

Поиск подстроки

Николай Вяххи

vyahhi@bioinformaticsinstitute.ru

Computer Science Club
Санкт-Петербург, 2013



**Институт
Биоинформатики**



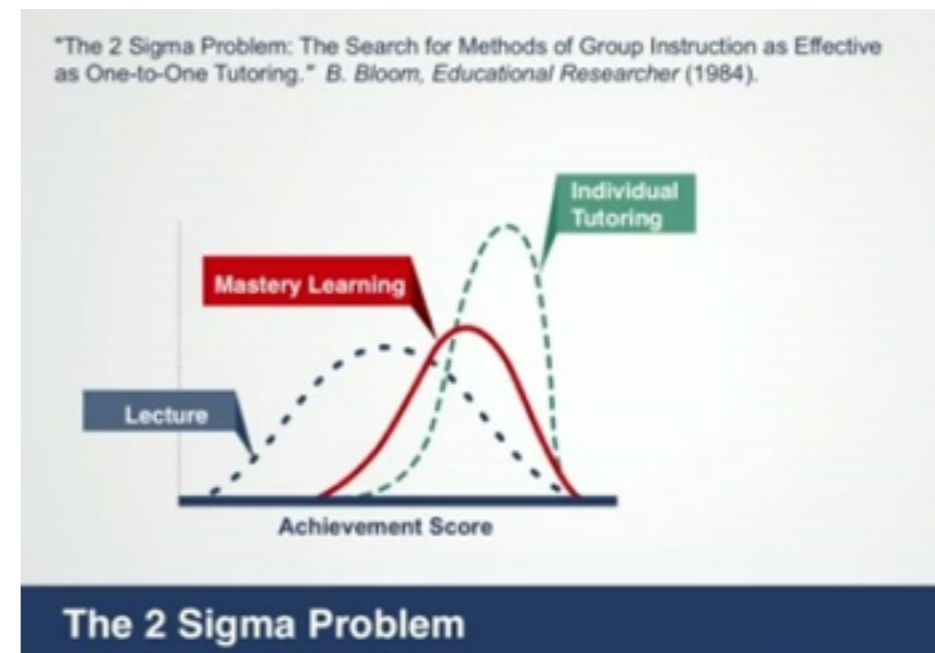
Формат обучения

12 лекций по воскресеньям

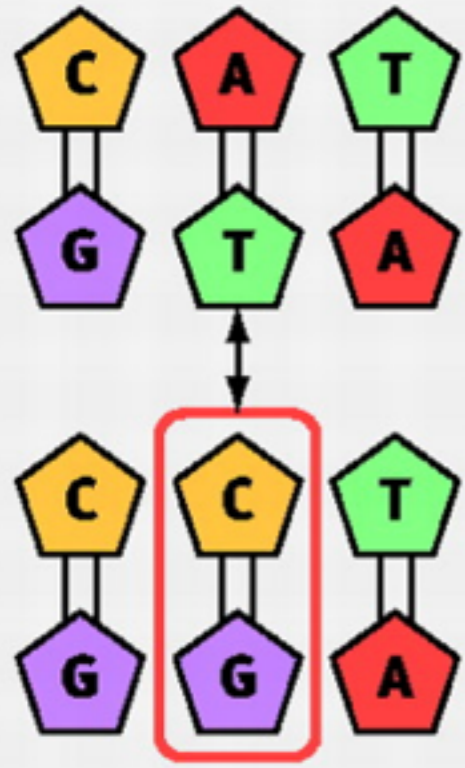
Квизы для самопроверки

Домашние задания и вопросы онлайн

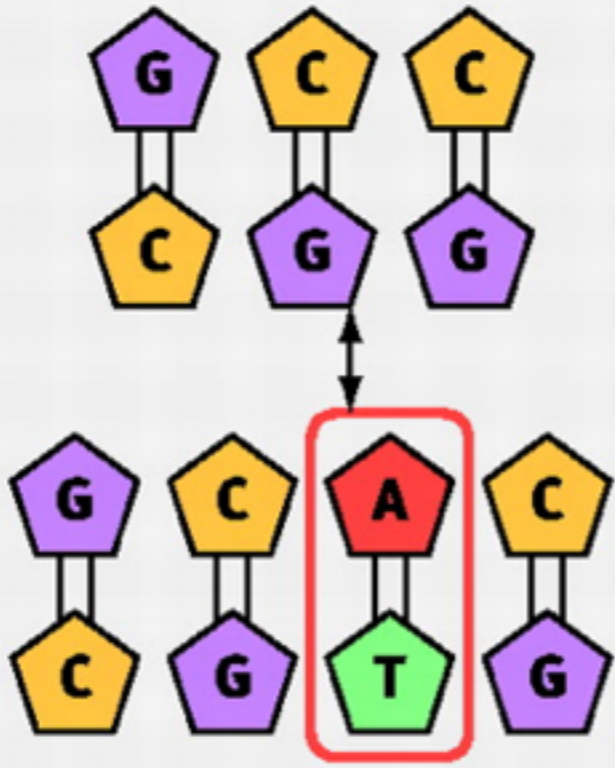
<http://rosalind.info/classes/enroll/ff45302de4/>



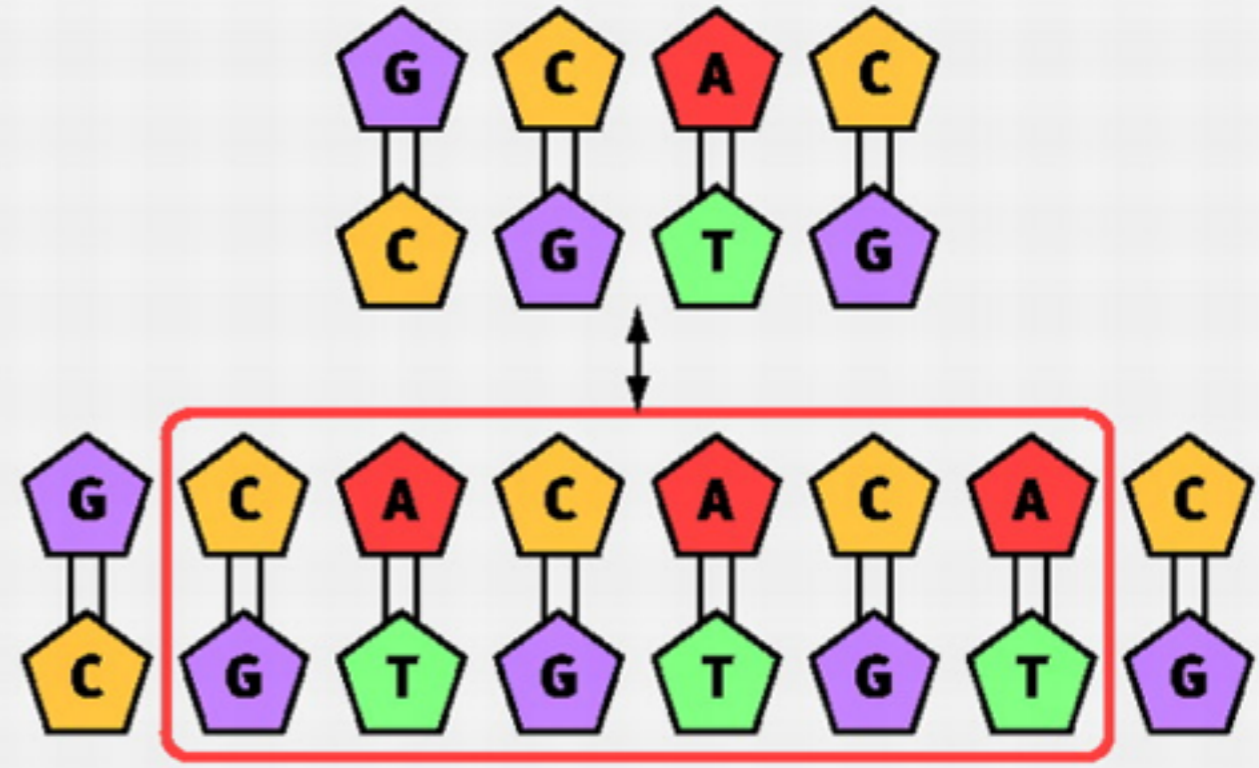
Single nucleotide polymorphism (SNP)



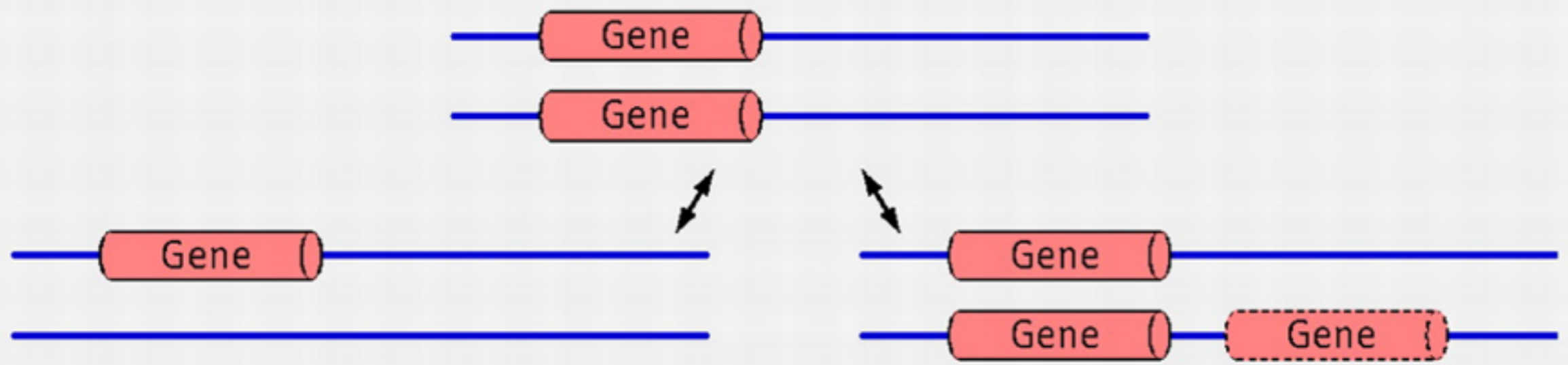
Insertion and deletion polymorphism (indel)



Nucleotide repeat polymorphism



Copy number variation



Deletion

Duplication

Геномные перестройки

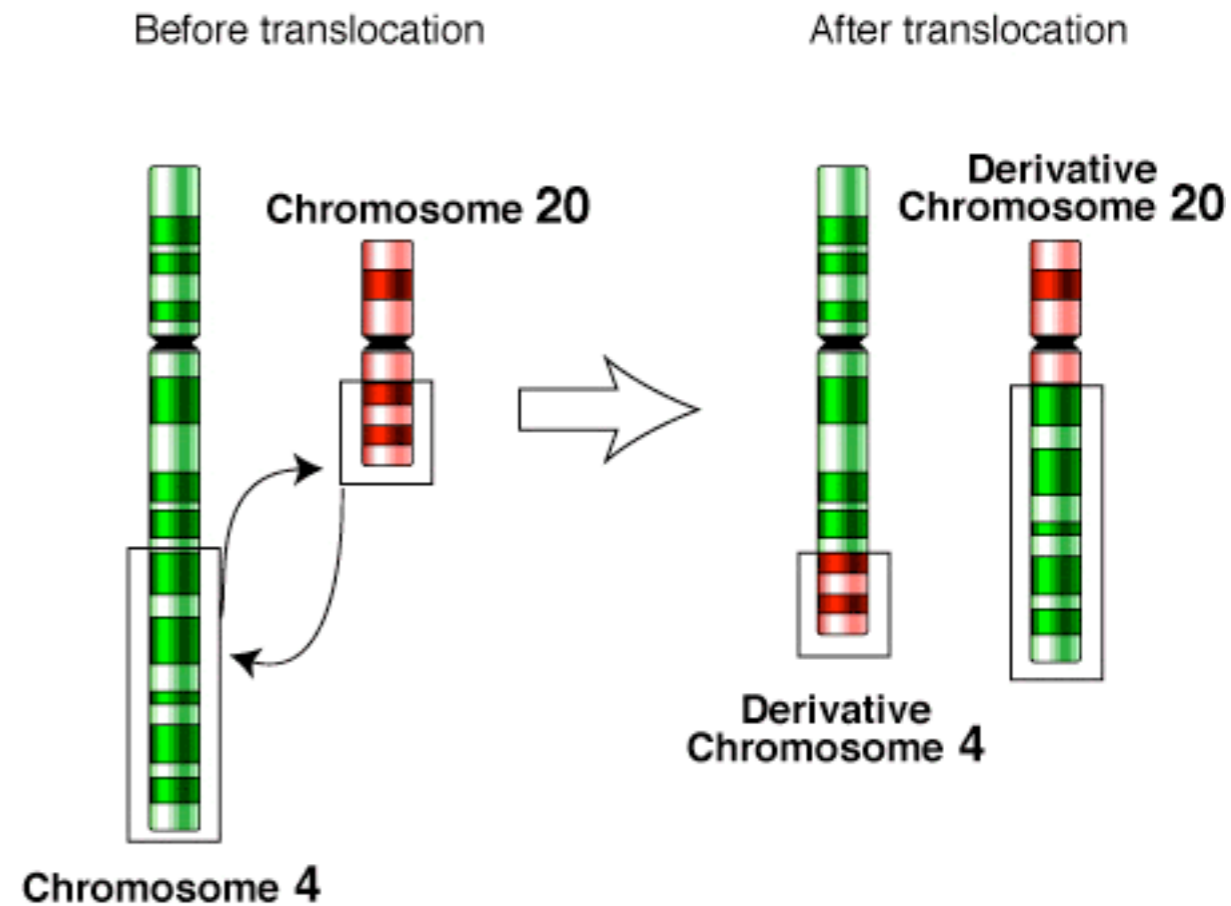
Genome Rearrangements:

реверсия (инверсия)

транслокация

слияние

расщепление



ОТВЕТЫ НА КВИЗ

1. Single Nucleotide Polymorphism, indels, короткие повторы, Copy Number Variations (дубликации, делеции), геномные перестройки (реверсии, транслокации, слияния, расщепления).
2. Геномные перестройки
3. $O(N^3N!)$
4. Понимание эволюции, функций, рака
5. International Union of Pure and Applied Chemistry
6. Локальное выравнивание.

Symbol	Description	Bases represented				
A	adenosine	A				1
C	cytidine		C			
G	guanosine			G		
T	thymidine				T	
U	uridine				U	
W	weak	A			T	2
S	strong		C	G		
M	amino	A	C			
K	keto			G	T	
R	purine	A		G		
Y	pyrimidine		C		T	3
B	not A (B comes after A)		C	G	T	
D	not C (D comes after C)	A		G	T	
H	not G (H comes after G)	A	C		T	
V	not T (V comes after T and U)	A	C	G		4
N or -	any base (not a gap)	A	C	G	T	

Сегодня

Наивный алгоритм

Кнут-Моррис-Пратт

Рабин-Карп

Бойер-Мурр

Неточный поиск с заменами

Задача

Поиск подстроки в строке:

Даны:

шаблон P длины N ,

текст T длины M .

Найти все позиции вхождения P в T .

Наивный алгоритм

АТСА

АТАТГААСТGAGATCAAT

Наивный алгоритм

АТСА

АТАТГААСТGAGATCAAT

Наивный алгоритм

ATCA

ATATGAACTGAGATCAAT

Наивный алгоритм

АТСА

АТАТГААСТGAGATCAAT

Наивный алгоритм

АТСА

АТАТГААСТGAG**АТСА**АТ

Кнут–Моррис–Пратт

Префикс функция — длина наибольшего префикса строки, который не совпадает с этой строкой и одновременно является её суффиксом.

$$\begin{aligned} S &= \text{abacaba} \\ \pi(S) &= 0010123 \end{aligned}$$

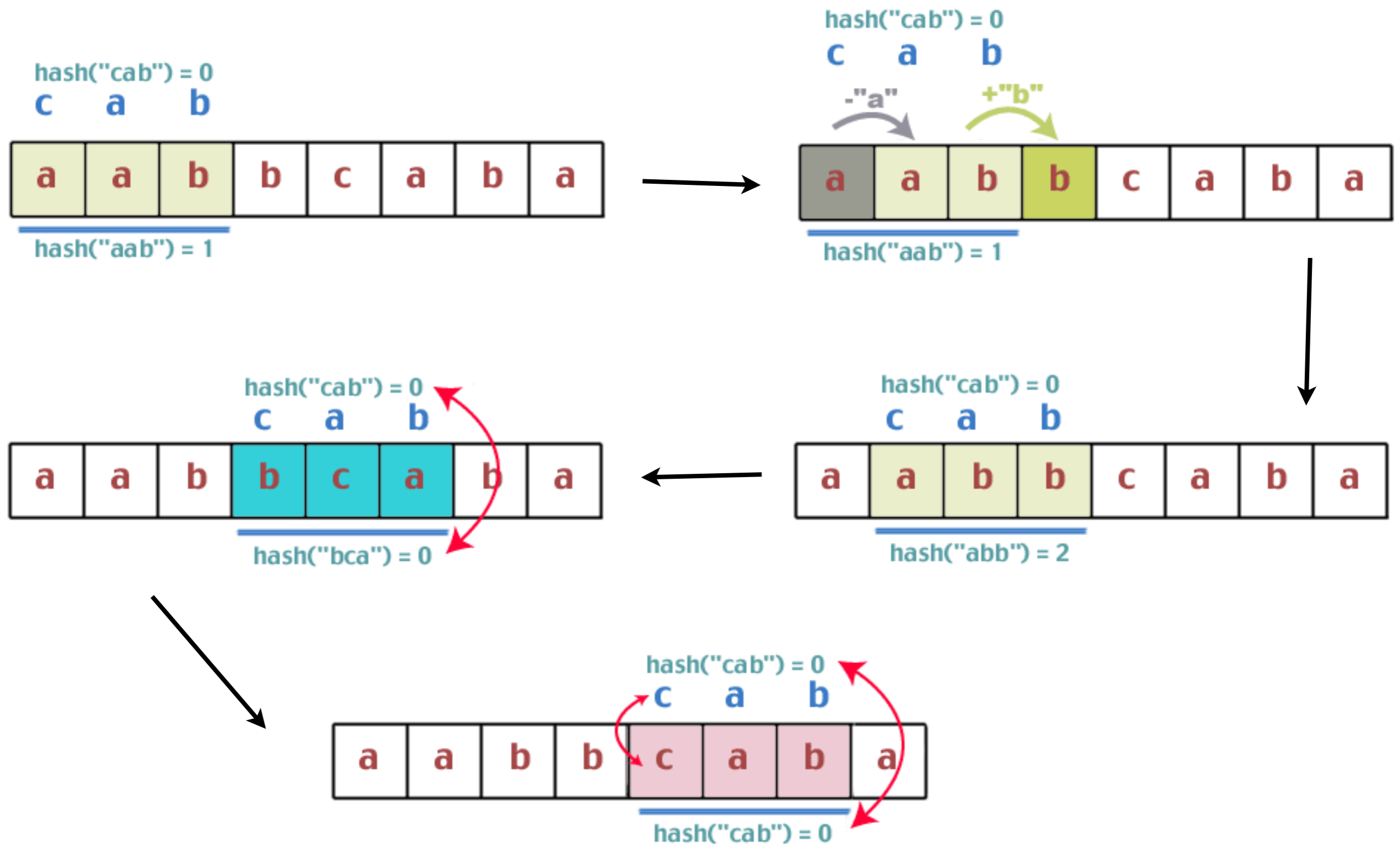
Кнут–Моррис–Пратт

Префикс функция — длина наибольшего префикса строки, который не совпадает с этой строкой и одновременно является её суффиксом.

$$\begin{aligned} S &= \text{abacaba} \\ \pi(S) &= 0010123 \end{aligned}$$

Построим префикс функцию от P\$T...

Рабин–Карп



Бойер–Мурр

Сканирование слева направо.

Сравнение справа налево.

Бойер–Мурр

Сканирование слева направо.

Сравнение справа налево.

The Bad Character Rule

(Эвристика стоп-символа)

The Good Suffix Rule

(Эвристика совпавшего суффикса)

The Bad Character Rule

A N P A N M A N A M -
- N N A A M A N - - -
- - - N N A A M A N -

The Bad Character Rule

A N P A N M A N A M -
- N N A A M A N - - -
- - - N N A A M A N -

Препроцессинг: таблица стоп-символов.

The Bad Character Rule

A N P A N M A N A M -
- N N A A M A N - - -
- - - N N A A M A N -

Препроцессинг: таблица стоп-символов.

N	A	M
1	5	4

Бойер–Мур–Хорспул

Изменённая эвристика стоп-символа.

A N P A N M A N A M -
- N N A A M A N - - -
- - - - - N N A A M

The Good Suffix Rule

M A N P A N A M A N A P -
A N A M P N A M - - - - -
- - - - A N A M P N A M -

The Good Suffix Rule

MANP ANAM ANAP -
ANAMP NAM - - - -
- - - - ANAMP NAM -

Препроцессинг: таблица суффиксов

The Good Suffix Rule

M A N P A N A M A N A P -
A N A M P N A M - - - - -
- - - - A N A M P N A M -

Препроцессинг: таблица суффиксов

abacaba

4444462

Поиск с ошибками

Даны:

шаблон P длины N ,

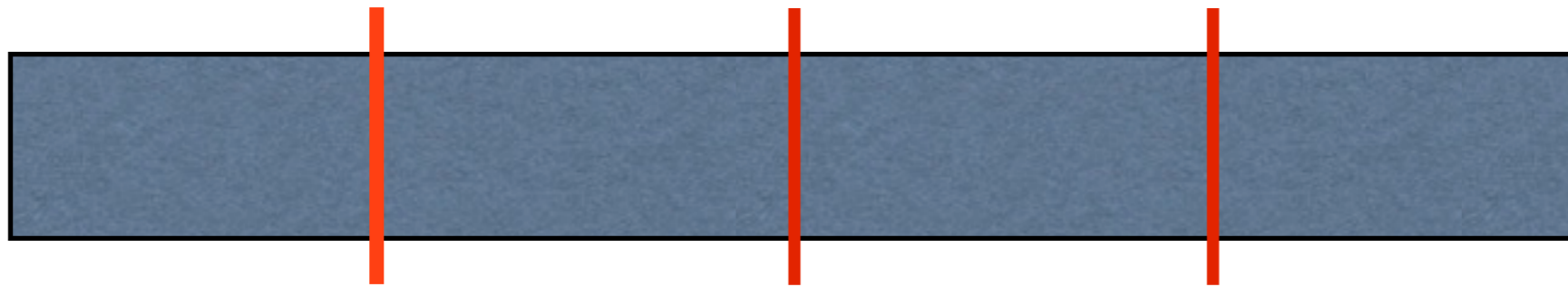
текст T длины M ,

целое число K .

Найти все позиции вхождения P в T с максимум K ошибками (заменами).

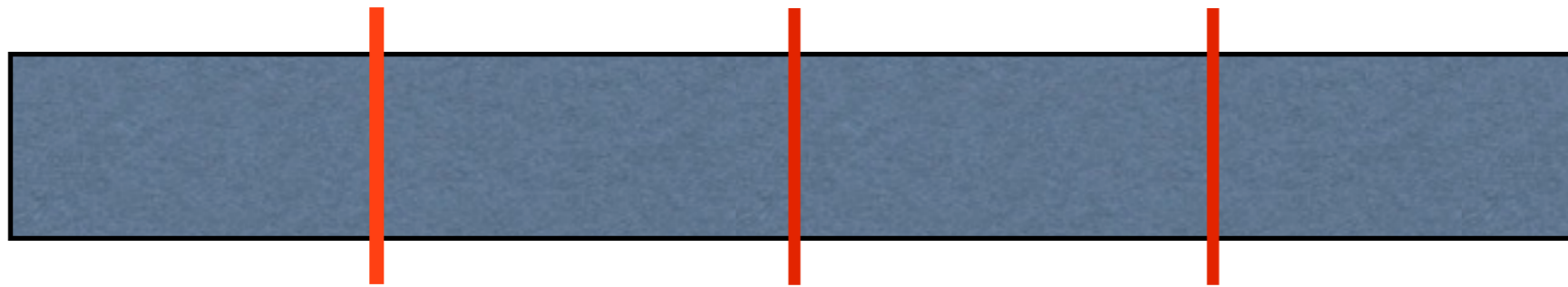
Решение

Разобьём шаблон на $K+1$ фрагментов.



Решение

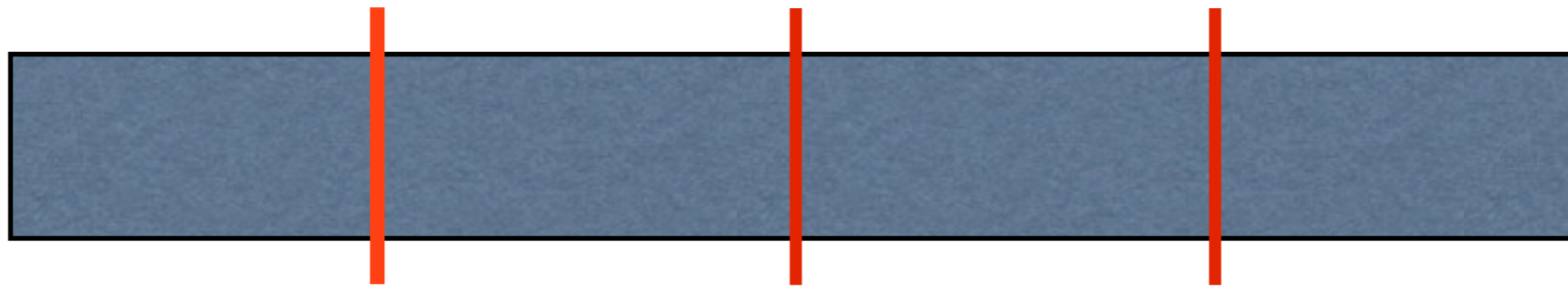
Разобьём шаблон на $K+1$ фрагментов.



Один из фрагментов должен
встречается в тексте точно.

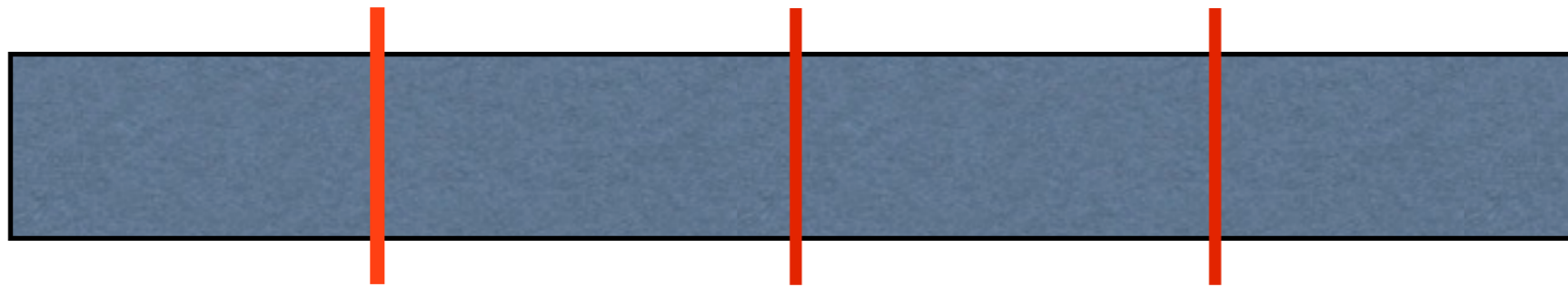
Решение

Найдём вхождения всех $K+1$ фрагментов
с помощью любого точного алгоритма.



Решение

Найдём вхождения всех $K+1$ фрагментов с помощью любого точного алгоритма.



Расширим все вхождения и проверим количество ошибок.

СЛОЖНОСТЬ

Найти фрагменты:

$$O(K \cdot (M + N/K)) = O(KM + N)$$

Проверить одно вхождение: $O(M)$

Чем больше K , тем больше ложных вхождений нам проверять...

Новая задача

Даны:

шаблон P длины N ,

текст T длины M .

Можно заранее обработать T .

Найти все позиции вхождения P в T .

Что мы узнали

Наивный алгоритм

Кнут-Моррис-Пратт

Рабин-Карп

Бойер-Мурр

Неточный поиск с заменами

Формат обучения

12 лекций по воскресеньям

Квизы для самопроверки

Домашние задания и вопросы онлайн

<http://rosalind.info/classes/enroll/ff45302de4/>

