



## Часть 2

# Как выравнивать очень длинные последовательности?



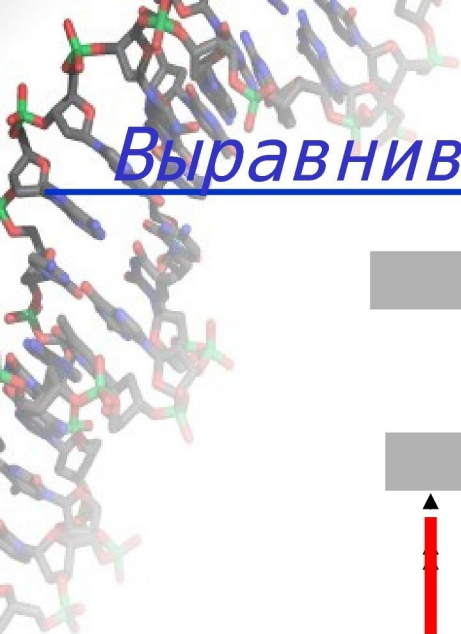
## Выравнивание

- Предполагаем наличие общего предка
- Оптимизируем вес выравнивания
- Есть эталонные выравнивания (для подбора весов и оценки качества)
- Время  $\sim L^2$ ;  $L \sim 10^2 - 10^3$

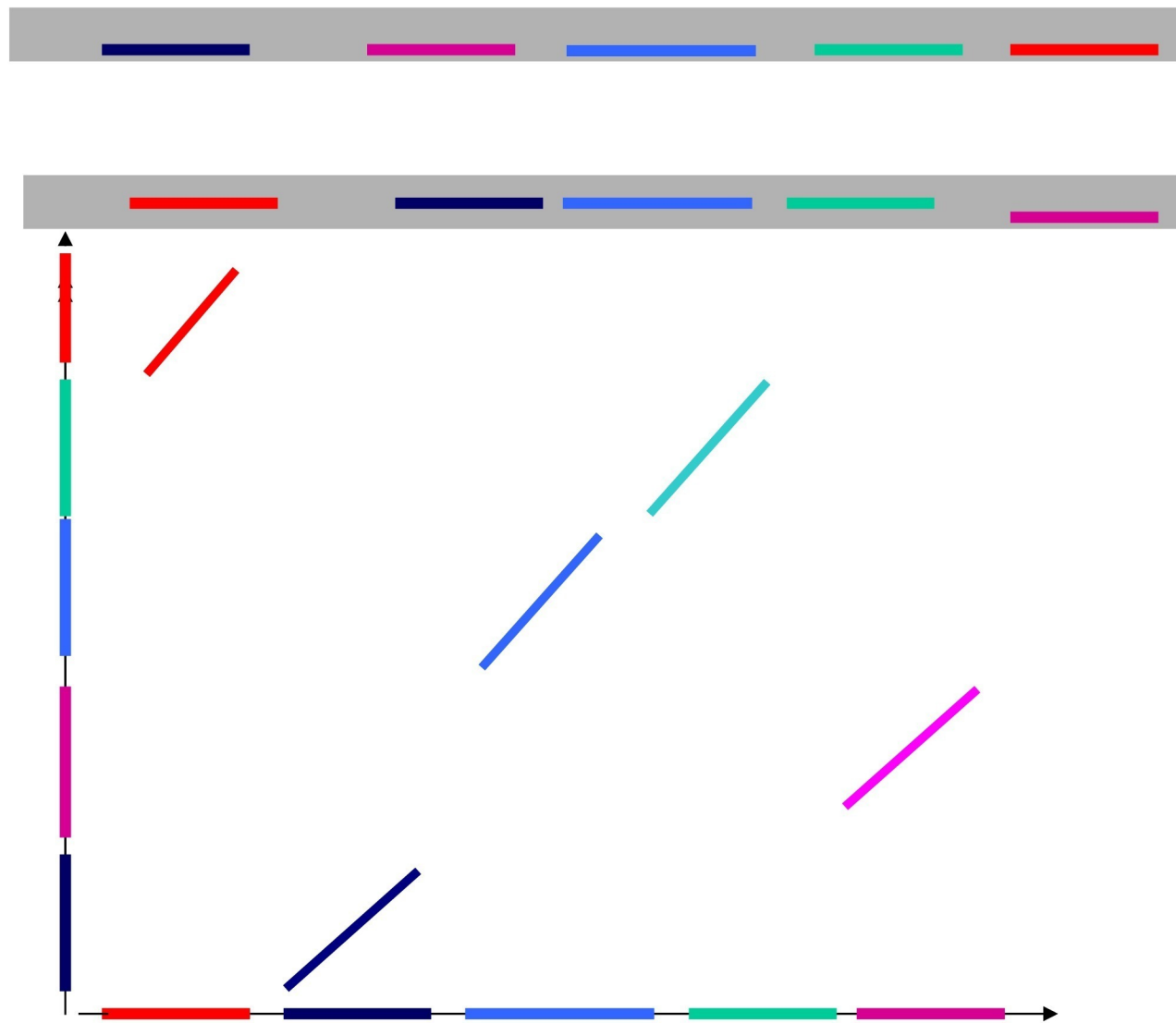
### **ОЧЕНЬ ДЛИННЫЕ ПОСЛЕДОВАТЕЛЬНОСТИ**

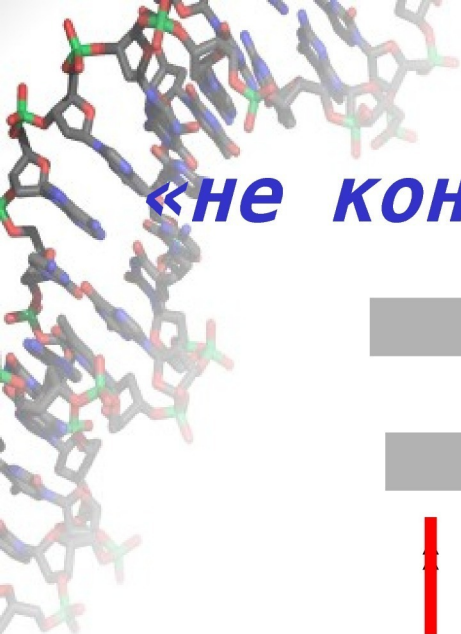
- $L \sim 10^7$
- Эталонных выравниваний нет

## **ЧТО ДЕЛАТЬ?**



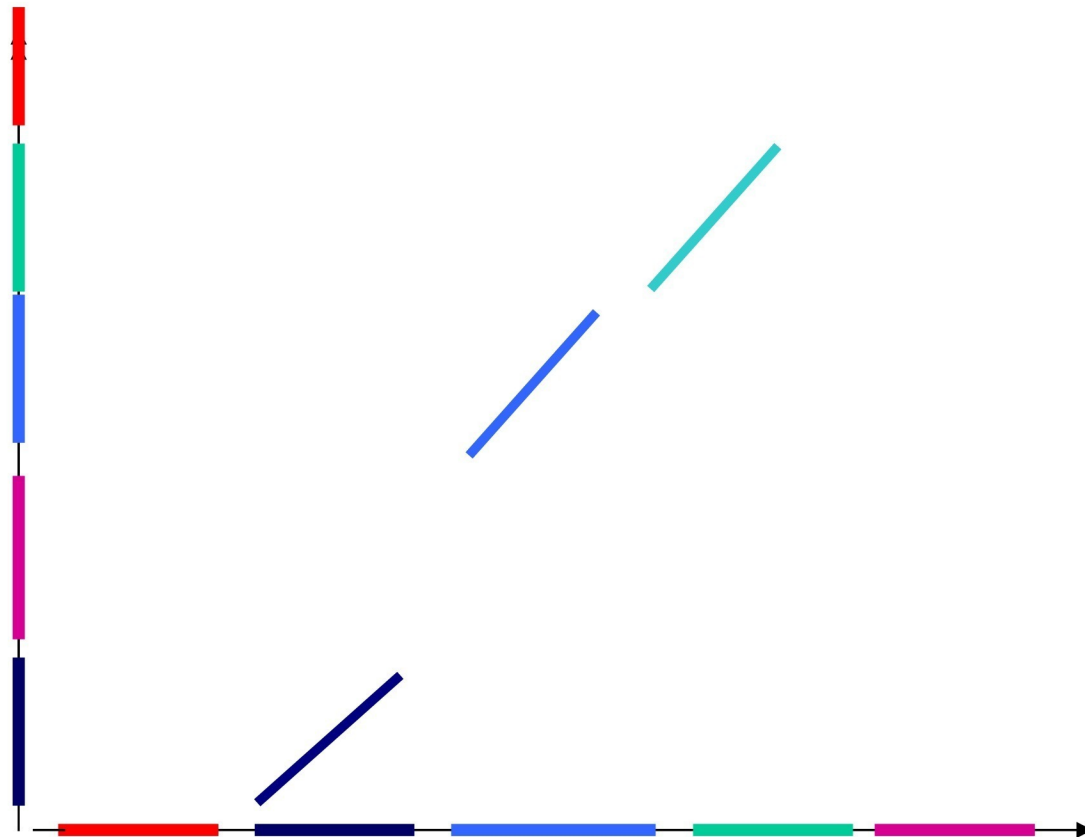
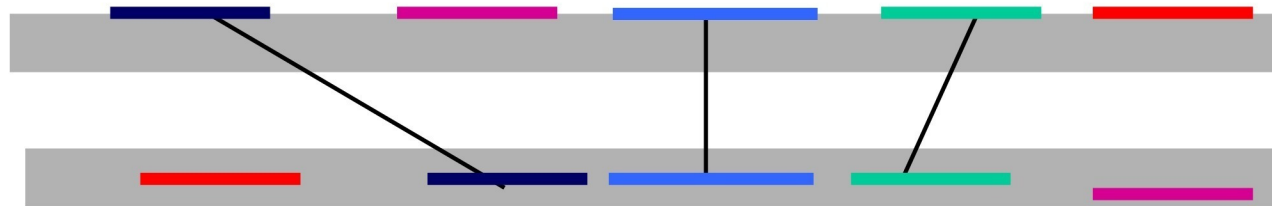
# Выравнивание на основе локальных сходств

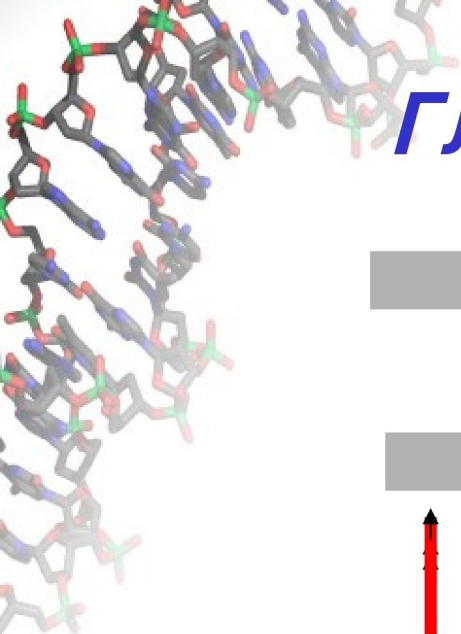




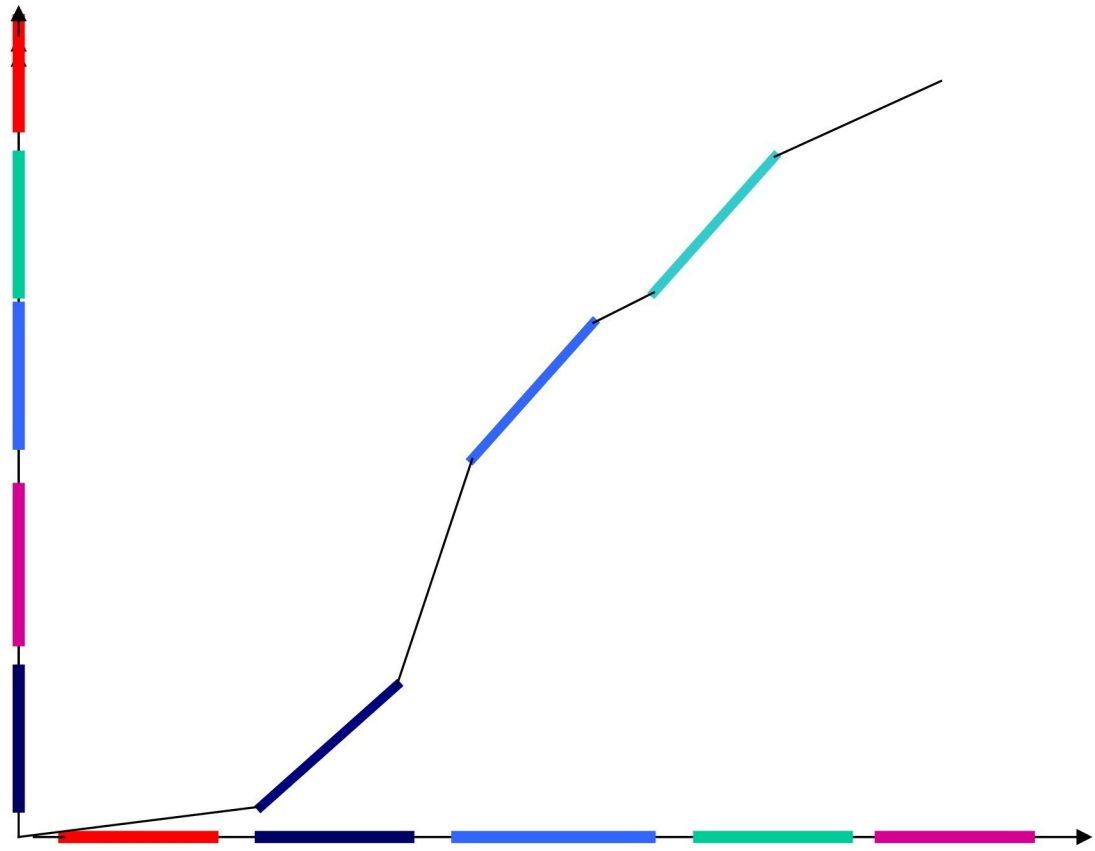
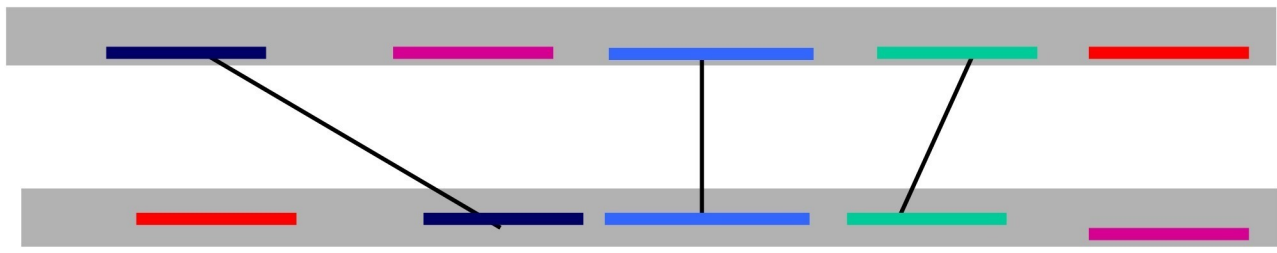
# Цепь

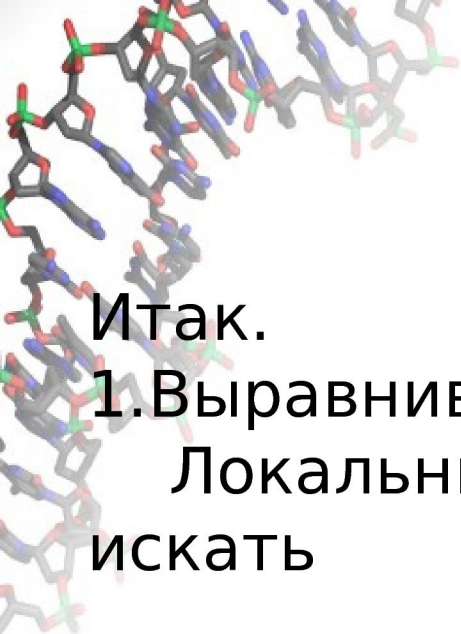
**«не конфликтующих» локальных сходств**





# *Глобальное выравнивание*



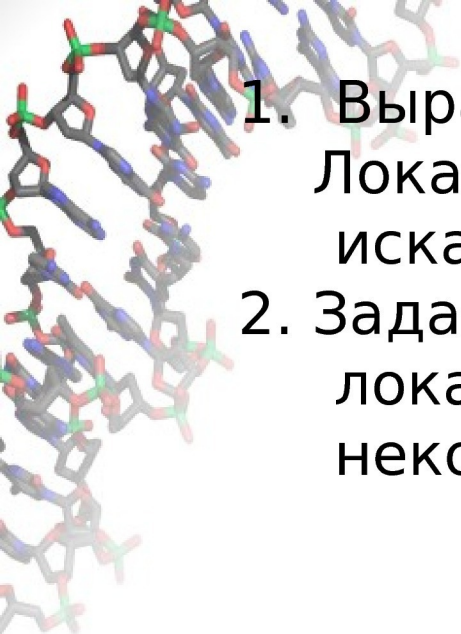


Итак.

1. Выравнивание – это цепочка локальных сходств.

Локальные (неточные) сходства можно легко  
искать

2. Задача выравнивания – выделить в множестве  
локальных сходств «разумную» цепочку  
неконфликтующих локальных сходств

- 
1. Выравнивание – это цепочка локальных сходств. Локальные (неточные) сходства можно легко искать
  2. Задача выравнивания – выделить в множестве локальных сходств «разумную» цепочку неконфликтующих локальных сходств

## **ПРОБЛЕМА:**

Чтобы получить детальное выравнивание, исходных локальных сходств должно быть много

## **РЕШЕНИЕ:**

**Иерархическое выравнивание**



## *Иерархическое выравнивание*

- отказ от глобальной оптимизации какой-либо весовой функции;
- разрешение конфликта между сходствами производится локально – путем сравнения их **статистической значимости**
- иерархическая процедура





## Статистическая значимость локального сходства

- Каждому сходству  $S$  приписывается вес  $W(S)$ .  
*Статистическая значимость локального сходства  $S$  в последовательностях длин  $m$  и  $n$  – это вероятность того, что в независимых случайных последовательностях длин  $m$  и  $n$  есть сходство того же или большего веса.*

**Нужно уточнить распределение  
вероятностей на последовательностях**

# Как вычислить статистическую значимость $p(W)$

## локального сходства по его весу $W$

Обозначения:

$Z(v_1, v_2)$  – вес оптимального безделеционного локального сходства последовательностей  $v_1$  и  $v_2$ .

Весовая матрица  $M[a, b]$  считается заданной

$P(w)$  – вероятность того, что для независимых случайных бернуллиевских последовательностей  $v_1$  и  $v_2$

длины  $m$  и  $n$  соответственно выполнено:  $Z(v_1, v_2) > w$ .

Теорема (Karlin, Altschul, 1999).

$$P(w) \approx \exp(-Kmn e^{-\lambda w})$$

здесь  $K, \lambda$  зависят только от весовой матрицы  $M$ .

Термин: Распределение экстремальных значений.



## Условия применимости распределения экстремальных значений

■ Теорема (Karlin, Altschul, 1999).

$$P(w) \approx \exp(-Kmn e^{\lambda w}) \quad (*)$$

Теория: Необходимо выполнение следующего:

1)  $\text{Exp}(M[x, y]) < 0$

2)  $\text{Max}(M[x, y]) > 0$

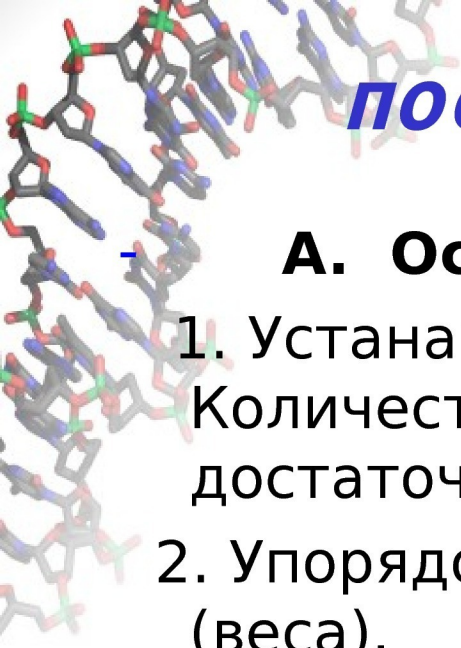
3) *распределение вероятностей – бернуллиево*

4) *рассматриваются только безделеционные  
сходства*

Практика:

Условия 1) и 2) в реальных случаях всегда выполняются

Если условия 3) и 4) не выполняются, то (\*) можно использовать приближенно, а коэффициенты **K**, **λ** определить методом Монте-Карло



# Иерархическая процедура построения цепочки локальных СХОДСТВ

## А. Основной шаг (жадный алгоритм)

1. Устанавливаем порог статистической значимости. Количество локальных сходств –кандидатов достаточно мало.
  2. Упорядочиваем сходства по убыванию значимости (веса).
  3. Берем очередное сходство. Если оно ни с кем не конфликтует – включаем в цепочку.  
Иначе – отбрасываем [на самом деле – откладываем для дальнейшего изучения]
- \*\* Техническая проблема: повторы (“low complexity”).  
Нужно отдельно фильтровать.



# *Иерархическая процедура построения цепочки локальных*

## *СХОДСТВ*

### **Б. Иерархический переход**

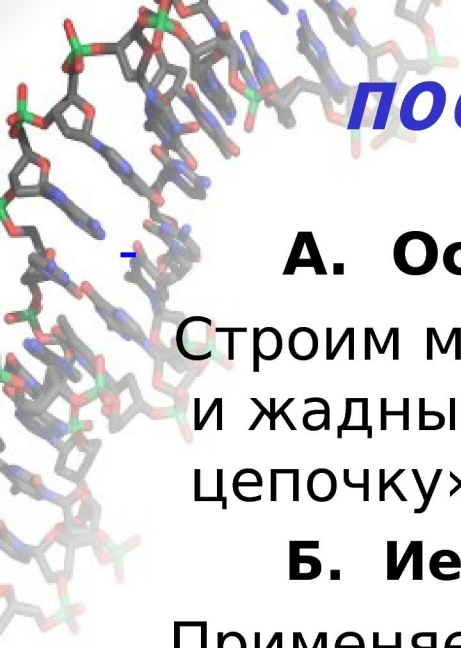
Применяем основной шаг отдельно к каждой из областей поиска между отобранными локальными сходствами.

!!! Из-за того, что области стали меньше, при том же значении порога

***СХОДСТВА СТАНОВЯТСЯ ЗНАЧИМЫМИ***

***при меньшем весе***

Поэтому промежутки между сходствами в цепочке будут постепенно заполняться



# *Иерархическая процедура построения цепочки локальных сходств*

## **А. Основной шаг**

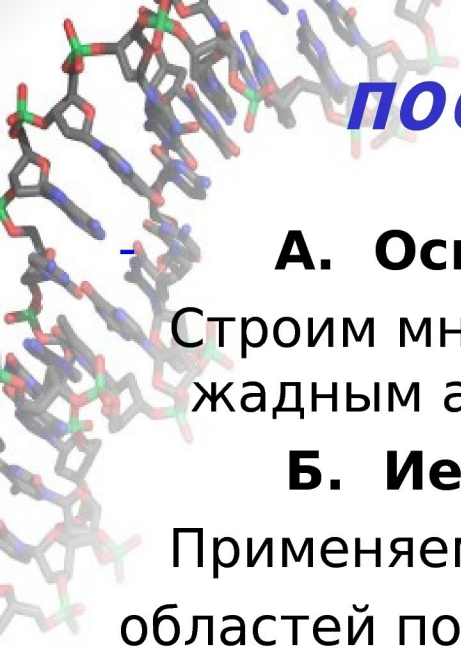
Строим множество статистически значимых сходств и жадным алгоритмом выбираем «основную цепочку».

## **Б. Иерархический переход**

Применяем основной шаг отдельно к каждой из областей поиска между отобранными локальными сходствами.

## **В. Итерирование**

Повторяем шаг **Б** до тех пор, пока появляются новые значимые сходства.



# *Иерархическая процедура построения цепочки локальных сходств*

## **А. Основной шаг**

Строим множество статистически значимых сходств и жадным алгоритмом выбираем «остовную цепочку».

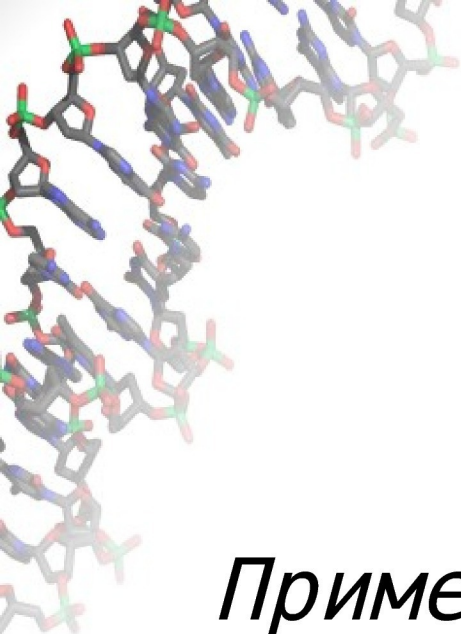
## **Б. Иерархический переход**

Применяем основной шаг отдельно к каждой из областей поиска между отобранными локальными сходствами.

## **В. Итерирование**

Повторяем шаг **Б** до тех пор, пока появляются новые значимые сходства.

*Время определяется принятым уровнем значимости (а не минимальным допустимым весом сходства)*



*Пример: иерархическое геномное  
выравнивание (длина  $\sim 10^7$ )*







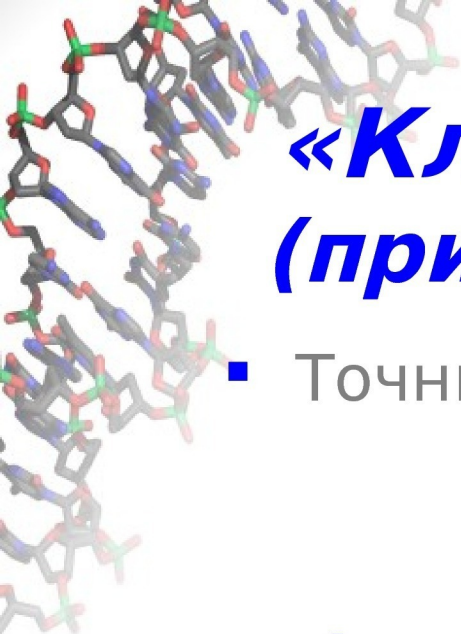




## **Часть 3. Поиск всех разумных локальных сходств**

- 3.1. Введение
- 3.2. Использование затравок (seed)
- 3.3. Избирательность и чувствительность
- 3.4. Типы затравок (seed model)





# «Классическая затравка» (пример: 6 совпадений подряд)

- Точные совпадения :  
ATCAGT  
| | | | |  
ATCAGT

Затравка («затравочное слово», описание затравочных сходств) : #####

**Вес : 6** [количество #]

- Пример : 16 совпадений из 20

#####

ATCAGT**GCAATG**СТСАТГАА

| | | . | . **| | | | |** : | | . | | |

ATCGGC**GCAATG**CGCAAGAA

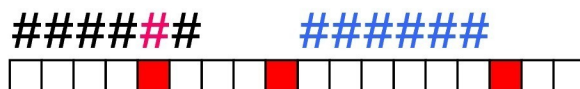
# Затравка ЛОВИТ СХОДСТВО (затравка соответствует сходству)

■ Затравка ##### ← *seed*

Затравочное сходство (... выравнивание)

ATGCAA

ATGCAA



Затравка <sup>1</sup>соответствует <sup>10</sup>сходству в позиции 10

Затравка **не соответствует** сходству в позиции 1

Затравка **ЛОВИТ** сходство







# Итак:

- “Избирательность”

Затравка может НЕ быть частью ***важного (для нас) сходства***

- “Чувствительность”

***Важное (для нас) сходство*** **МОЖЕТ НЕ** содержать ни одной затравки

**Нужно уточнить:**

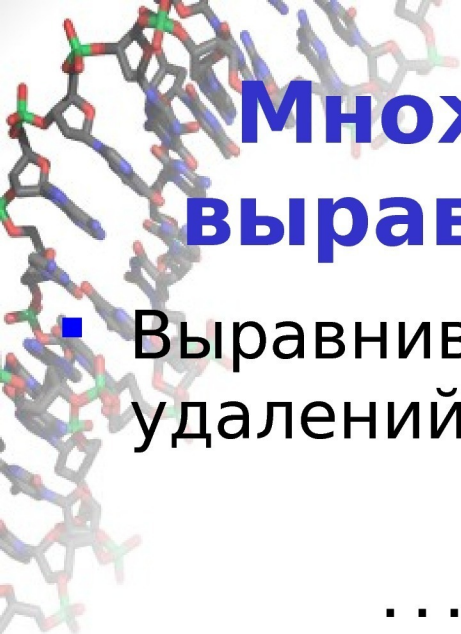
- **Что такое «важное сходство»?**

# Что может быть мерой избирательности и чувствительности

- Избирательность заправки:  $\sim 4\text{-weight}$   
вероятность ее обнаружения при  
сравнении независимых случайных  
последовательностей
- Чувствительность заправки:  
***вероятность*** того, что заправка попадет в  
***важное сходство***.

**Нужно уточнить:**

- **Что такое «важное сходство»?**
- **Каково распределение вероятностей для важных сходств?**



# Множество важных [целевых] выравниваний и их вероятности

- Выравнивания фиксированной длины без удалений

GCTACGACTTCGAGCTGC



...СТСАГСТАТГАССТСГАГСГГССТАТСТА.  $L=18$

- Вероятностная модель: **Бернулли** ;  
*Случайные* выравнивания:  $Prob(match) = 0.25$   
*Целевые* выравнивания:  $Prob(match) \gg 0.25$

*Обобщения: Марковские модели, скрытые марковские модели (сегодня не рассматриваем)*



# Разрезанные затравки

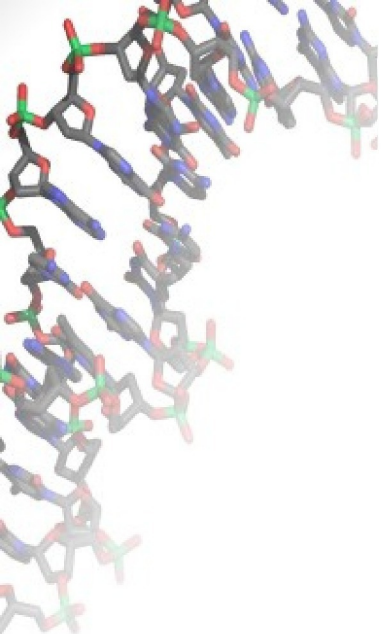
Ma, Tromp, Li 2002 (PatternHunter)

- Затравка: ### - -# - ##  
{ '#' : должно быть совпадение  
'-' : «джокер» (“все равно, что” )

**Вес : 6** [количество #]

- **Пример:**

```
### - -# - ##  
ATCAGTGC AATGCTCAAGA  
| | | | . | | . | | | | : | | | |  
ATCAGCGC GATGCGCAAGA
```



#####

ATCAGTGC**A**ATGCT**T**CAAGA

|||||.|||.||||:|||||

ATCAG**C**GC**G**ATGC**G**CAAGA

#####

#####

#####

#####

#####

#####

#####

#####

#####

#####

#####

#####

#####

#####

###- -#-##

ATCAGTGC**A**ATGCT**T**CAAGA

|||||.|||.||||:|||||

ATCAG**C**GC**G**ATGC**G**CAAGA

###- -#-##

###- -#-##

###- -#-##

###- -#-##

###- -#-##

###- -#-##

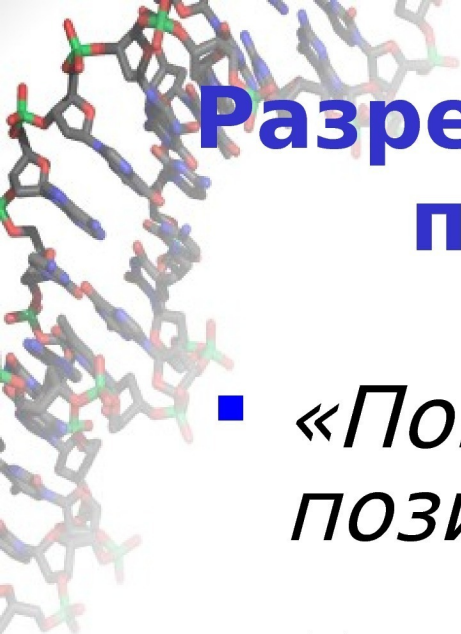
###- -#-##

###- -#-##

###- -#-##

###- -#-##

###- -#-##

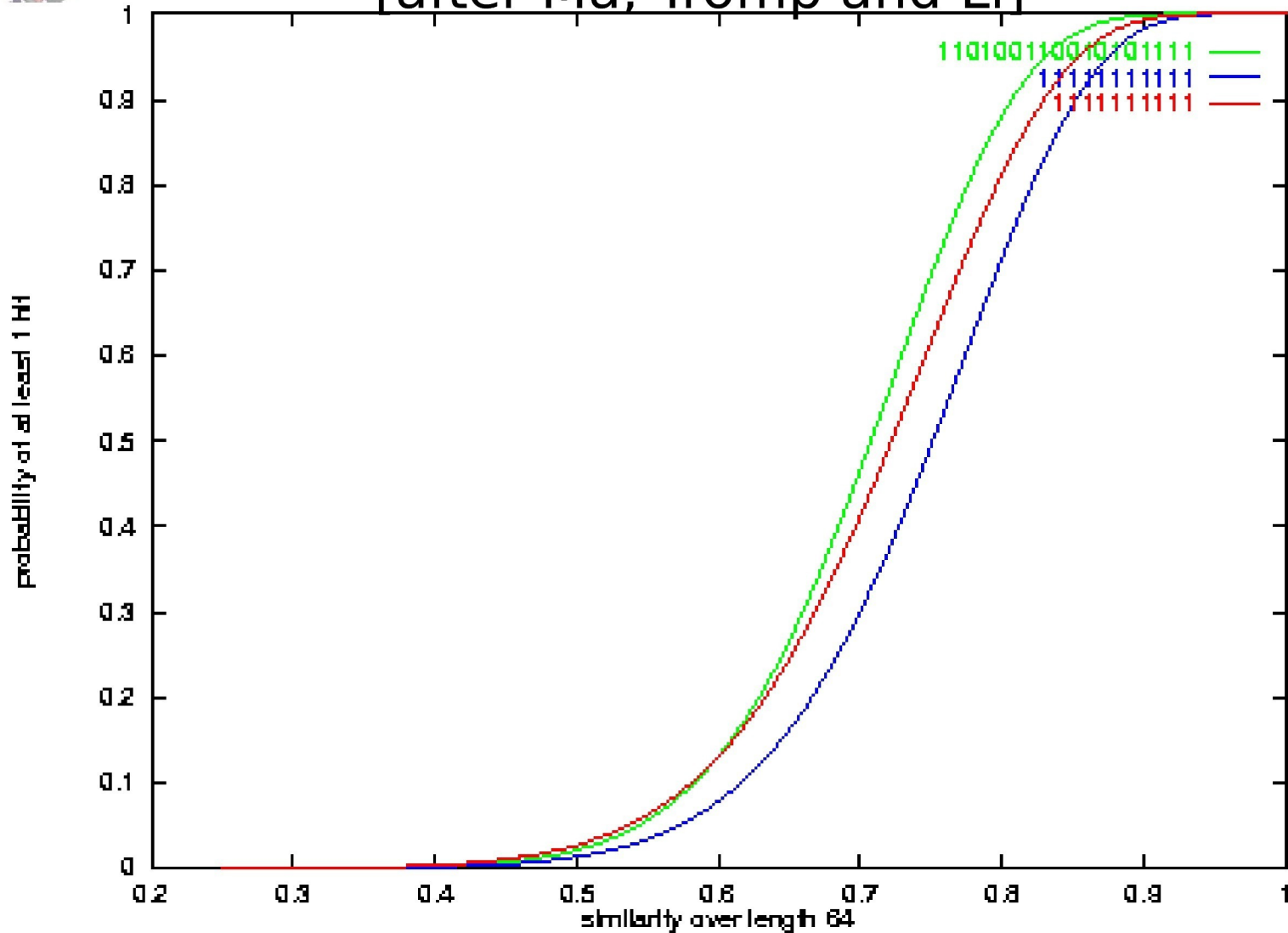


## Разреженные затравки: в чем преимущество перед классическими?

- «Попадания» в соседних позициях менее зависимы
- *For contiguous vs. spaced seeds of the same weight, the expected number of hits is (basically) the same but the probabilities of having **at least one hit** are very different*

# Sensitivity: PH weight 11 seed vs BLAST 11 & 10

[after Ma, Tromp and Li]





## 3.4.1. Семейства затравок

- single filter based on several distinct seed patterns
- each seed pattern detects a part of interesting similarities but together they detect [almost] all of them
- Li, Ma, Kisman, Tromp 2004 (PatternHunter II)
- Kucherov, Noe, Roytberg, 2005
- Sun, Buhler, RECOMB 2004





## Пример: ВСЕ (18,3)

«Ловить» **все** сходства длины **18**,  
в которых не более **3** несовпадений

**Чувствительность = 1.0**

## **Избирательность**

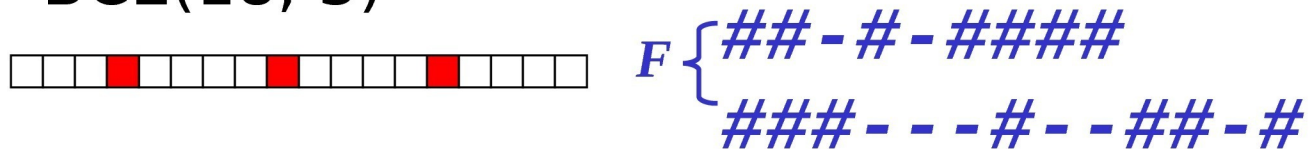
*(вероятность случайного появления*

*затравочного сходства) -> **MIN***

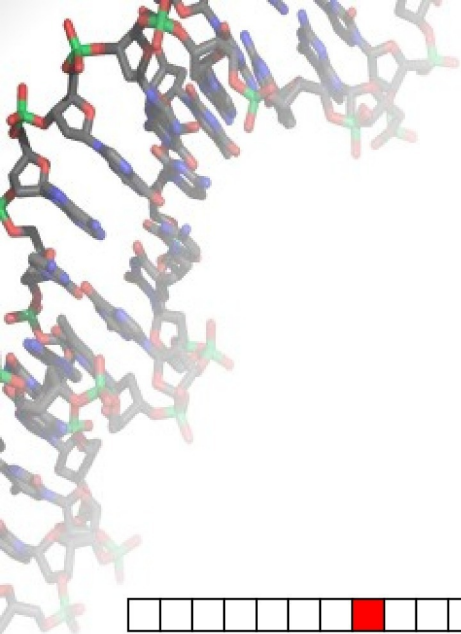
# Пример: ВСЕ (18,3)

Обнаружить **все** сходства длины **18**,  
в которых не более **3** несовпадений

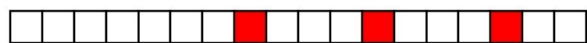
Множественная заставка  $F$  решает проблему  
ВСЕ(18, 3)



Заставка  $F$  состоит из двух простых  
заставок, каждая из них имеет вес 7

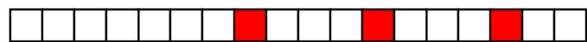


# Пример: ВСЕ (18.3)



### - - -# - -## -#

### - - -# - -## -#



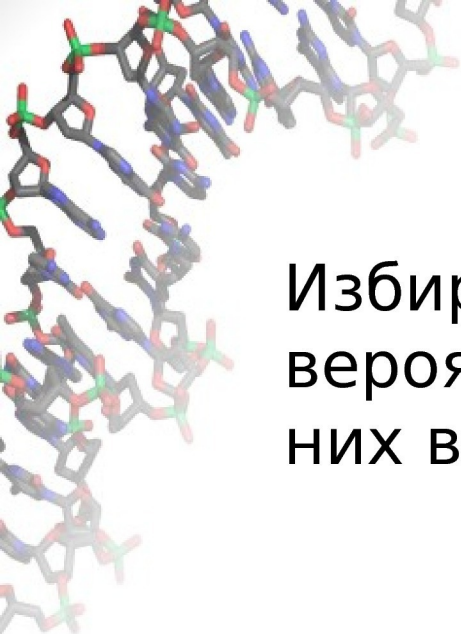
### - ## - - -# - ###

{ ## - # - #####  
 ##### - - -# - -## - #

**w=7**

{ ## - ## - #####  
 ##### - ##### - -##  
 ##### - ## - - -# - #####  
 ## - - - - ##### - #####  
 ##### - - -# - # - ## - ##  
 ##### - # - # - # - - - - - #####

**w=9**



## Пример: ВСЕ (18.3). Избирательности

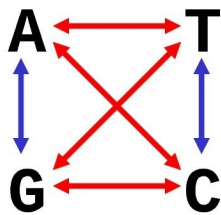
Избирательность семейства затравок –  
вероятность встретить хотя бы одну из  
них в случайном месте ( $p(\text{match}) = 1/4$ )

####	$w=4$	$\sim 39. \cdot 10^{-4}$
###-##	$w=5$	$\sim 9.8 \cdot 10^{-4}$
{ ##-#-#### ###-#-#-#-#	$w=7$	$\sim 1.2 \cdot 10^{-4}$
{ ##-##-##### ###-#####-## ###-##-#-### ##-#####-### ###-#-#-##-## ###-#-#-#-#####	$w=9$	$\sim 0.23 \cdot 10^{-4}$

## 3.4.2. ДА, НЕТ И МНОГОЕ ДРУГОЕ

Different mutational events have different probabilities

*transversions* ■



*transitions* ■

ATCAGTGC AATGCTCAAGA  
| | | | | · | | · | | | | : | | | |  
ATCAGCGCGATGCGCAAGA

*Transitions are usually over-represented.*



# *Extended seed alphabet*

- seed: ##@# - #@ - ###

{ '#': obligatory match position

{ '-' : joker position ("don't care" position)

'@' : ***transition-constrained*** position

position that corresponds to either a match or a transition.

##@# - #@ - ###

ATCAGTGC**A**ATGCT**T**CAAGA

|||||.|||.||||:|||||

ATCAGCGC**G**ATGCG**G**CAAGA



# *Subset letters and seeds*

- Seed **letter** is a subset of aligned pairs.
- # = {(A,A), (C,C), (G,G), (T,T)}
- @ = { (A,A), (C,C), (G,G), (T,T),  
(A,G), (G,A), (T,C), (C,T)}
- - = {all pairs}

##@# - #@ - ###

ATCAGTGCAATGCTCAAGA

|||||.|||.||||:||||

ATCAGCGCGATGCGCAAGA



## Затравки. Выводы.

- Классические заправки не оптимальны
- Применяя заправки, разберитесь, насколько они адекватны интересующему Вас классу сходств.

*О чем мы не говорили:*

- Как строить хорошие заправки?
- Как вычислить чувствительность заправки?
- И еще о многом...