sub>1</sub> a a c g \$ **a**<sub>2</sub>

**c**<sub>2</sub> g \$ a c a **a**<sub>3</sub>

**g**<sub>1</sub> \$ a c a a **c**<sub>2</sub>

**g** **c** a

# Pattern search with errors

$\$$ <sub>1</sub>	a	c	a	a	c	g <sub>1</sub>
a <sub>1</sub>	a	c	g	\$	a	c <sub>1</sub>
a <sub>2</sub>	c	a	a	c	g	\$ <sub>1</sub>
a <sub>3</sub>	c	g	\$	a	c	a <sub>1</sub>
c <sub>1</sub>	a	a	c	g	\$	a <sub>2</sub>
c <sub>2</sub>	g	\$	a	c	a	a <sub>3</sub>
g <sub>1</sub>	\$	a	c	a	a	c <sub>2</sub>

g c a

# Pattern search with errors

$\$$ <sub>1</sub>	a	c	a	a	c	g	<sub>1</sub>	
a	<sub>1</sub>	a	c	g	$\$$	a	c	<sub>1</sub>
a	<sub>2</sub>	c	a	a	c	g	$\$$	<sub>1</sub>
a	<sub>3</sub>	c	g	$\$$	a	c	a	<sub>1</sub>
c	<sub>1</sub>	a	a	c	g	$\$$	a	<sub>2</sub>
c	<sub>2</sub>	g	$\$$	a	c	a	a	<sub>3</sub>
g	<sub>1</sub>	$\$$	a	c	a	a	c	<sub>2</sub>

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  $g_1$

$a_1$  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

$\$$ <sub>1</sub> a c a a c **g**<sub>1</sub>

**a**<sub>1</sub> a c g \$ a **c**<sub>1</sub>

**a**<sub>2</sub> c a a c g **\$**<sub>1</sub>

**a**<sub>3</sub> c g \$ a c **a**<sub>1</sub>

**c**<sub>1</sub> a a c g \$ **a**<sub>2</sub>

**c**<sub>2</sub> g \$ a c a **a**<sub>3</sub>

**g**<sub>1</sub> \$ a c a a **c**<sub>2</sub>

**g** c a

# Pattern search with errors

$\$$ <sub>1</sub>	a	c	a	a	c	g	$g_1$
a <sub>1</sub>	a	c	g	\$	a	c	$c_1$
a <sub>2</sub>	c	a	a	c	g	\$	$\$$ <sub>1</sub>
a <sub>3</sub>	c	g	\$	a	c	a	$a_1$
c <sub>1</sub>	a	a	c	g	\$	a	$a_2$
c <sub>2</sub>	g	\$	a	c	a	a	$a_3$
g <sub>1</sub>	\$	a	c	a	a	c	$c_2$

g c a





# Pattern search with errors

$\$_1$  a c a a c  $g_1$

$a_1$  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  $g _1$

$a _1$  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  $c _2$

c c a

# Pattern search with errors

$\$_1$	a	c	a	a	c	$g_1$
$a_1$	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><math>a_2</math></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	$g_1$
$a_1$	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

# Pattern search with errors

$\$_1$	a	c	a	a	c	$g_1$
$a_1$	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

# Pattern search with errors

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

$\$$ <sub>1</sub>	a	c	a	a	c	g	$g$ <sub>1</sub>
$a$ <sub>1</sub>	a	c	g	$\$$	a	$c$ <sub>1</sub>	
$a$ <sub>2</sub>	c	a	a	c	g	$\$$ <sub>1</sub>	
$a$ <sub>3</sub>	c	g	$\$$	a	c	$a$ <sub>1</sub>	
$c$ <sub>1</sub>	a	a	c	g	$\$$	$a$ <sub>2</sub>	
$c$ <sub>2</sub>	g	$\$$	a	c	a	$a$ <sub>3</sub>	
$g$ <sub>1</sub>	$\$$	a	c	a	a	$c$ <sub>2</sub>	

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

A red arrow points from the first 'a' of the second string 'c a a' to the first 'a' of the first string 'c a a c a'.

# 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



$\leq 2$

$\geq 0$

# Seeds



$\leq 2$

$\geq 0$

# Seeds

1.

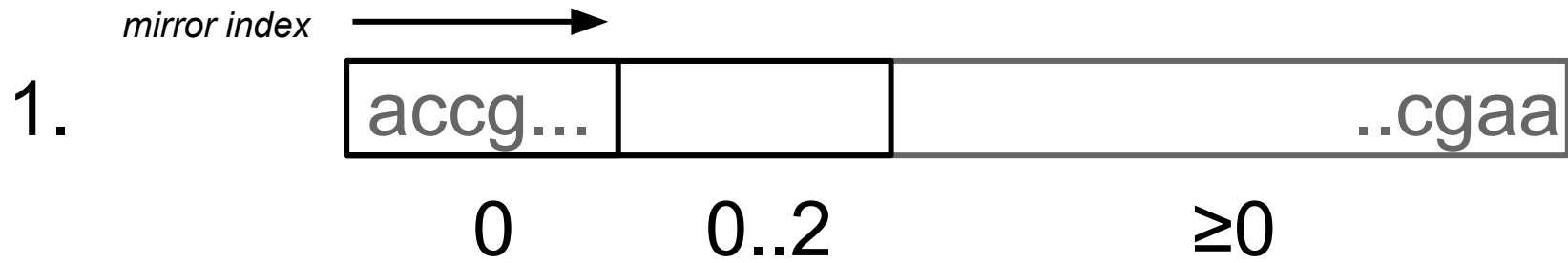


0

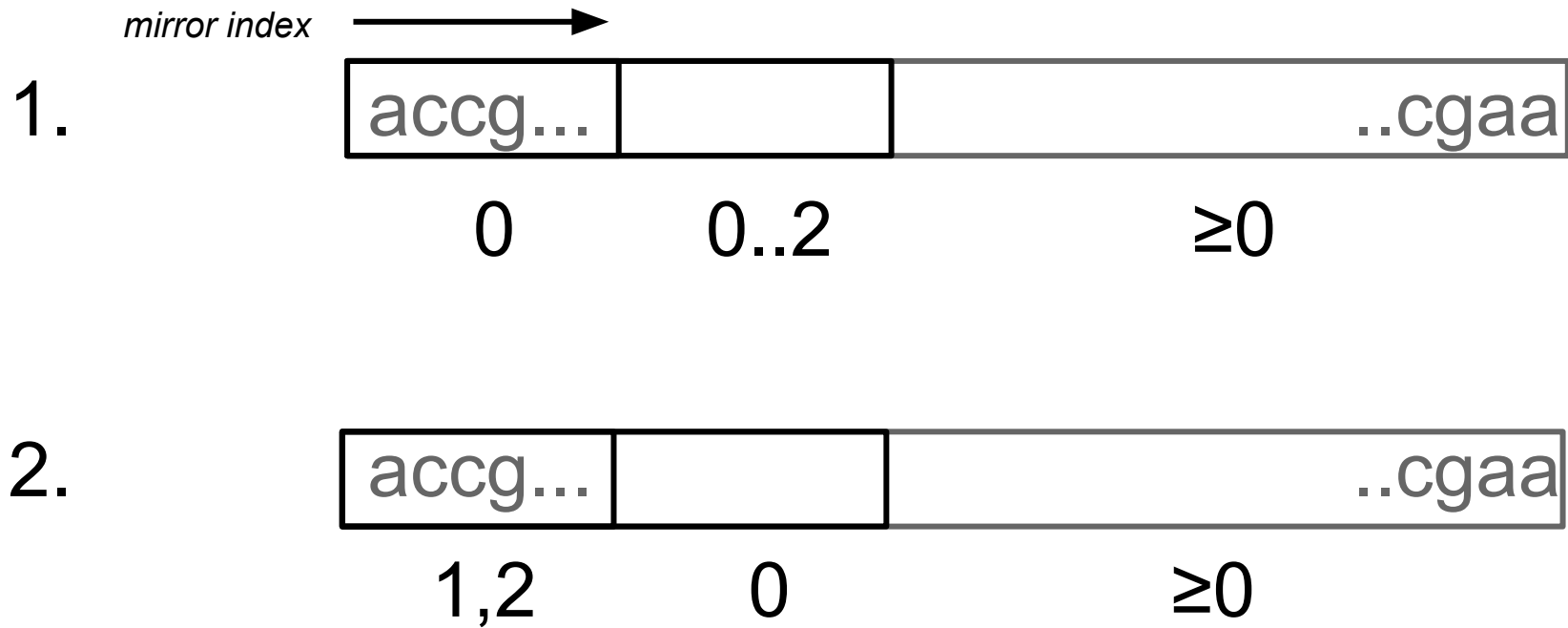
0..2

$\geq 0$

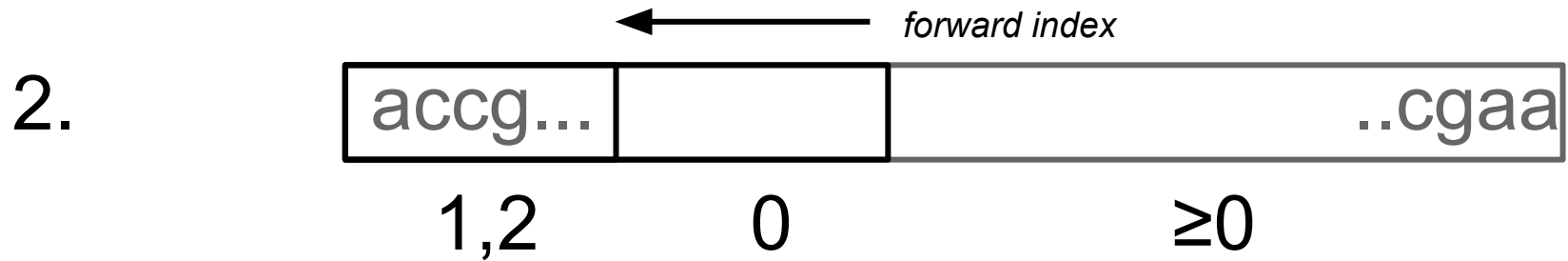
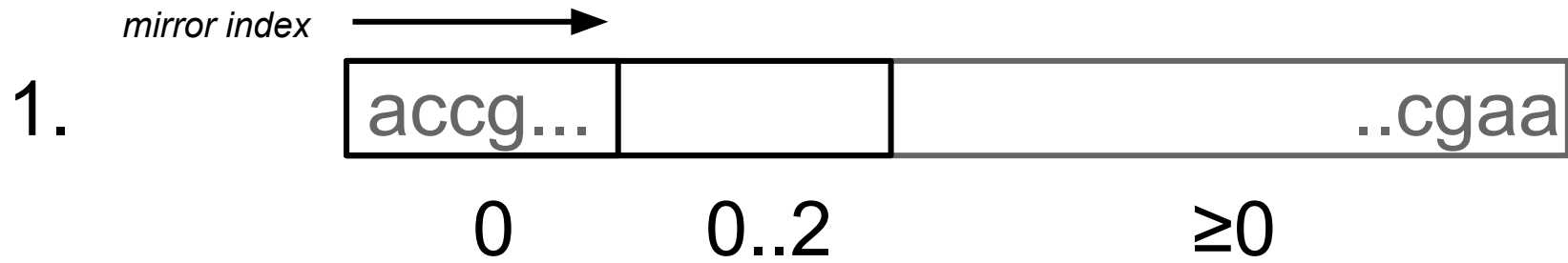
# Seeds



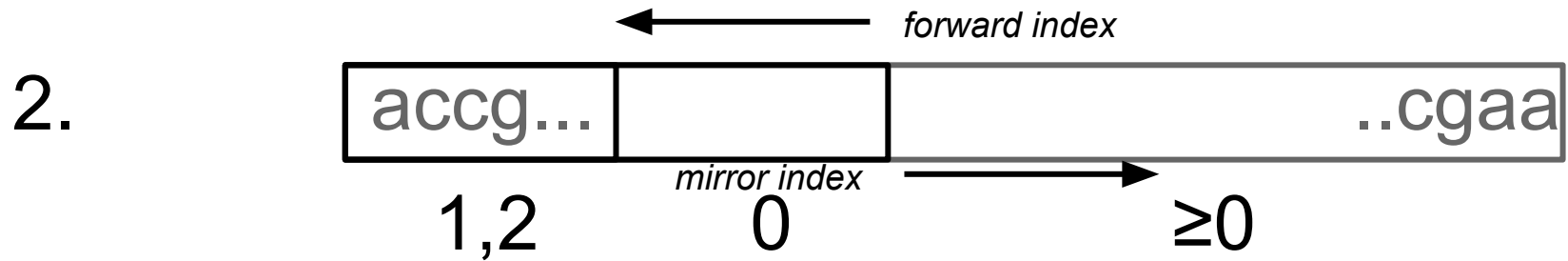
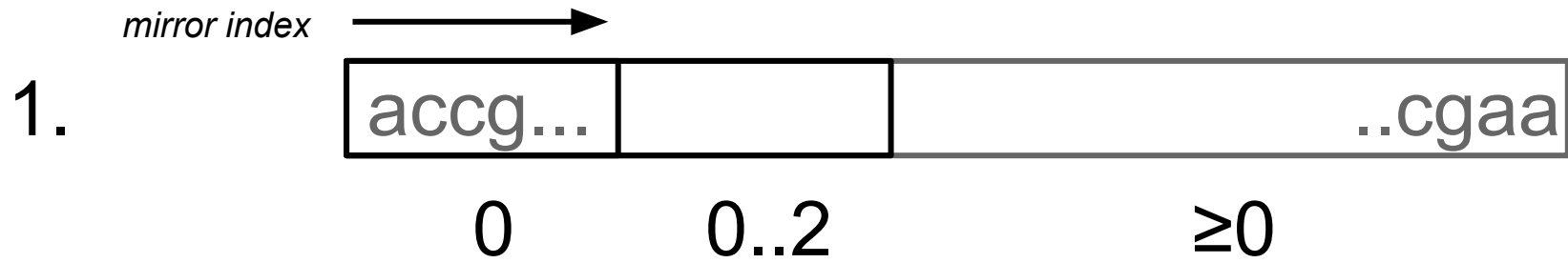
# Seeds



# Seeds

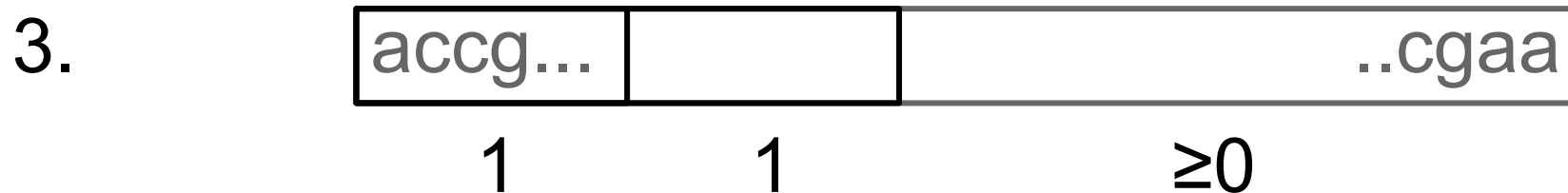
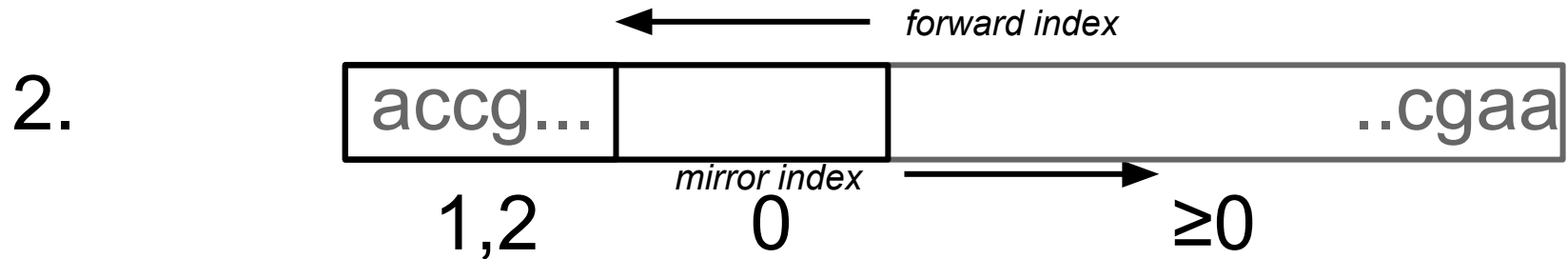
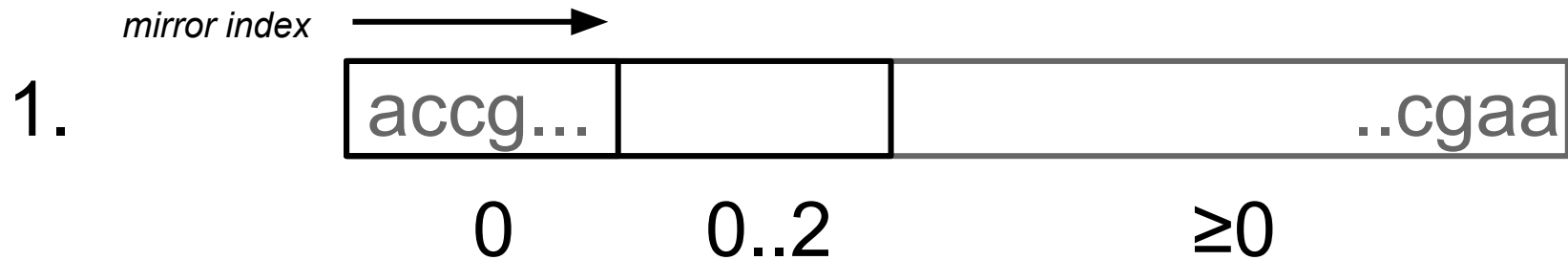


# Seeds

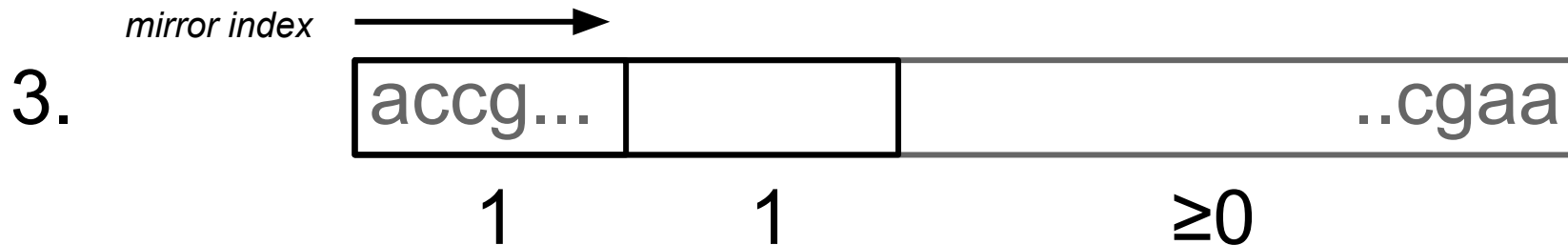
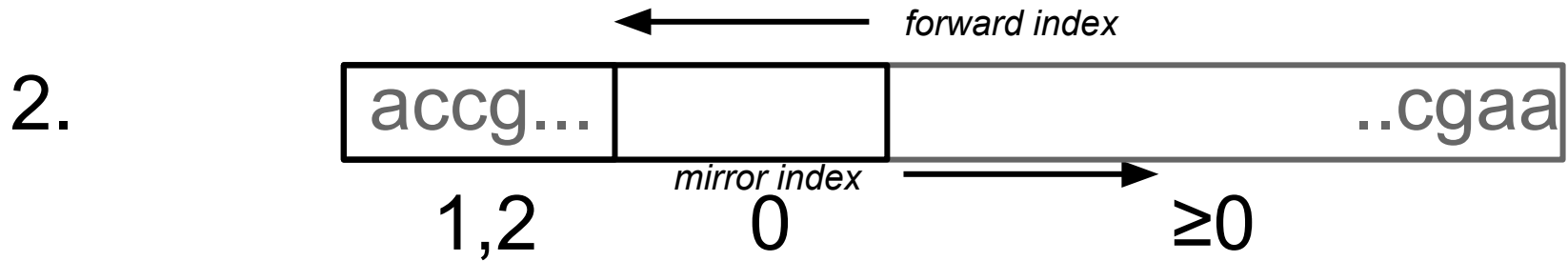
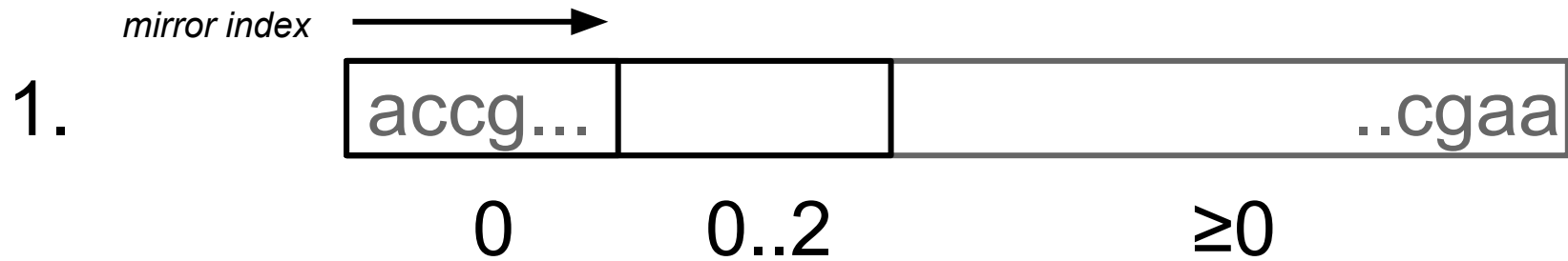




# Seeds



# Seeds



# Working with alignments

# Bowtie steps

- Build index from the genome

**bowtie2-build** <genome.fasta> <index>

- Map reads using index

**bowtie2** -x <index>

-1 <left.fastq> -2 <right.fastq>

-S <output.sam>

# BWA tool

- **bwa index** — index construction
- **bwa aln** — short 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



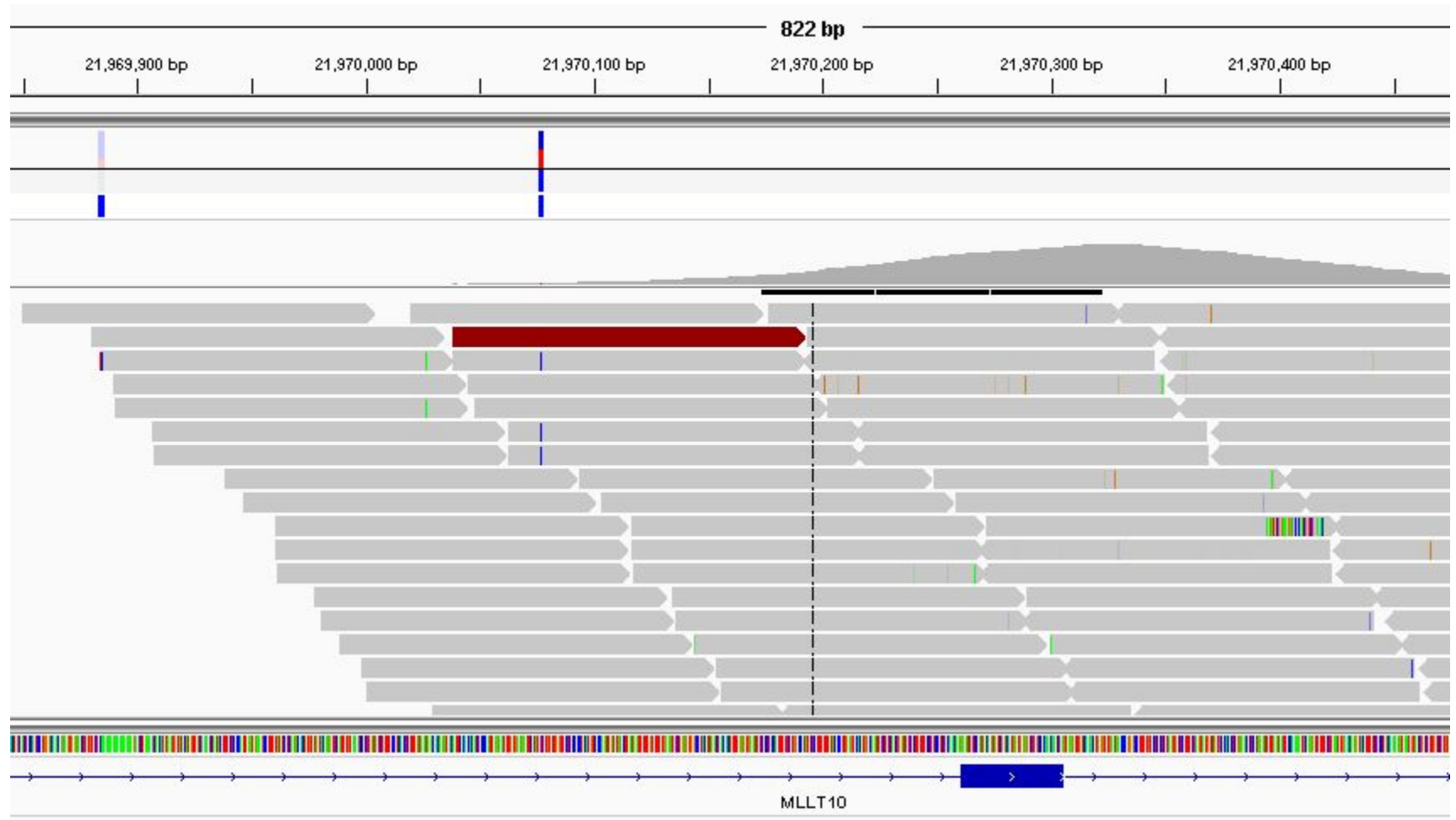
# SAM/BAM files

- 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
      CCGTGTTTAAAGGTGGATGCGGTCACCTTCCCAGCTAGGCTTAGGGATTCTTAGTTGGCCTAGGAAATCCAGCTAGTCCTGTCTCTCAGTCCCCCTCT
C      BBDCCDDCCDDDDCDDDDDDCDDCCDBC?DDDDDDDDDDDDDDDDCCDDDDDDDDDDCCCCEDDDC?DDDDDDDDDDDDDDDDDDDDDDBDHFFFFDC@@
      AS:i:-15      XM:i:3      XO: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
      TGCTGGATCATCTGGTTAGTGGCTTCTGACTCAGAGGACCTTCGTCCCCTGGGGCAGTGGACCTTCCAGTGATTCCCCTGACATAAGGGGCATGGACGA
G      DCDDDDDEDDDDDDCDDDDDDCDDDDDEEC>DFFFEJJJJIGJJJJIHGBHHGJIJJJJJGJJJIJJJJJIHJJJJJJHHHHHFFFFFCCC
      AS:i:-16      XM:i:3      XO: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
      GGCTTTATTGGTAAAAAAGGAATAGCAGATTTAATCAGAAATTCCCACCTGGCCCAGCAGCACCAACCAGAAAGAAGGGAAGAAGACAGGAAAAAACCA
C      DDDDDDDDDCDDDDDDDDDDDEEEEEEEFFFEFFEGHHHHFGDJJIHJJJIJJJJIIIGGFJJIIHIIIIJJJJJJIGHHFAHGFHJHFGGHFFFDDBB
      AS:i:-11      XM:i:2      XO: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      GTGGCTCTTCCACAGGAATGTTGAGGATGACATCCATGTCTGGGGTGCACTTGGGTCTCCGAAGCAGAACATCCTCAAATATGACCTCTCG
accepted_hits.sam
```

# Alignment visualization with IGV



# Alignment visualization (Tablet)





# Samtools: working with SAM/BAM

- samtools **flagstat** <input.bam>
- samtools **view** <input.bam> | less -S
- samtools **sort** <input.bam> -o <out.bam>
- samtools **index** <input.bam>
- samtools **view** <input.bam> “chr1:100-500”

<http://www.htslib.org/doc/samtools.html>

# Pysam: SAM/BAM in Python

```
import pysam
samfile = pysam.AlignmentFile("ex1.bam", "rb")

for read in samfile.fetch('chr1', 10000, 20000):
    print('%s %s %d' %
          (read.query_name,
           read.reference_name,
           read.reference_start))
```

<https://pysam.readthedocs.io/en/latest/api.html>

# SNP calling

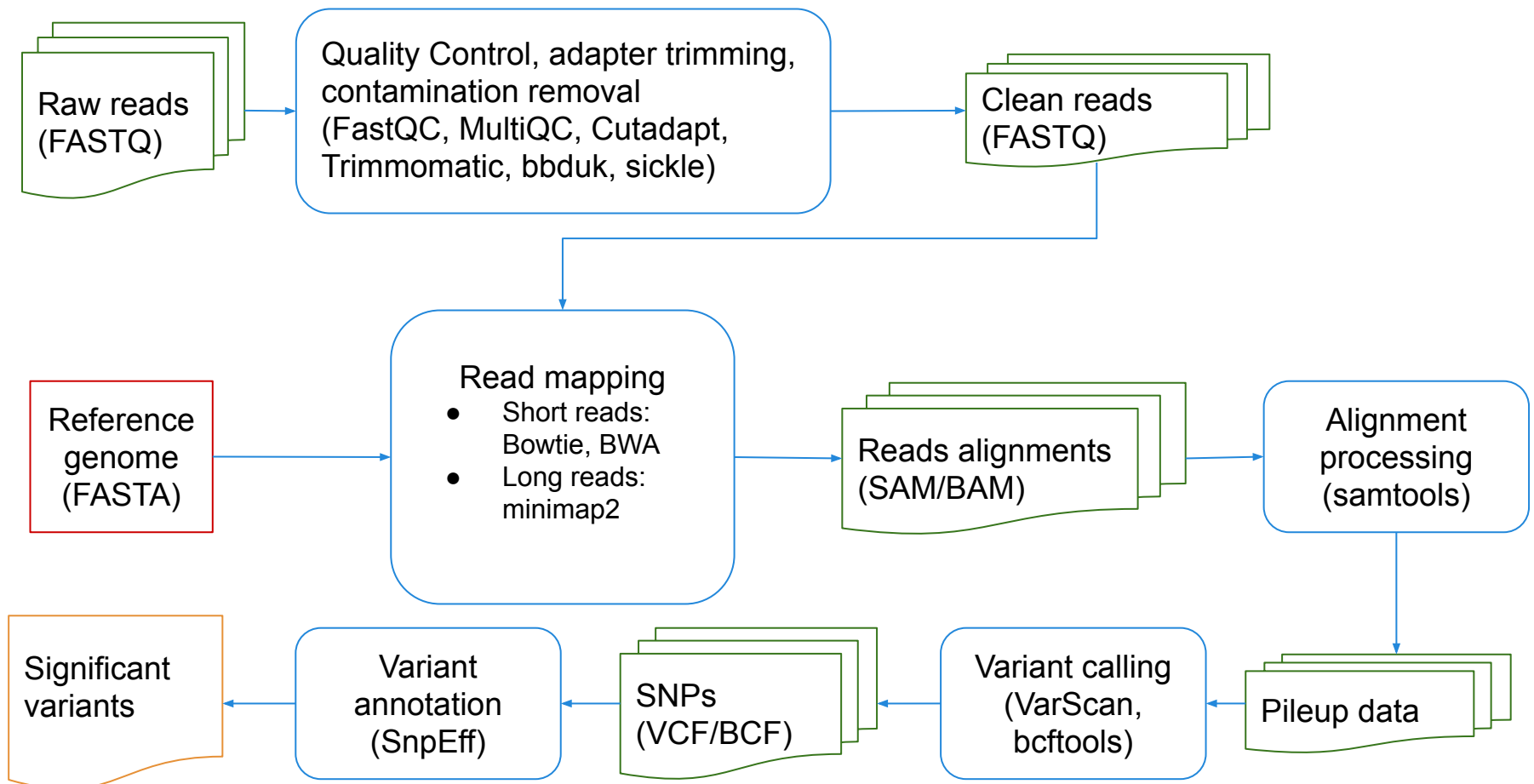
Process of detecting positions that differ from the reference genome

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

We need enough coverage to distinguish real variants from sequencing errors

# Reference based analysis

## SNP calling pipeline





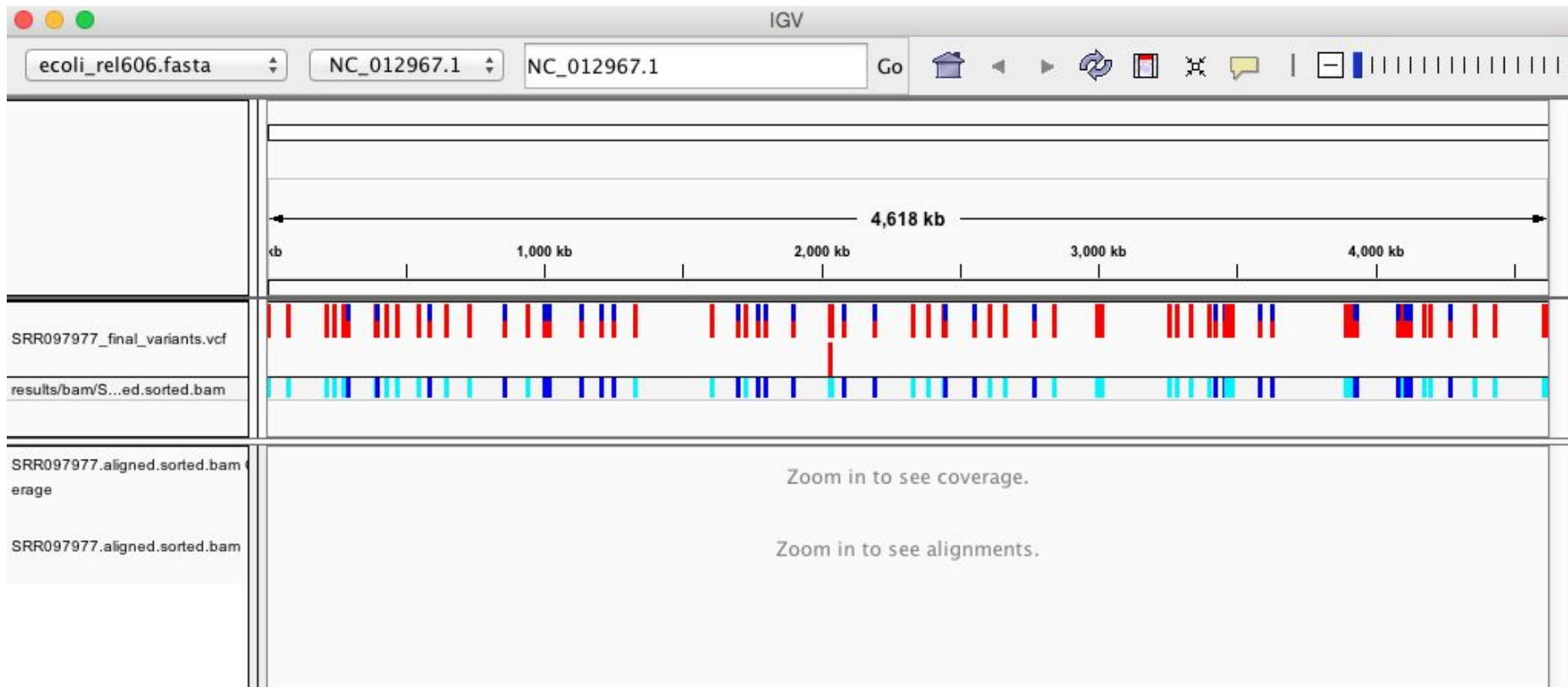
# VCF/BCF files

- 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	CCGCCGTTGCAAAGGCGCGCCG	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

# SNP visualization with IGV



# Tools

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

# Thank you!

## Questions?