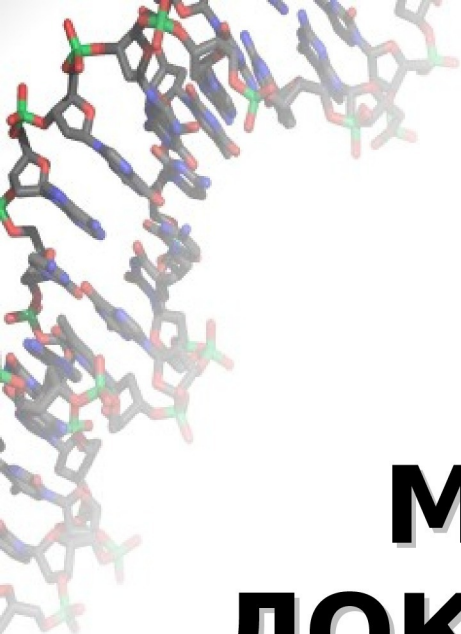




План занятия*

- 1. Множественное выравнивание и НММ
- 2. Поиск локальных множественных сходств.
- 3. MEME и Gibbs sampler
- 4. Иерархическое глобальное множественное выравнивание.



Часть 1

МНОЖЕСТВЕННЫЕ ЛОКАЛЬНЫЕ СХОДСТВА

.



Множественные локальные сходства.

- Дано: тексты S_1, \dots, S_n
- *Множественное локальное сходство* – это набор фрагментов этих последовательностей F_1, \dots, F_k таких, что фрагменты похожи.



Множественные локальные сходства. Нужны уточнения

- Дано: тексты S_1, \dots, S_n
- Множественное локальное сходство – это **набор фрагментов** ЭТИХ последовательностей F_1, \dots, F_k таких, что фрагменты **ПОХОЖИ**.



Множественные локальные сходства. Уточнение.

- Дано: S_1, \dots, S_n
- Множественное локальное сходство – это набор фрагментов этих последовательностей F_1, \dots, F_k таких, что фрагменты **ПОХОЖИ**.
- 1. Количество фрагментов в тексте:
 - ровно 1 в каждом,
 - не менее 1 в каждом,
 - сколько получится (= неточные повторы в объединенном тексте)



Множественные локальные сходства. Уточнение

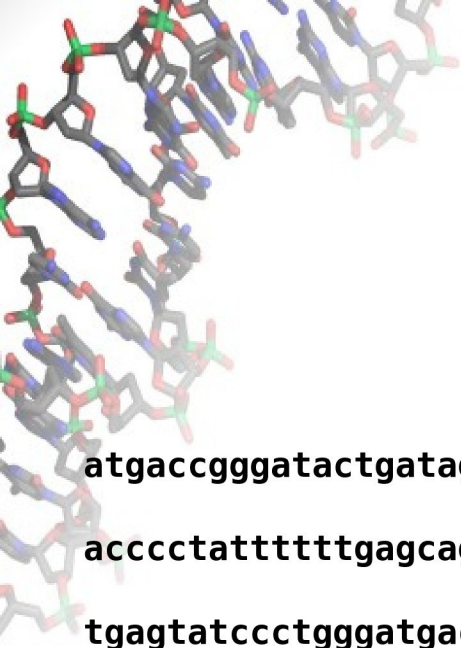
- Дано: S_1, \dots, S_n
- Множественное локальное сходство – это набор фрагментов этих последовательностей F_1, \dots, F_k таких, что фрагменты *похожи*.
- 2. Вид сходства:
 - допускаются ли делеции/вставки;
 - порог на качество сходств – парный или интегральный



Поиск неизвестного мотива The Motif Finding Problem

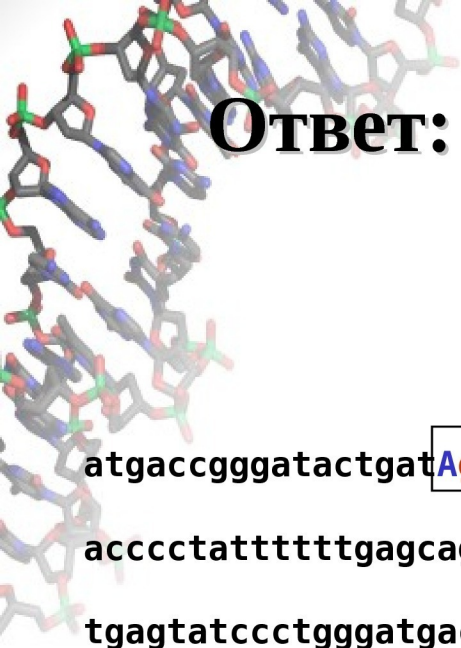
- Дано N последовательностей
Предположение: это случайные последовательности;
в каждую «встроен» мотив – данное слово с вариациями (без делеций)

```
cctgatagacgctatctggctatccacgtacgtaggtcctctgtgCGAATctatgCGTTTTccaacat  
agtactgggtgtacatttgatacgtacgtacaccggcaacctgaaacaacgctcagaaccagaagtgc  
aaacgtacgtgcaccctctttcttcgtggctctggccaacgagggctgatgtataagacgaaaatttt  
agcctccgatgtaagtcatagctgtaactattacctgccaccctattacatcttacgtacgtataca  
ctgttatacaacgcgctcatggcggggatgCGTTTTGGTCgtcgtacgctcgatCGTTAACGTACGTC
```



Пример. Где мотив? 😊

atgaccgggatactgatagaagaaagggtggggggtacacattagataaacgtatgaagtacgtagactcggcgccgcccg
accctatTTTTTgagcagatttagtgacctggaaaaaaaaatttgagtacaaaactTTTccgaatacaataaaacggcgggg
tgagtatccctgggatgacttaaataatggagtggtgctctccgattTTTgaatatgtaggatcattcgccagggtccga
gctgagaattggatgcaaaaaagggttgtccacgcaatcgcgaaccaacgcggacccaaaggcaagaccgataaaggaga
tccTTTTTgcggtaatgtgcccgggaggctggttacgtagggaagccctaacggacttaataataaaggaagggcttatag
gtcaatcatgttcttTgtgaatggatttaacaataagggctgggaccgcttggcgcacccaaattcagtgTgggcgagcgcaa
cggTTTTTggcccttTtagaggccccgtataaacaaggaggggccaattatgagagagctaattctatcgcgTgcgtgttcat
aacttgagTtaaaaaatagggagccctggggcacatacaagaggagtcttcttatcagTtaatgctgtatgacactatgta
ttggccattggctaaaagcccaacttgacaaatggaagatagaatccttgcatactaaaaaggagcggaccgaaaggggaag
ctggTgagcaacgacagattcttacgtgcattagctcgcttccggggatctaatagcacgaagcttactaaaaaggagcgga

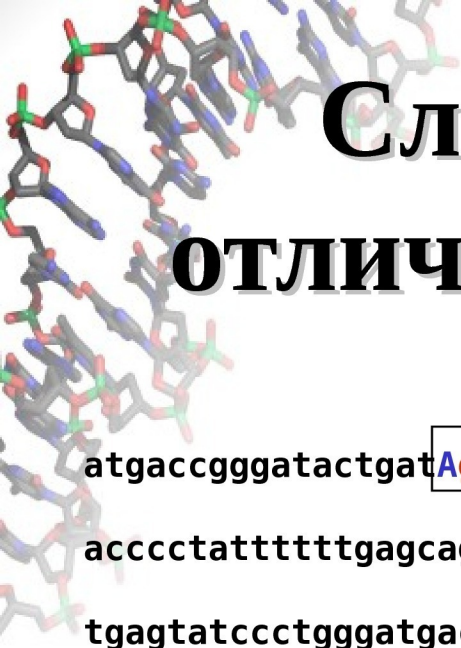


Ответ: **AAAAAAAAAGGGGGGG** (длина 15)

с 4 заменами

(15, 4)- МОТИВ

atgaccgggatactgat**AgAAgAAAGGttGGGgg**cgtagacattagataaacgtatgaagtacgtagactcggcgccgscg
accctatTTTTtgagcagatttagtgacctggaaaaaaaaaatttgagtacaaaactTTTccgaata**cAAtAAAAcGGcGGGa**
tgagtatccctgggatgact**AAAAtAAtGGaGtGGt**gctctccgattTTTtgaatatgtaggatcattcgccagggtccga
gctgagaattggatg**cAAAAAAGGGattGt**ccacgcaatcgcgaaccaacgcggaccsaaaggcaagaccgataaaggaga
tccTTTTgCGgtaatgtgcccgggaggctggttacgtagggaagccctaacggacttaat**AtAAtAAAGGaaGGG**cttatag
gtcaatcatgttcttTgtgaatggatt**AAcAAtAAGGGctGGg**accgcttggcgcaccsaaattcagtgTgggCGagCGcaa
cggTTTTggcccttgtagaggccccgT**AtAAAcAAGGaGGGc**caattatgagagagctaattctatcgcgTgcgTgttcat
aacttgagtt**AAAAAAtAGGGaGcc**ctggggcacatacaagaggagtcttcttatcagttaatgctgtatgacactatgta
ttggccattggctaaaagccsaaacttgacaaatggaagatagaatccttgcat**ActAAAAAGGaGcGGa**ccsaaagggaag
ctggTgagcaacgacagattcttacgtgcattagctcgcttccggggatctaatagcgaagctt**ActAAAAAGGaGcGGa**



Сложность; 2 экземпляра отличаются в $8 > 15/2$ позициях

atgaccgggatactgat**AgAAgAAAGGttGGG**ggcgctacacattagataaacgtatgaagtacgtttagactcggcgccgscg
 acccctatTTTTTgagcagatttagtgacctggaaaaaaatttgagtacaaaacttttccgaatac**AAtAAAAcGGcGGG**a
 tgagtatccctgggatgactt**AAAAtAAtGgaGtGG**tgctctccgatttttgaatatgtaggatcattcgccagggtccga
 gctgagaattggatg**cAAAAAAAGGGattG**tccacgcaatcgcgaaccaacgcggaccsaaaggcaagaccgataaaggaga
 tcccttttgcggtaatgtgcccgggaggctggttacgtagggaagccctaacggacttaat**AtAAtAAAGGaaGGG**cttatag
 gtcaatcatgttcttgtgaatggattt**AAcAAtAAGGGctGG**gaccgcttggcgcacccaaattcagtggtggcgagcgc
 cggttttggcccttgtagaggcccccg**AtAAAcAAGGaGGGc**caattatgagagagctaattctatcgcgtgctgttcat
 aacttgagtt**AAAAAAtAGGGaGcc**ctggggcacatacaagaggagtcttccttatcagttaatgctgtatgacactatgta
 ttggccattggctaaaagcccaacttgacaaatggaagatagaatccttgcat**ActAAAAAGGaGcGG**accgaaagggaag
 ctggtgagcaacgacagattcttacgtgcattagctcgcttccgggatctaatagcacgaagcctt**ActAAAAAGGaGcGGa**
 AgAAgAAAGGttGGG
 ..||..||..||..||..||..||..||..||
 cAAtAAAAcGGcGGG

Определения

- Будем задавать мотив вектором начальных позиций его вхождений в каждую из последовательностей

$$\mathbf{s} = (s_1, s_2, s_3, \dots, s_t)$$

Как оценить качество мотива?





Подготовка: (Частотный) Профиль и консенсус

Alignment

a	G	g	t	a	c	T	t
C	c	A	t	a	c	g	t
a	c	g	t	T	A	g	t
a	c	g	t	C	c	A	t
C	c	g	t	a	c	g	G

Profile

A	3	0	1	0	3	1	1	0
C	2	4	0	0	1	4	0	0
G	0	1	4	0	0	0	3	1
T	0	0	0	5	1	0	1	4

Consensus

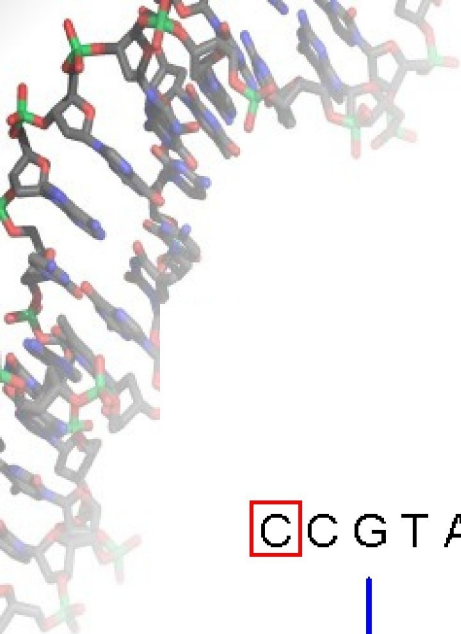
A C G T A C G T

- Расположим вхождения друг над другом
 $\mathbf{s} = (s_1, s_2, \dots, s_t)$
- Строим матрицу частот (частотный профиль)
- Консенсус-нуклеотид – наиболее частый нуклеотид в своей колонке



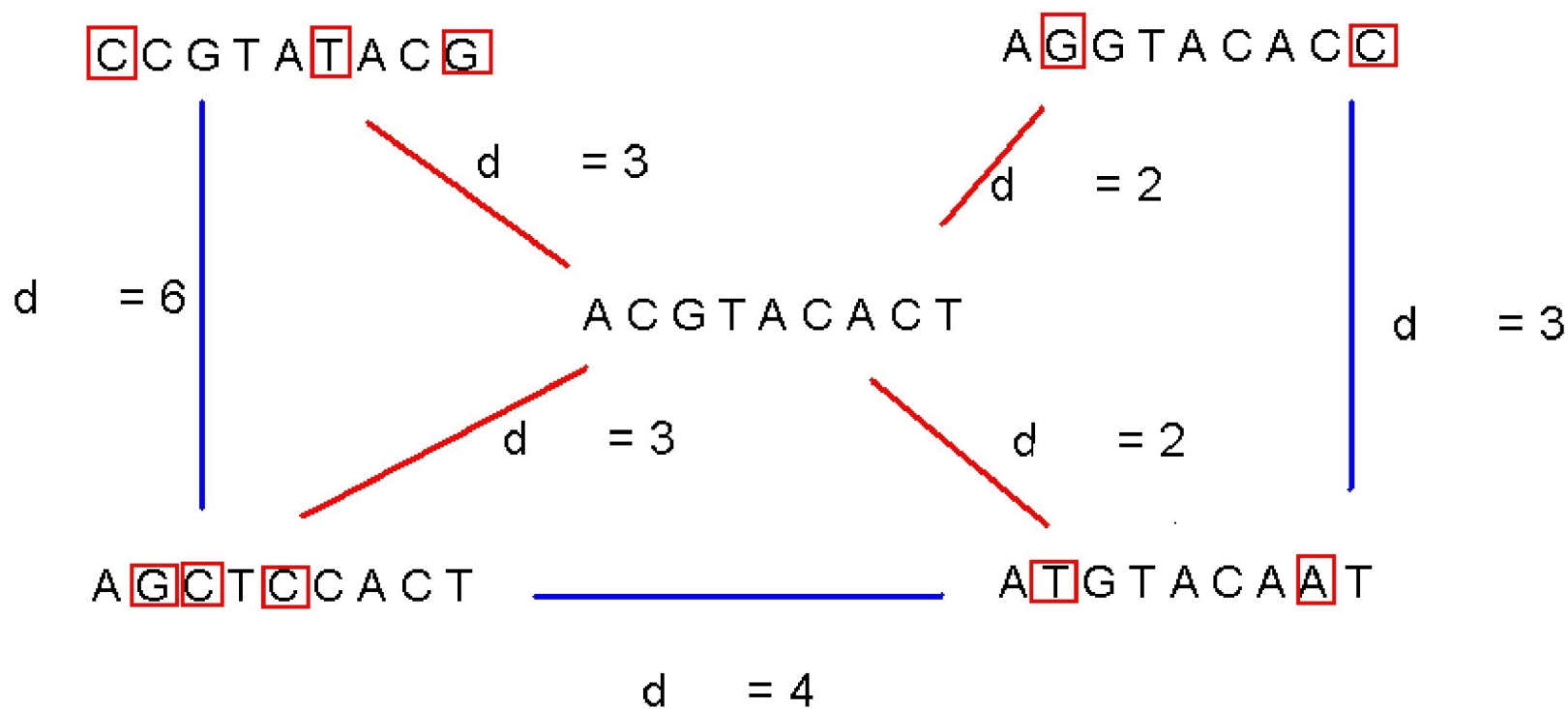
Консенсус

- Можно рассматривать консенсус, как предковый мотив
- Расстояние от консенсуса до каждого из вхождений *как правило* МЕНЬШЕ, чем расстояние между наиболее непохожими вхождениями мотива



Пример

$$d = 4$$

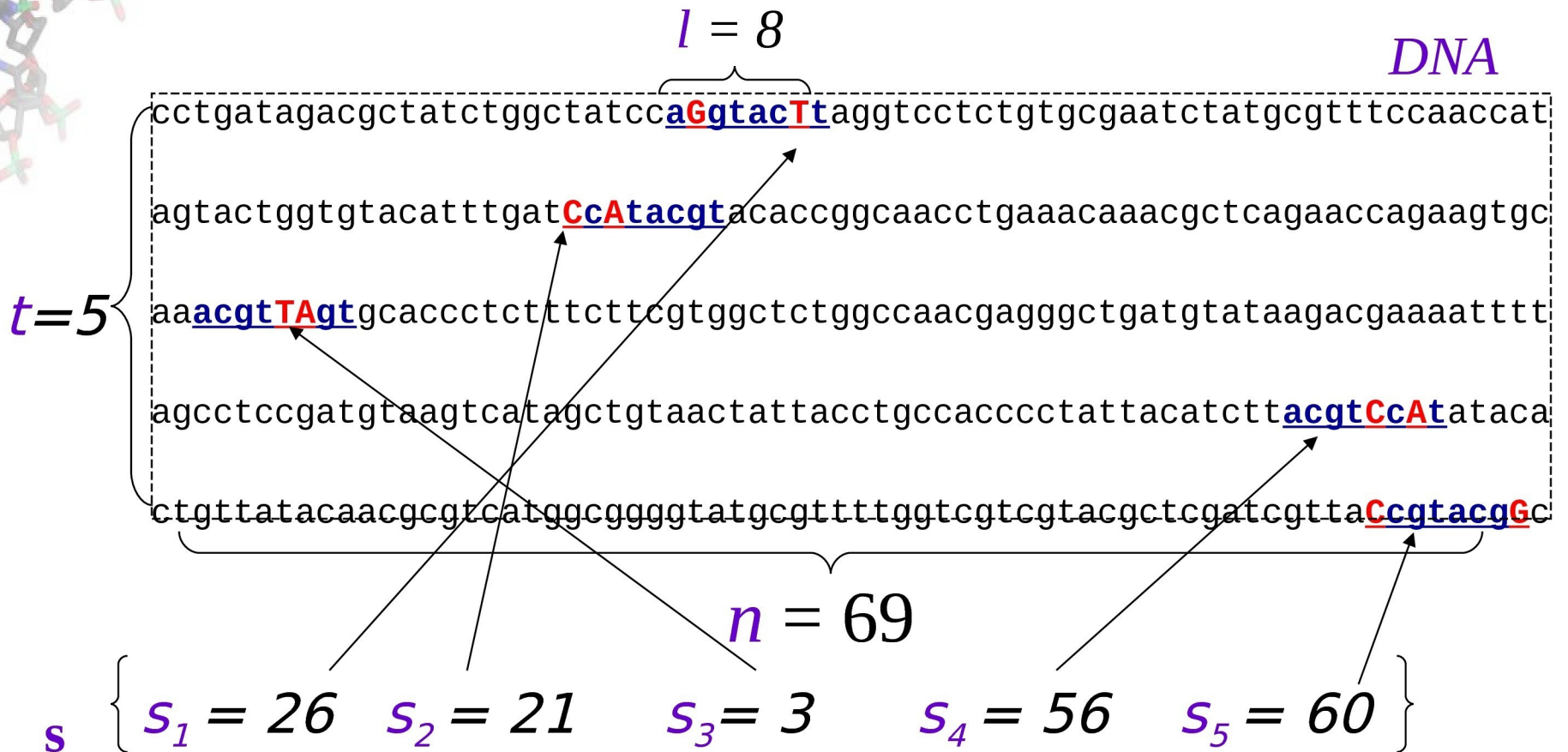
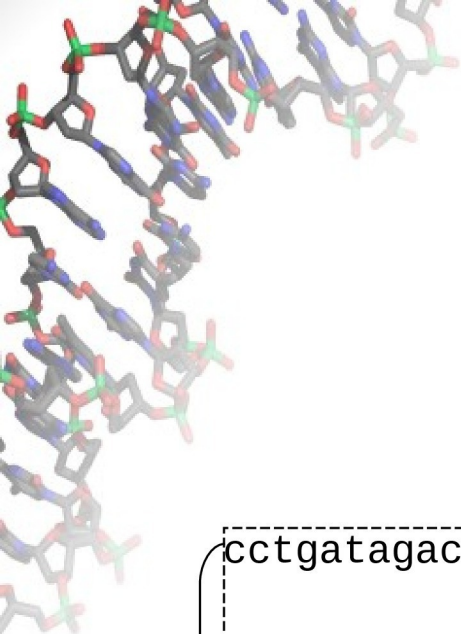


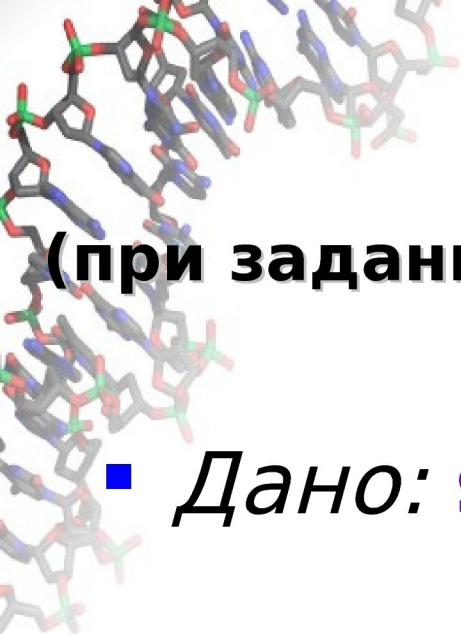


Качество мотива = качество консенсуса

- We have a guess about the consensus sequence, but how “good” is this consensus?
- Need to introduce a scoring function to compare different guesses and choose the “best” one.

Пример





Определение веса мотива

(при заданном наборе из t последовательностей)

■ Дано: $\mathbf{s} = (s_1, \dots, s_t)$:

	a	G	g	t	a	c	T	t	}
	C	c	A	t	a	c	g	t	
	a	c	g	t	T	A	g	t	
	a	c	g	t	C	c	A	t	
	C	c	g	t	a	c	g	G	

$$Score(\mathbf{s}) = \sum_{i=1}^l \max_{k \in \{A, T, C, G\}} count(k, i)$$

A	3	0	1	0	3	1	1	0
C	2	4	0	0	1	4	0	0
G	0	1	4	0	0	0	3	1
T	0	0	0	5	1	0	1	4

Consensus
Score

a c g t a c g t
3+4+4+5+3+4+3+4
=30



Постановка задачи (Motif Finding Problem)

- Дано: t последовательностей длины n ;
длина мотива L
- Найти: Вектор стартовых позиций $\mathbf{s} = (s_1, s_2, \dots, s_t)$, для которого $Score(\mathbf{s})$ будет максимальным

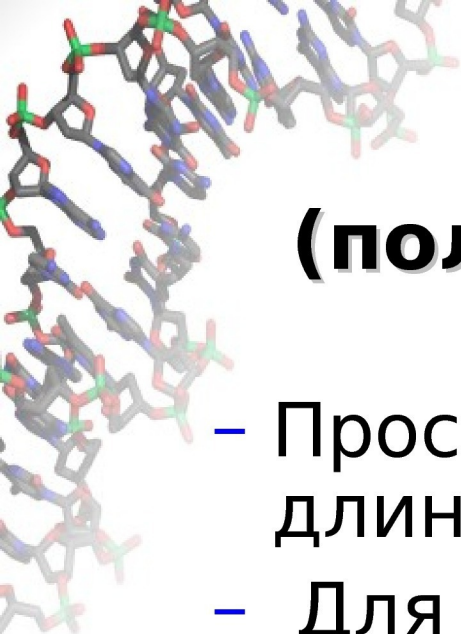


Решение 1 (полный перебор по данным)

– Перебираем все возможные вектора s

$$s_i = [1, \dots, n- \\ l+1]$$

$$i = [1, \dots, t]$$



Решение 2 (полный перебор по мотивам)

- Просматриваем все возможные слова длины L
- Для каждого слова w длины L
 - находим его наилучшее вхождение в каждой строке;
 - оцениваем полученный профиль



CONSENSUS (*Hertz, Stormo (1989)*): Greedy Motif Search

- Find two closest l -mers in sequences 1 and 2 and forms $2 \times l$ alignment matrix with $Score(\mathbf{s}, 2)$
- At each of the following $t-2$ iterations CONSENSUS finds a “best” l -mer in sequence i from the perspective of the already constructed $(i-1) \times l$ alignment matrix for the first $(i-1)$ sequences
- In other words, it finds an l -mer in sequence i maximizing $Score(\mathbf{s}, i)$
under the assumption that the first $(i-1)$ l -mers have been already chosen
- CONSENSUS sacrifices optimal solution for speed: in fact the bulk of the time is actually spent locating the first 2 l -mers



Some Motif Finding Programs

- **CONSENSUS**

Hertz, Stromo (1989)

- **GibbsDNA**

Lawrence et al (1993)

- **MEME**

Bailey, Elkan (1995)

- **RandomProjections**

Buhler, Tompa (2002)

- **MULTIPROFILER** *Keich, Pevzner (2002)*

- **MITRA**

Eskin, Pevzner (2002)

- **Pattern Branching**

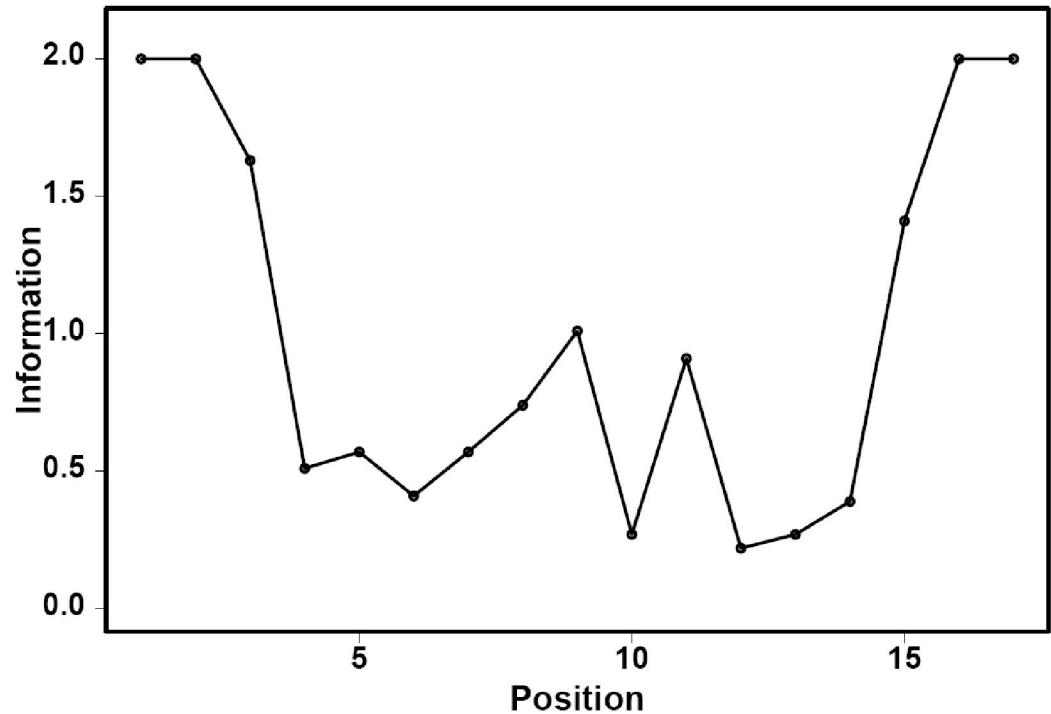
Price, Pevzner (2003)

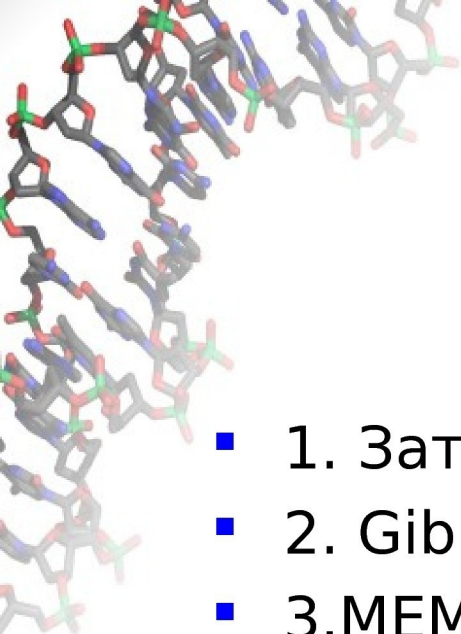
Gal4 Motif Information Content

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17
A	0	0	0	$\frac{7}{14}$	$\frac{1}{14}$	$\frac{1}{14}$	$\frac{9}{14}$	0	$\frac{6}{14}$	$\frac{1}{14}$	0	$\frac{3}{14}$	$\frac{1}{14}$	$\frac{2}{14}$	0	0	0
C	$\frac{14}{14}$	0	$\frac{1}{14}$	$\frac{3}{14}$	$\frac{3}{14}$	$\frac{6}{14}$	$\frac{3}{14}$	$\frac{8}{14}$	0	$\frac{5}{14}$	$\frac{3}{14}$	$\frac{7}{14}$	$\frac{5}{14}$	$\frac{3}{14}$	$\frac{12}{14}$	$\frac{14}{14}$	0
G	0	$\frac{14}{14}$	$\frac{13}{14}$	$\frac{4}{14}$	$\frac{9}{14}$	$\frac{6}{14}$	$\frac{1}{14}$	$\frac{5}{14}$	0	$\frac{6}{14}$	$\frac{1}{14}$	$\frac{2}{14}$	$\frac{6}{14}$	$\frac{1}{14}$	$\frac{2}{14}$	0	$\frac{14}{14}$
T	0	0	0	0	$\frac{1}{14}$	$\frac{1}{14}$	$\frac{1}{14}$	$\frac{1}{14}$	$\frac{8}{14}$	$\frac{2}{14}$	$\frac{10}{14}$	$\frac{2}{14}$	$\frac{2}{14}$	$\frac{8}{14}$	0	0	0

Information =

$$2 + \sum_{b \in \{A, C, G, T\}} p(b) \log_2 \frac{1}{p(b)}$$





Как искать МЛС

- 1. Затравки
- 2. Gibbs Sampler (случай 1 из каждого)
- 3. MEME



Затравки

- Так же, как и в парном случае.
- Общая идея: сокращение пространства поиска (фильтрация).
- Разреженные затравки и пр.



Gibbs Motif Sampling

- Gibbs Motif Sampler
 - Bayesian model, prior distribution
- Algorithm (MCMC)

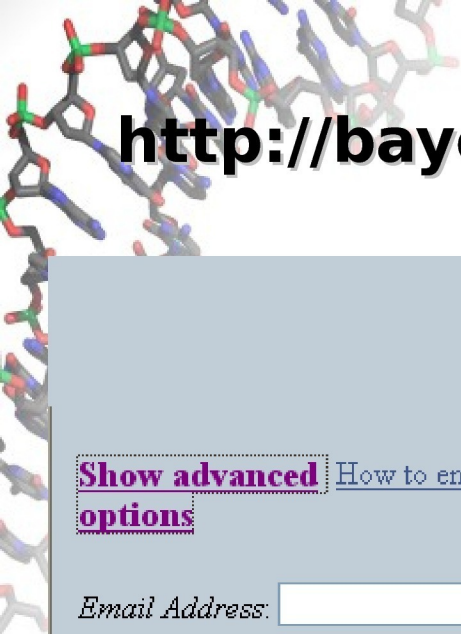
Initialization: Randomly select motif start-positions in each sequence

Iterations:

Remove randomly selected sequence k'

- Update frequency matrix
- Randomly select a motif start-position j for k' proportional to:

$$\frac{\text{Probability under motif model}}{\text{Probability under background model}} = \frac{p_i(b_{j+i-1}^k)}{p_0(b_{j+i-1}^k)}$$



Gibbs Motif Sampler

<http://bayesweb.wadsworth.org/gibbs/gibbs.htm>



The Gibbs Motif Sampler

(for DNA)

Show advanced options [How to enter data?](#)

Email Address:

Please enter the data sequence: ([FASTA format](#)) *

Prokaryotic Defaults

Sampler Mode: Site Sampler

No. of different motifs (patterns):

Motif Width(s):*

Eukaryotic Defaults

Motif Sampler Recursive Sampler

Max sites per seq: (recursive sampler)

Est. total sites for each motif type:

Expectation-Maximization (EM)

- MEME

Algorithm

- Missing data problem: Expectation-Maximization (EM) Algorithm to obtain maximum likelihood estimates

- EM Algorithm

Initialization: Set frequency matrix p and p_0

Iterations:

- E-step: Calculate probability of motif start-positions

For each sequence k and position j

$$W_{kj} = \Pr(\text{motif start-position} = j \mid p)$$

- M-step: Update frequency matrix estimate

$$\hat{p}_i(b) = \frac{\sum_k \sum_j W_{kj} \mathbb{1}(\text{sequence } k, \text{ position } j+i-1 = b)}{N}, \quad b = A, C, G, T$$


MEME

<http://meme.sdsc.edu/meme/website/meme.htm>


File Edit View Favorites Tools Help

Back Forward Stop Home Search Favorites Refresh Print Mail

Address <http://meme.sdsc.edu/meme/meme.html> Go Links David BBC CNN Yahoo Mail



Version 3.5.0

Hosted by 

Data Submission Form

Use this form to submit DNA or protein sequences to MEME. MEME will analyze your sequences for similarities among them and produce a description ([motif](#)) for each pattern it discovers. Your results will be sent to you by e-mail.

Your **e-mail address**:

Re-enter **e-mail address**:

[Optional] **Description** of your sequences:

Please enter the **sequences** which you believe share one or more motifs. The sequences may contain no more than **60,000 characters** total in any of a large number of **formats**.

- Enter the **name of a file** containing the sequences here:
- or the **actual sequences** here ([Sample Input Sequences](#)):

```
>YBR018C_176 433 649
TTGGTAAAGTAGAGGGGTAATTTTTCCCTTATTGTTTCATACATT
CTTAAATTGCTTTGCCCTCCCTTTGGAAAGCTATACTTCGGAGCACTG
TTGAGCGAAGGCTCATTAGATATATTTCTGTCATTTCCCTTAACCCAA
AAATAAGGGAAAGGGTCCAAAAAGCGCTCGGACAACGTTGACCGTAT
```

How do you think the occurrences of a single motif are **distributed** among the sequences?

One per sequence

Zero or one per sequence

[Optional] MEME will find the optimum **number of sites** for each motif within the limits you specify here: **Minimum sites** (>= 2)

MEME will find the optimum **width** of each motif within the limits you specify here: **Minimum width** (>= 2)

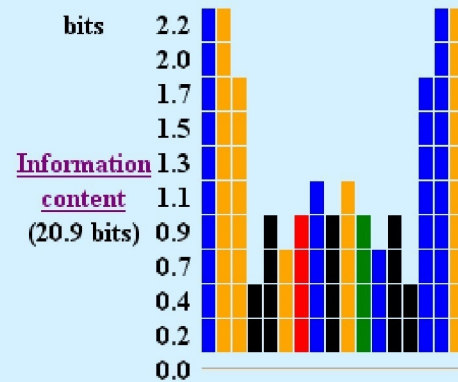
Internet

MEME Output

FN

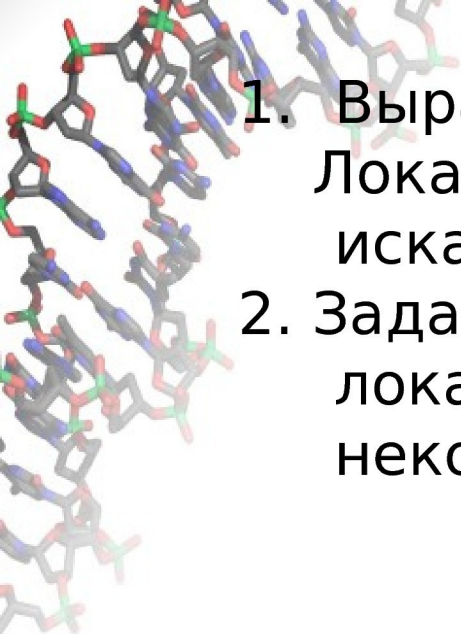
MOTIF 1 width = 17 sites = 14 llr = 203 E-value = 2.7e-025

Simplified A : : : 5: 1615: : 1: 1: : :
pos.-specific C a: 1152: 7: 246549a:
probability G : a945642: 7: 2511: a
matrix T : : : 1: 1: : 5161: 5: : :



Multilevel C G G A G G A C T G T C C C T C C G
consensus G C C G A C G G C
sequence

NAME	START	P-VALUE	SITES
YBR020W_74	93	1.11e-09	AGCCGCCGAG C G G G C G A C A G C C C T C C G AC G G A A G A C T
YBR020W_153	153	1.11e-09	AGCCGCCGAG C G G G C G A C A G C C C T C C G AC G G A A G A C T
YBR020W_154	136	1.11e-09	AGCCGCCGAG C G G G C G A C A G C C C T C C G AC G G A A G A C T
YBR019C_94	176	1.44e-09	AGTCTTCGGT C G G A G G G C T G T C G C C C G C T C G G C G G C T

- 
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

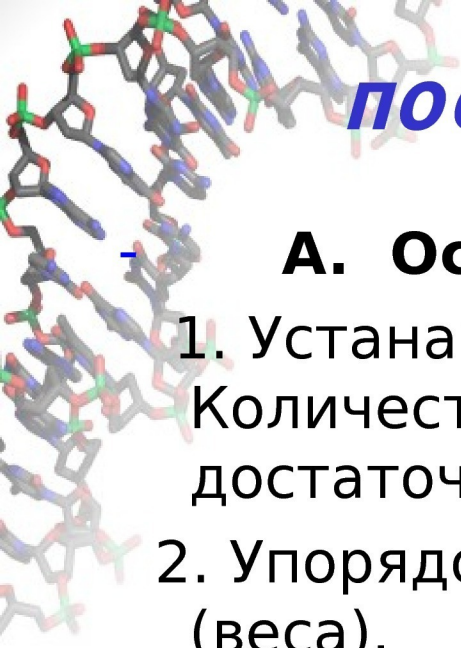
РЕШЕНИЕ:

Иерархическое выравнивание



Иерархическое выравнивание

- отказ от глобальной оптимизации какой-либо весовой весовой функции;
- разрешение конфликта между сходствами производится локально
- иерархическая процедура



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

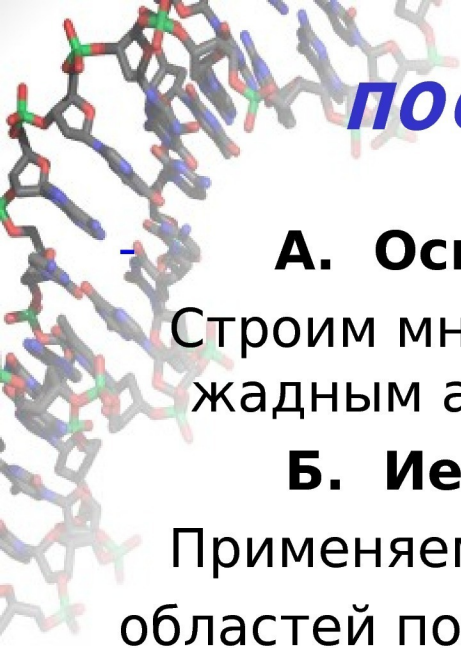
Применяем основной шаг отдельно к каждой из областей поиска между отобранными локальными сходствами.

!!! 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 «остовную цепочку».

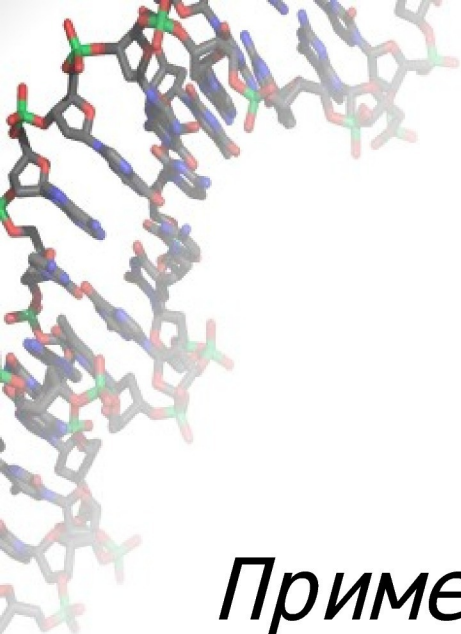
Б. Иерархический переход

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

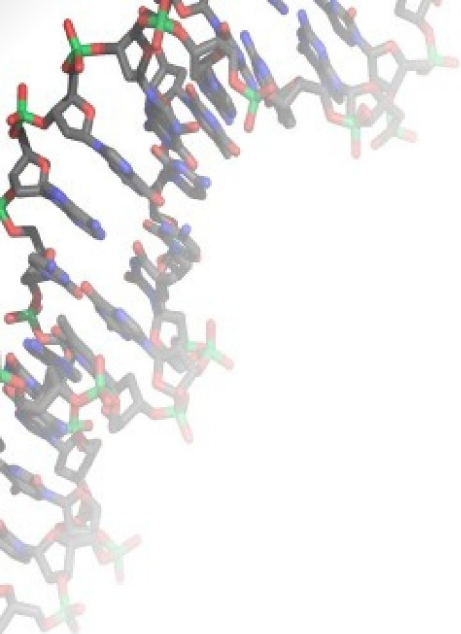
В. Итерирование

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

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



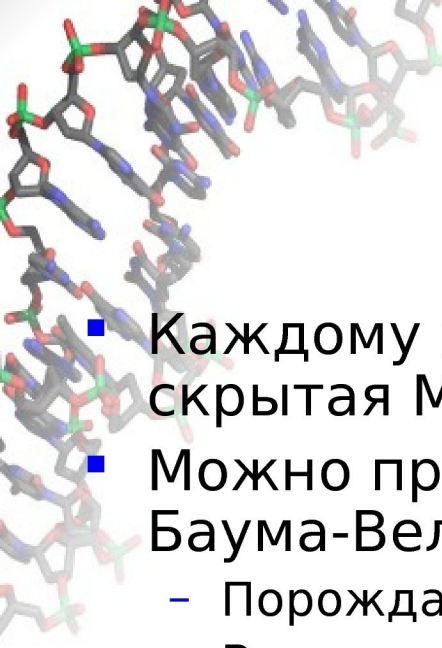
*Пример: иерархическое геномное
выравнивание (длина $\sim 10^7$)*



3. Вес множественного выравнивания - сумма весов столбцов

*++++**+
- ИВАН - - - ОВ - - - -
- ИВАН - - ЦОВ - - - -
- ИВАН - - КОВ - - - -
- ИВАНЧУКОВ - - - -
ДИВАН - - - ОВ - - - -
ДИВАНЧИКОВ - - - -
- ИВАН - - - ОВСКИЙ
- - ВАН - ЪКОВ - - - -

Недостаток: не учитываем, что «работать» можно не с символами, а с фрагментами.
Вариант: выделяем «значимые» столбцы; все остальное учитываем отдельно в каждой последовательности.



Множественное выравнивание с помощью НММ

- Каждому множественное выравнивание соответствует скрытая Марковская модель.
- Можно применить алгоритм максимизации ожидания Баума-Велча:
 - Порождаем случайные параметры НММ.
 - Выравниваем все последовательности с этой моделью
 - Переоцениваем параметры.
- Проблема: легко попасть в локальный максимум
- Обход проблемы: время от времени параметры НММ возмущаются.
- Другой вариант – использование искусственного отжига.
- Достоинство подхода: одновременно анализируются все последовательности. Нет проблемы необратимости, характерной для прогрессивного выравнивания.