## В прошлый раз

**DNA** Sequencing

Бор

Ахо-Корасик

#### Сегодня

Таблица К-меров

Суффиксное дерево

Суффиксный массив

FM-index

# Short Read Alignment



#### Задача

Даны:

шаблон Р длины N, текст Т длины М.

Можно заранее обработать Т.

Найти все позиции вхождения Р в Т.

## К-мер

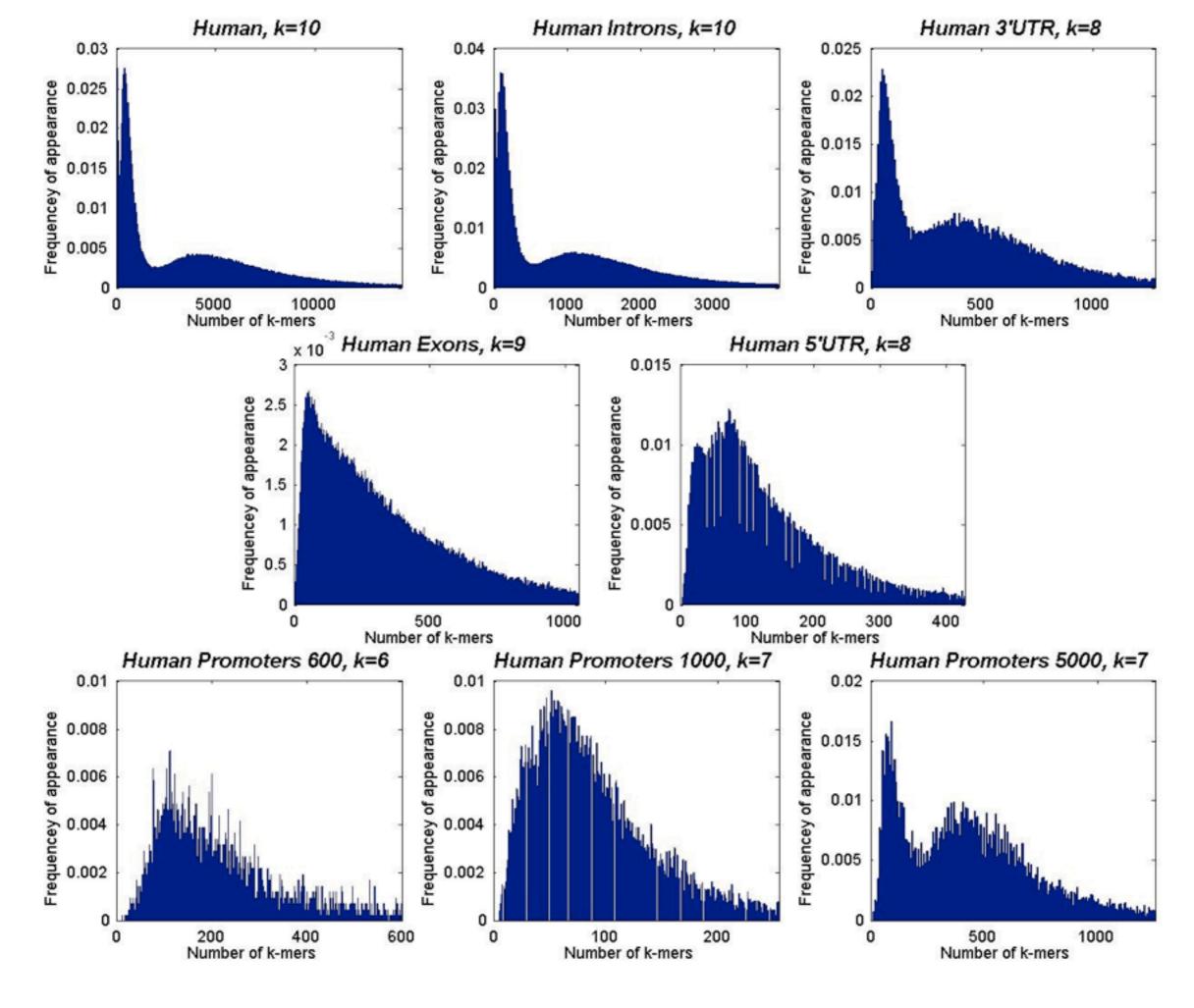
К-мер — слово длины К

I-меры: А С G Т

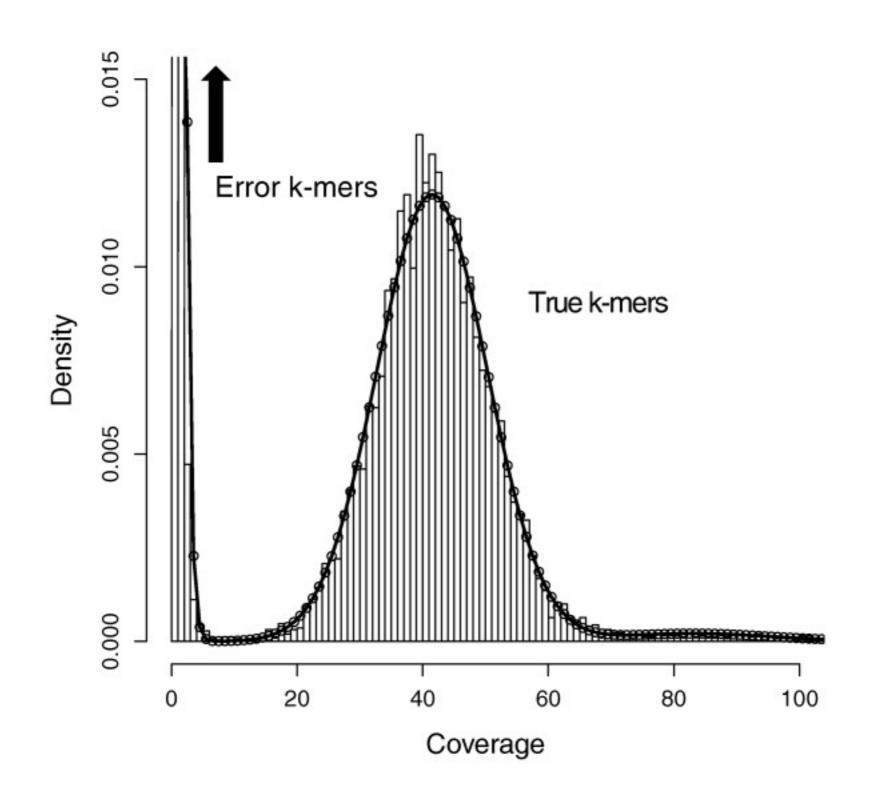
2-меры: AA AC AG AT CA CC CG CT

GA GC GG GT TA TC TG TT

Сколько всего существует К-меров из алфавита {A, C, G, T}?



# Ошибки в ридах



## К-мер

К-мер — слово длины К

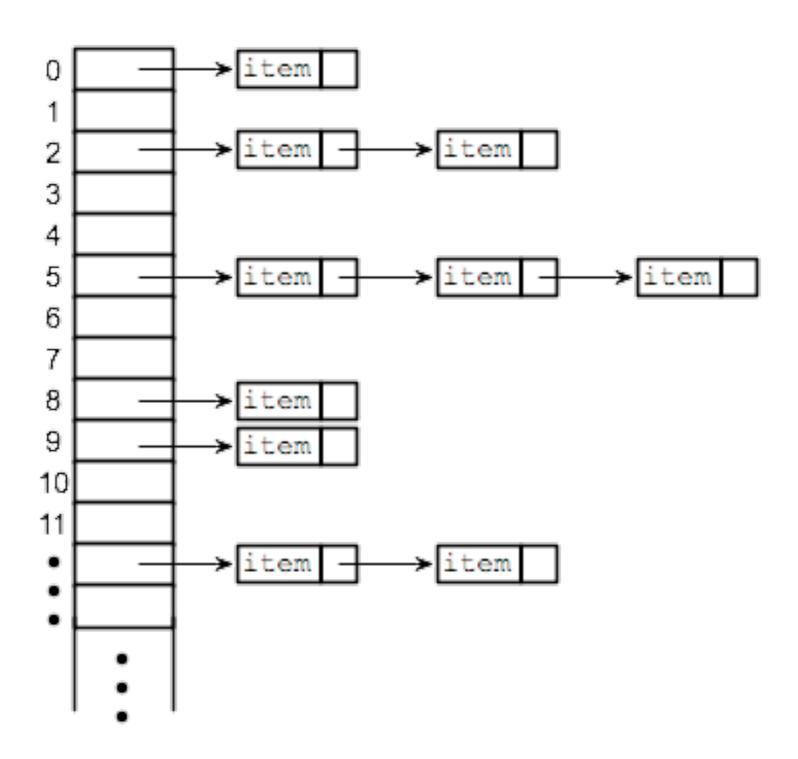
Алфавит  $\{A, C, G, T\} \rightarrow \{0, 1, 2, 3\}$ 

index(10-мер) =  $s_0 \cdot 4^9 + s_1 \cdot 4^8 + ... + s_9 \cdot 4^0$ 

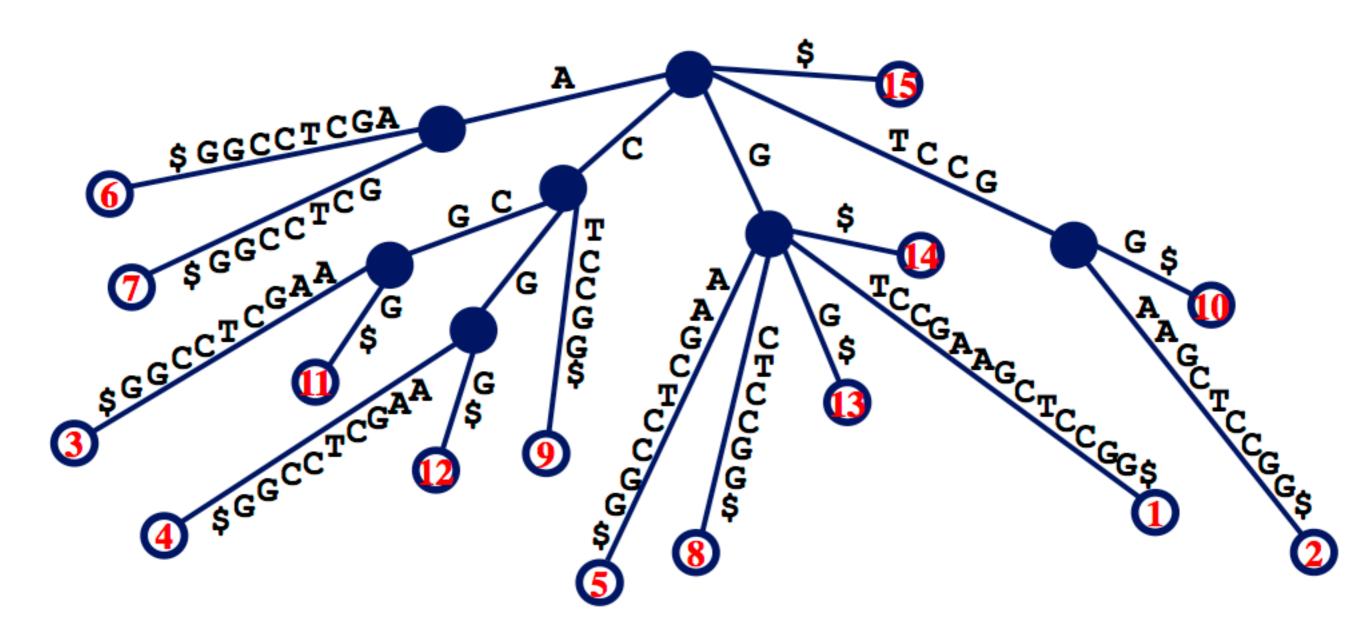
Минимальное значение index?

Максимальное значение index?

## Таблица К-меров



## Суффиксное дерево



#### GTCCGAAGCTCCGG\$

## Суффиксный массив

\$	15
AAGCTCCGG\$	6
AGCTCCGG\$	7
CCGAAGCTCCGG\$	3
CCGG\$	11
CGAAGCTCCGG\$	4
CGG\$	12
CTCCGG\$	9
G\$	14
GAAGCTCCGG\$	5
GCTCCGG\$	8
GG\$	13
GTCCGAAGCTCCGG\$	1
TCCGAAGCTCCGG\$	2
TCCGG\$	10

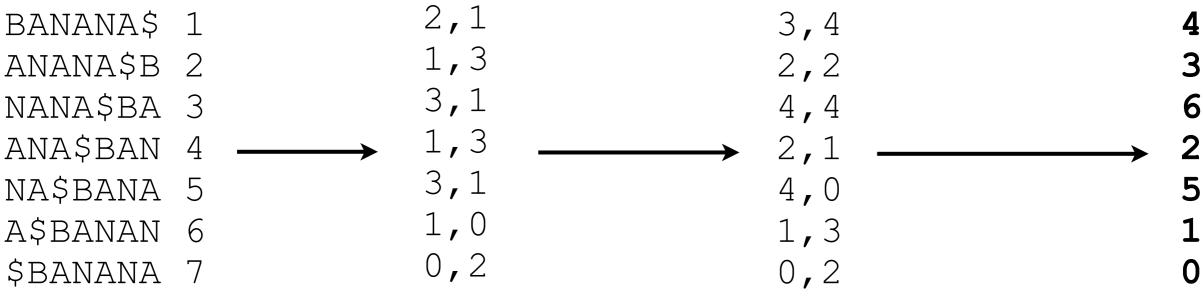
#### GTCCGAAGCTCCGG\$

#### LCP

\$	15	_
AAGCTCCGG\$	6	0
AGCTCCGG\$	7	1
CCGAAGCTCCGG\$	3	0
CCGG\$	11	3
CGAAGCTCCGG\$	4	1
CGG\$	12	2
CTCCGG\$	9	1
G\$	14	0
GAAGCTCCGG\$	5	1
GCTCCGG\$	8	1
GG\$	13	1
GTCCGAAGCTCCGG\$	1	1
TCCGAAGCTCCGG\$	2	0
TCCGG\$	10	3

#### GTCCGAAGCTCCGG\$

#### Суффиксный массив



$$0 = \$$$
  $0 = 02 = \$B$   
 $1 = A$   $1 = 10 = A\$$   
 $2 = B$   $2 = 13 = AN$   
 $3 = N$   $3 = 21 = BA$   
 $4 = 31 = NA$ 

BANANA\$
7 6 4 2 1 5 3