



Обзор задач биоинформатики, связанных с анализом и обработкой текстовых последовательностей

Орлов Ю.Л.

Кафедра информационной биологии ФЕННГУ

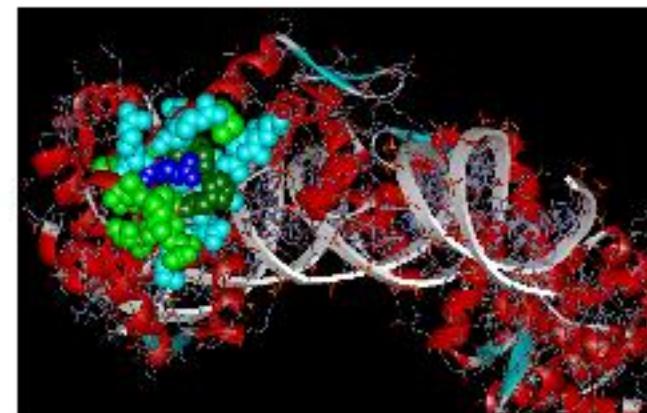
Обзор задач биоинформатики как задач компьютерного анализа генетических текстов.

Проблемы компьютерного анализа генетических текстов

Распознавание и функциональная аннотация регуляторных последовательностей генов эукариот имеет большое значение для молекулярной биологии в наступившую "пост-геномную" эру.

До начала эпохи массового секвенирования многим исследователям казалось, что функциональные участки будут закодированы однозначными последовательностями, скорее всего изменяющими локальные физические свойства ДНК. В действительности, проблема определения функции по последовательности ДНК гораздо сложнее, что связано с неоднозначностью кодирования генетической информации.

Пример – сайты рестриктаз и сайты связывания ТФ.



Рекомендуемая литература:

David W. Mount «Bioinformatics. Sequence and genome analysis»

«Компьютерный анализ генетических текстов» (Ред. Франк-Каменецкий), 1990

Основные задачи компьютерного анализа генетических текстов:

- 1) поиск гомологии и выравнивание генетических текстов, множественное выравнивание**
- 2) статистический анализ генетических текстов, исследование структуры повторов и модели порождения символьных последовательностей, сегментация геномов;**
- 3) предсказание кодирующих участков генов и открытых рамок считываания;**
- 4) предсказание функциональных сигналов (функциональных сайтов и регуляторных районов);**
- 5) анализ вторичной структуры РНК и сигналов трансляции;**
- 6) анализ аминокислотных последовательностей белков, предсказание вторичной структуры, функциональных сайтов и доменов глобуллярных белков по их аминокислотным последовательностям;**
- 7) филогенетические сравнения;**
- 8) ДНК-чибы – экспрессионные кривые**
- 9) задачи оперирования с большими массивами информации и управления (Интернет-навигации) разрозненными специализированными базами данных**

Основные задачи компьютерного анализа генетических текстов:

1) поиск гомологии и выравнивание генетических текстов, множественное выравнивание

Рассматриваемые вопросы

- Дот-матрица или метод диаграмм для сравнения последовательностей
- Выравнивание последовательностей с помощью динамического программирования
- Поиск локального выравнивания последовательностей.
- Множественное выравнивание последовательностей.
- Поиск гомологии в базах данных. Методы FASTA и BLAST для поиска в базах данных.

Дот-матрица или метод диаграмм для сравнения последовательностей

1970 - A.J.Gibbs and G.A.McIntyre

Сравнение последовательностей.

Точечная матрица

совпадений –

dot-matrix

(Дот-матрица).

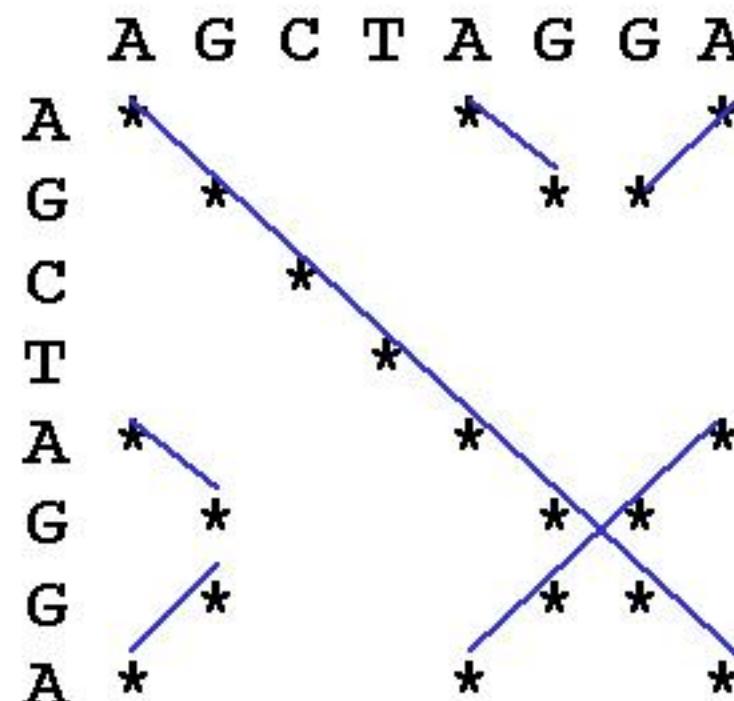
	A	G	C	T	A	G	G	A
G	*				*	*		
A	*				*			*
C			*					
T				*				
A	*				*			*
G		*				*	*	
G		*				*	*	
C			*					

Рассмотрим две последовательности: AGCTAGGA и GACTAGGC. Диагональ точек выявляет общую подпоследовательность CTAGG

Рассмотрим одну последовательность AGCTAGGA.

В этом случае матрица всегда симметрична.

Диагонали точек выявляют прямые и симметричные повторы



Короткий повтор: AG

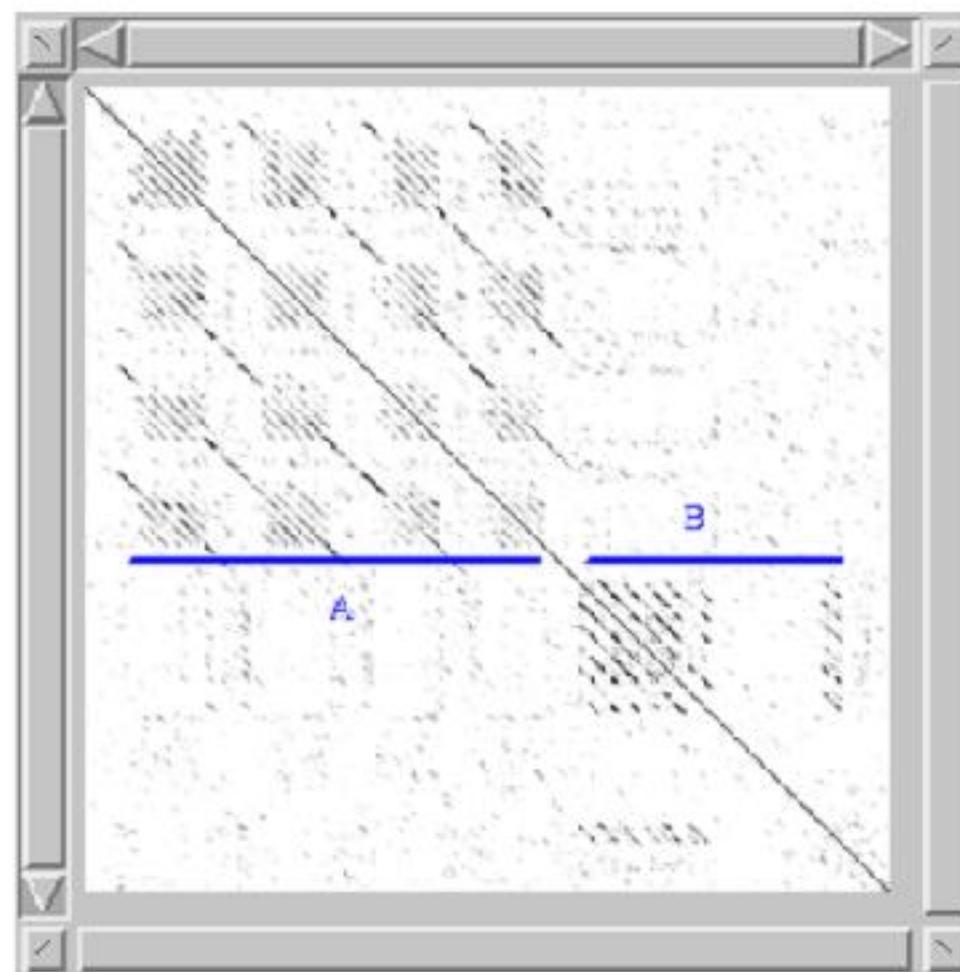
Симметричный повтор: AGGA

Пример произвольной символьной последовательности.



Матрица позволяет находить общие подслова в последовательностях символов любого алфавита – нуклеотидном, аминокислотном, на естественных языках.

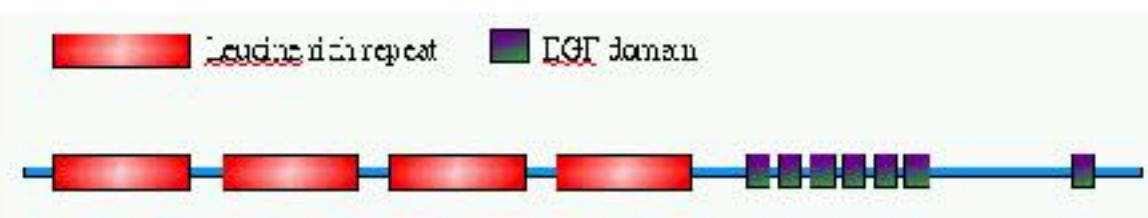
Усовершенствование сравнения.



На данном примере представлен график (точечная матрица) для белка SPLIT *D.melanogaster*. Видны 4 повтора в N-концевом домене и (А) и 6 в С-концевом (В)



Результат поиска повторов с помощью точечной матрицы гомологии. Белок SPLIT *D.melanogaster*. Аминокислотная последовательность и структура повторов.



- SPLIT_DROME (P24L14) :

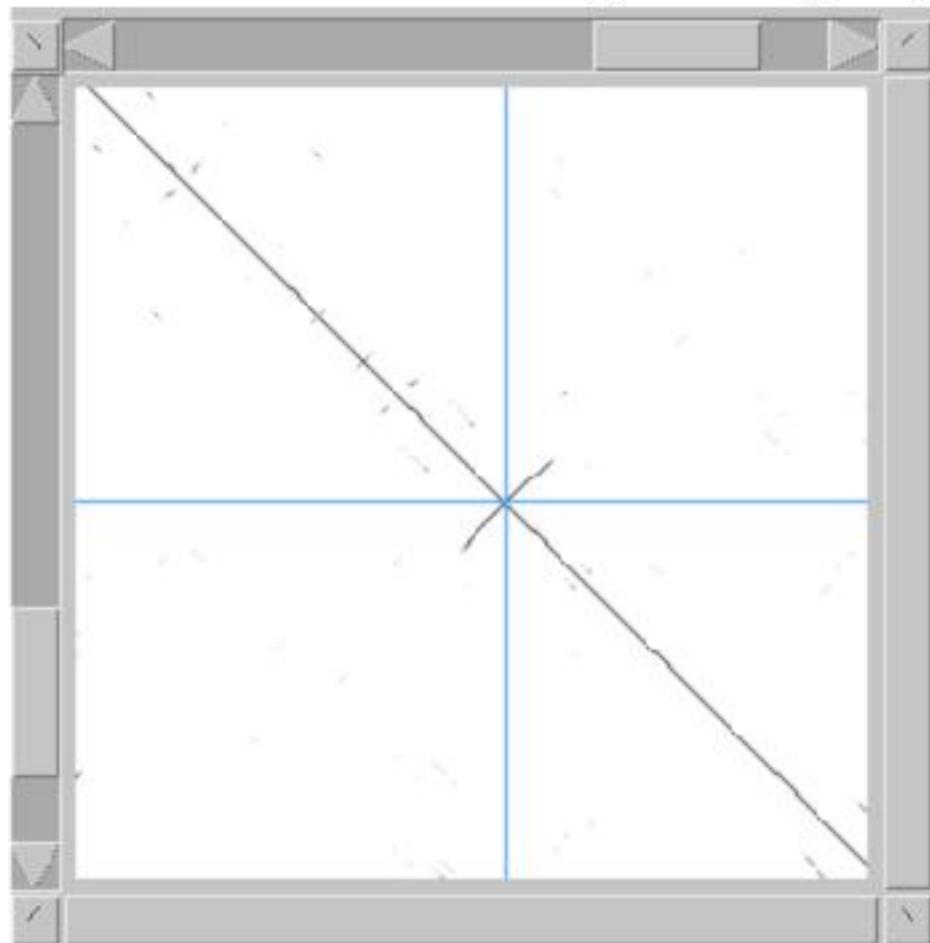
```
MAAFSRTTLKPPFFRLQLLELLPILLLRDAVHAEPYEGGFG32A723GGLG2VGHIHIFGGGVGVCTEARCPRVCSCT  
GLNVDQSERHGLTSEVPEKISADVERLEQGNKLTVIYETDFOQLTKLRLMQLTDNQCHTEEENSFQDLVSLERLDISMMV  
TTRARRUFKQIAQELRELQIDNKGQTTCATLGEHRKQIATLTTTMMMMTTSIFHNTFGGLGETRARIISNPTDQCSHISPI  
LQRFLRCATHLAFTTYCQCPSQLKCQNAQDIIQDFEKCQCLTEIIAEPECGAENGSPRIFCRCADGCDVDCREKOLTGVPVTL  
FDOTTOVRLQEKFTELPPESFSSFRRLRRIDLAMNNMSEKEDALSGLEQLTIVLVYGNKIKILP2GVFEGLGSLRLLL  
TMMKETPCTEKDPRFLIHEIISTIATIITQSIATNGTETRANKEPKTVAHAKNPFTTGSQMT  
SPKEMHIRRRIEGLREEKFKCSUGELRMKLOGDCRMDCDCFAMCPSCEGTTVQCTGRRLEIFRDIPLIETTELLNDNLGR  
ISSDGLFGRPLPHLWALELERNQLTGEEPMNAFEGASHIQELQLGENKIZIEISNMFLGLHQLETLNLYVDNQISCVMFGSFE  
FIMPLTSMIATASNPTRKNCMHTAPTFECURKKSINGGAIRGAGSKVRDUQTKDIFPNSFFKQSSMSFEGI GQGTCFEPSC  
CTGTUVVACSENQLKEIPRGCPACTCCLYLEEONIEQIHYERIRIILSOLTRDLSNNQITILQNYTPANLTKLSTLISYN  
KLUUQKHALSGLMNLWVUSLEGNAKLSLPEGSMEDLKSLVHILLSNPLYLCDCGLKAFSTIUKLQYVEFGLANCAEPQ  
MEKDLILSTF332FVCRGEVRNDILAKXHACFEGDODNOAOOCVALFOREYDOLCOOPGYHGKHCEFMIDACTGPNPCRMAT  
CTYLEECSRFSQCAFCYTCACTEWDDCLCEIEKOMNATCICDCCESVHCECQDCFSCEPCDTKICQFCSTEFNTPCANCAC  
LDEFTKTYSCDQYQIGEHDINUTDNIDWQNEPDQMGFIVUDGIMDQSCPBDYTGKCYCQGHNWLSKDVQTCRQDNEC  
KEGIVCFOPNAOGDYLCRCHPGTCKYCEYLTSISFVENNSFVELEPLERTRFEANVTIVF23AEONGCMTYCGODAHAV  
ELFNCRIRVSEYDQVNERVSTHYSPEHVADQKTHAVELLAIKKNFTLIVDRLARSQINECSDVYLKLTTTIPLOGLPVR  
AQOQAYAKUQ1BNL175FEGKQKAKVWV1HHLV7FGR1QKQK_TFGCAL_EGEGQQEEEDIEEQFMEETBKEEFVDFC_LEN  
KCREG3RCUEN3KARDGYCQCKKHCORGRTYDOGEGRTEPPVTAASTCRKEOVREYTTENDCRSROPLETAKC7GGCGN  
QCCAALKIVQPERKVRMVC8MNRKTYIHLDEYVRCOSCCTKCY
```

Пример взят: <http://www.isrec.isb-sib.ch/java/dotlet/repeats.html>

Поиск гомологии. Дот-матрица

Точечная матрица гомологий. Представлен инвертированный повтор (X-форма)

Bacillus subtilis UTP-glucose-1-phosphate uridylyltransferase gene



BSGTABX | 1018

ACTAAACAAAGAAGAAATCTAAAACAAAAAGGCTATTTGACATTCTATCCATAACCCCTTTTATTTCACACATCAAGGTCAATGT
TACTAAACAAAGAAGAAATCTAAACAAAAAGGCTATTTGACATTCTATCCATAACCCCTTTTATTTCACACATCAAGGTCAATGT

3SGTABX | 1017

BSGTABX (revcomp'd) | 1018

TACATTGACTTTGATGTTGAAATAAAGGGCTAATTGGATGAATGTCCAAATGCCCTTTTGTAGATTTCTTCCTTTGTTAG
TACTAACAAAGAAGAAATGCTAACAAAGGCTAATTGGACATTCAACCAATGCCCTTTTATTTCAACATCAGTCATG

SSGTABX | 101 /

Точечная матрица гомологий.

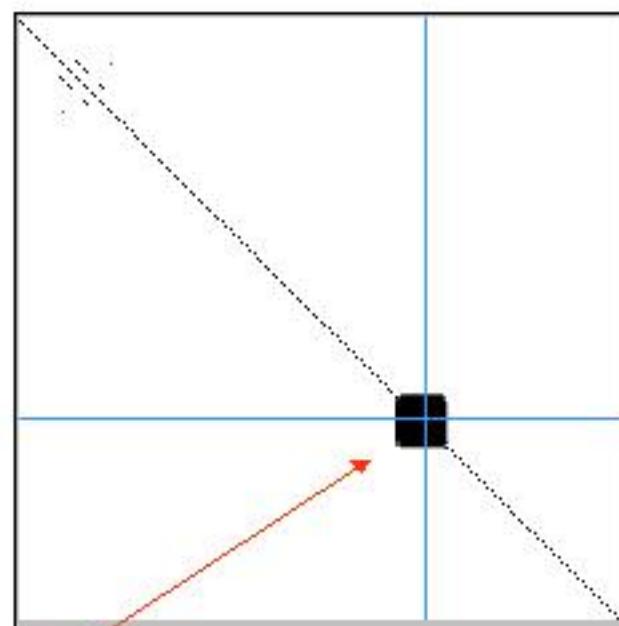
Участки низкой сложности

Белок SERA PLAFG (P13823)

маллярийного плазмодия

содержит серин-богатый участок

Расшифровать структуру квадрата



- SERA PLAFG (P13823) :

Выравнивание последовательностей с помощью динамического программирования

Выравнивание последовательностей с помощью динамического программирования. По этапам

Рассмотрим две последовательности: GATCTA и GATCA

Needleman-Wunsch	G	A	T	C	T	A	G	A	T	C	T	A
Алгоритм Нидльмана-Вунша	G	1	0	0	0	0	G	1	0	0	0	0
	A	0	1	0	0	1	A	0	2	0	0	1
	T	0	0	1	0	1	T	0	0	3	0	1
	C	0	0	0	1	0	C	0	0	0	4	0
	A	0	1	0	0	0	A	0	1	0	0	5

Итоговое выравнивание двух последовательностей:

GATCTA

GATC-A

Поиск гомологии. Выравнивание последовательностей

Пример – выравнивание слов COELACANTH и PELICAN с использованием простой схемы подсчета: +1 для совпадающих букв (match), -1 для несовпадений (mismatches), и -1 для пробелов (gaps).

COELACANTH

P-ELICAN--

COELACANTH

-PELICAN--

Пример матрицы
выравнивания

The diagram illustrates the sequence alignment between "COELACANTH" and "PELICAN". Above the sequences, the words are written in red. Below them, the aligned sequences are shown: "COELACANTH" and "-PELICAN--". A small window titled "Scoring Matrix" shows the scoring of each character pair. The matrix has "P" (PELICAN) on the top row and "E" (ELICAN) on the left column. The diagonal from top-left to bottom-right shows matches ('C', 'O', 'E', 'L', 'A', 'C', 'A', 'N', 'T', 'H'). Off-diagonal entries include 'P' at (1,2), 'E' at (2,4), 'L' at (3,6), 'I' at (4,7), 'C' at (5,8), 'A' at (6,9), and '-' at (7,10). The matrix is a 10x10 grid.

	C	O	E	L	A	C	A	N	T	H
P	C	0								
E			E							
L				L						
I					A					
C						C				
A							A			
N								N	T	H
								N	-	-

Поиск гомологии. Выравнивание последовательностей

Инициализация матрицы
сравнения

	C	O	E	L	A	C	A	N	T	H
0	-1	-2	-3	-4	-5	-6	-7	-8	-9	-10
P	↑ -1									
E										
L										
I										
C										
A										
N										

Начало заполнения матрицы
сравнения

CO CO
-P P-

	C	O	E	L	A	C	A	N	T	H
0	-1	-2	-3	-4	-5	-6	-7	-8	-9	-10
P	↑ -1	↑ -1								

Заполнение второй строки
матрицы сравнения

	C	O	E	L	A	C	A	N	T	H
0	-1	-2	-3	-4	-5	-6	-7	-8	-9	-10
P	↑ -1	↑ -1	↑ -2							

Заполненная матрица сравнения для тестового примера

	C	O	E	L	A	C	A	N	T	H	
P	0	-1	-2	-3	-4	-5	-6	-7	-8	-9	-10
E	-1	-1	-2	-3	-4	-5	-6	-7	-8	-9	-10
L	-2	-2	-2	-1	-0	-3	-4	-5	-6	-7	-8
I	-3	-3	-3	-2	-2	-1	-2	-3	-4	-5	-6
C	-4	-4	-4	-3	-1	-1	-2	-1	-4	-5	-6
A	-5	-3	-4	-4	-2	-2	-0	-1	-2	-3	-4
N	-6	-4	-4	-5	-3	-1	-1	-1	-0	-1	-2
	-7	-5	-5	-5	-4	-2	-2	-0	-2	-1	-0

COELACANTH
-PELICAN-

Local Alignment: Smith-Waterman

Алгоритм Смита-Уотермана является модификацией алгоритма Нидльмана-Вуниша

Отличия

- Углы матрицы инициализируются 0 вместо увеличения штрафа за пробел.
- Максимальный счет никогда не бывает меньше нуля, и указатель не записывается до тех пор пока счет больше нуля.
- Обратный проход начинается с наибольшего счета в матрице (а не с конца) и заканчивается при счете 0 (а не в начале матрицы).

ELACAN

ELICAN

		C	O	E	L	A	C	A	N	T	H
P	0	0	0	0	0	0	0	0	0	0	0
	0	0	0	0	0	0	0	0	0	0	0
E	0	0	0	1	0	0	0	0	0	0	0
	0	0	0	0	2	1	0	0	0	0	0
L	0	0	0	0	1	2	1	0	0	0	0
	0	0	0	0	1	1	0	0	0	0	0
I	0	0	0	0	1	1	0	0	0	0	0
	0	0	0	0	1	1	0	0	0	0	0
C	0	1	0	0	0	0	2	0	0	0	0
	0	0	0	0	0	1	0	3	2	1	0
A	0	0	0	0	0	1	0	3	2	1	0
	0	0	0	0	0	0	0	1	4	3	2
N	0	0	0	0	0	0	0	1	4	3	2
	0	0	0	0	0	0	0	1	4	3	2

Поиск локального выравнивания последовательностей

1981 - Mike Waterman, Temple Smith

Наиболее биологически значимые районы в ДНК и белках – локальные районы, которые выравниваются хорошо, в то время как остающиеся участки менее значимы.

Две меры оценки выравнивания (подсчета, скора – score) – счет сходства (близости) последовательностей и счет расстояния между ними.

Smith T.F. and Waterman M.S. (1981) Identification of common molecular subsequences. *J.Mol.Biol.* **147**: 195-197.

Значимость счета выравнивания

Karlin S. and Altschul S.F. (1990) Methods for assessing the statistical significance of molecular sequence features by using general scoring schemes. *Proc. Natl. Acad. Sci. USA* **87**:2264-2268.

Множественное выравнивание последовательностей

Johnson and Doolittle, 1986

GCG PILEUP

CLUSTALW (Thompson et al., 1994)

(Baylor College of Medicine,

<http://dot.imgen.bcm.tmc.edu:9331/multi-align/multi-align.html>)

Gribskov et al., 1987

A00825_hsp17_L.Merr.	TTACAATCTCCCTAGTTTC---TAATCTCAGCTAAGAAAA-
ACCAAAAGA	
A00826_hsp17.5-M_L.Merr.	TTTGATCTCCAAGTTTC---
AAATCTCGCGAATAAAATATATCAAAGA	
A00662_mult.r.g.1b_Rat	GCCGCTGCTCCCATCTTC----
GAGGCTCAGCTCAACTCAGAGCTACTT-	
A00172_mult.r.g.1b_Mus	GCCGCTGCTTCCATCTTCT---
GAGGTTCCGCTCAACTCAGAGCTACT--	
A00597_M35021_70hsp_Mus	TCCAGAG---ACAAGCGAA---GACAAGAGAAGCAGAGC-
GAGCGGGCGC-	
A00597_M76613_70hsp_Mus	TCCAGAG---ACAAGCGAA---
GACAAGAGAACGAGCAGAGCAGAGCGGGCGC-	
A00551_heme_ox.1_Homo	GCACGAA----CGAGCCC----
AGCACCGGCCGGATGGAGCGTCCGCAA	
A00569_heme_ox.1_Rat	GCCGGAG----CAGAGCC----
ATCTCGAGCGGAGCCCGGAGCCTGAAG	
A00552_hsp70_Xenopus	TACTTAC----T

Поиск гомологии в базах данных. Методы FASTA и BLAST для поиска в базах данных

FASTA (Pearson and Lipman, 1988)

BLAST (Altschul et al., 1990)

www.ncbi.nlm.nih.gov/BLAST

Gapped-BLAST (в три раза быстрее)

PSI-BLAST (большая чувствительность)

Основная идея – поиск коротких полностью совпадающих фрагментов и расширение выравнивания

Поиск по базам данных нуклеотидных и белковых последовательностей с помощью программы BLAST

Основная задача, которая решается с помощью программы **BLAST** (Basic Local Alignment Search Tool) - быстрый поиск сходных (гомологичных) участков последовательностей в банках данных нуклеотидных или аминокислотных последовательностей. **BLAST не является** программой выравнивания.

Входными данными для программы **BLAST** является нуклеотидная или аминокислотная последовательность, для которой производится поиск сходных последовательностей среди всех последовательностей из банков данных. Выходными данными являются последовательности из банков данных, которые имеют участки статистически значимого сходства с тестируемой последовательностью.

The screenshot shows the NCBI BLAST search interface. At the top, there's a blue header bar with the NCBI logo and links for Genomic Biology, Human Genome Guide, and Human Sequence. Below the header is a search bar with dropdown menus for 'Search' (set to 'LocusLink') and 'for', and a 'Go' button. To the left, there's a sidebar with links for 'BLAST overview', 'FAQs', 'news', 'manual', and 'references'. The main content area has a title 'BLAST the Human genome' and a subtitle 'Compare your query sequence to the working draft sequence of the human genome or its mRNA and protein products.' Below this are two dropdown menus: 'Database' set to 'genome' and 'Program' set to 'blastn'. There's also a checked checkbox for 'use MegaBLAST' and a large 'Begin Search' button.



PubMed

Info

- FAQs
- News
- References
- Credits

Education

- Program selection guide
- Tutorials
- URI API guide

Download

- Executables
- Databases
- Source code

Support

- Helpdesk
- Mailing list

Entrez

NEW 15 August 2003 We will be reorganizing the executable directory of our FTP site.
[Read more...](#)

BLAST

OMIM

Taxonomy

Structure

Nucleotide

- Discontiguous megablast
- Megablast
- Nucleotide nucleotide BLAST (blastn)
- Search for short nearly exact matches
- Search trace archives with megablast or (dis)align in megatext

Protein

- Protein protein BLAST (blastp)
- PHI and PSI BLAST
- Search for short nearly exact matches
- Search the conserved domain database (cddblast)
- Search by domain and lecture (blastx)

Translated

- Translated query vs. protein database (blastp)
- Protein query vs. translated database (blastn)
- Translated query vs. translated database (tblastx)

Genomes

- Human, mouse, rat
- Fugu rubripes, zebrafish
- Flies, nematodes, plants, yeasts, malaria
- Microbial genomes & other eukaryotic genomes

Special

- Align two sequences (bl2seq)
- Screen for vector contamination (VecScreen)
- Immunoglobulin BLAST (IgBlast)

Meta

- Retrieve results by RID
- Get this page with javascript-free links

BLAST

Основные задачи компьютерного анализа генетических текстов:

2) статистический анализ генетических текстов, исследование структуры повторов и модели порождения символьных последовательностей, сегментация геномов

Статистическая обработка последовательностей

(частоты нуклеотидов и аминокислот, олигонуклеотидов, l -граммы)

Поиск паттернов

Поиск консервативных компонент в наборе последовательностей.

PROFILE

PSSM (Position Specific Score Matrix)

Stormo, Staden

Нейронные сети , скрытые марковские модели, максимизация ожидания, выравнивание по Гиббсу.

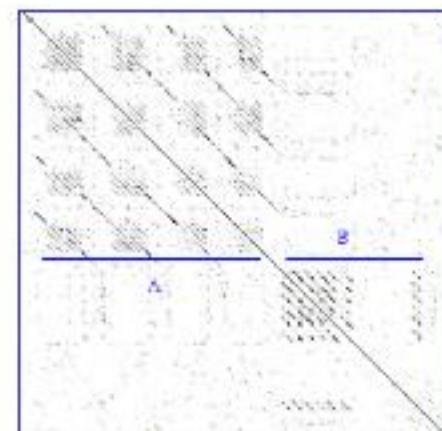
Статистический анализ...

Классификация повторов

Название	Пример	Направление	Комплémentарность
Прямой	AGCTTT TCGAAA	→ → ↑ ↑	Вперед Нет
Инвертированный	AGCTTT TCGAAA	→ → ↑ ↓	Назад Есть
Симметричный	AGCTTT TCGAAA	→ → ↓ ↓	Назад Нет
Прямой комплémentарный	AGCTTT TCGAAA	→ → ↑ ↓	Вперед Есть
Палиндром	AAGCCGAA TTCGGCTT	→ ← ↑ ↓	Назад Нет
Комплémentарный палиндром	AAGCGCTT TTCGCGAA	→ ← ↑ ↓	Назад Есть

Повторы могут пересекаться и накладываться друг на друга в последовательности.

Тандемные и диспергированные
повторы



Последовательность
ДНК:

GTAGTCTGATGCA

Прямой



Инвертированный



Симметричный



Комплементарный



Повторы могут пересекаться и накладываться друг на друга в последовательности.

Повторы могут иметь длину от нескольких пар до нескольких тысяч пар оснований.

Повторы могут быть совершенными (полное совпадение) и несовершенными (допускается частичное несовпадение). Степень вырожденности (несовершенности) повтора фиксированной длины вычисляется (1) по числу несовпадений и (2) по вероятности получить такое число несовпадений по случайным причинам.

Тандемные и



диспергированные повторы



Вопросы из комбинаторики:

Определить вероятность получить последовательность длины l ,
вероятность получить повтор длины l в последовательности
длины n

Модели порождения последовательностей

Анализ закодированных в последовательностях ДНК функциональных сигналов требует применения современных методов распознавания образов, статистических подходов и специальных оптимизированных алгоритмов для преодоления вычислительных трудностей, связанных с обработкой огромных массивов информации.

Вероятностные предсказания функции ДНК возможны лишь при достаточном объеме известных экспериментальных данных, накопленных в специализированных базах данных и при адекватном выборе моделей, описывающих выполнение функции. Надежность предсказания зависит от степени полноты априорной информации (наличия и полноты обучающей выборки, эволюционной близости анализируемых последовательностей).

$$P(\text{ATGCGG})=?$$

P или *Pr* – probability,
вероятность

Модели порождения последовательности:

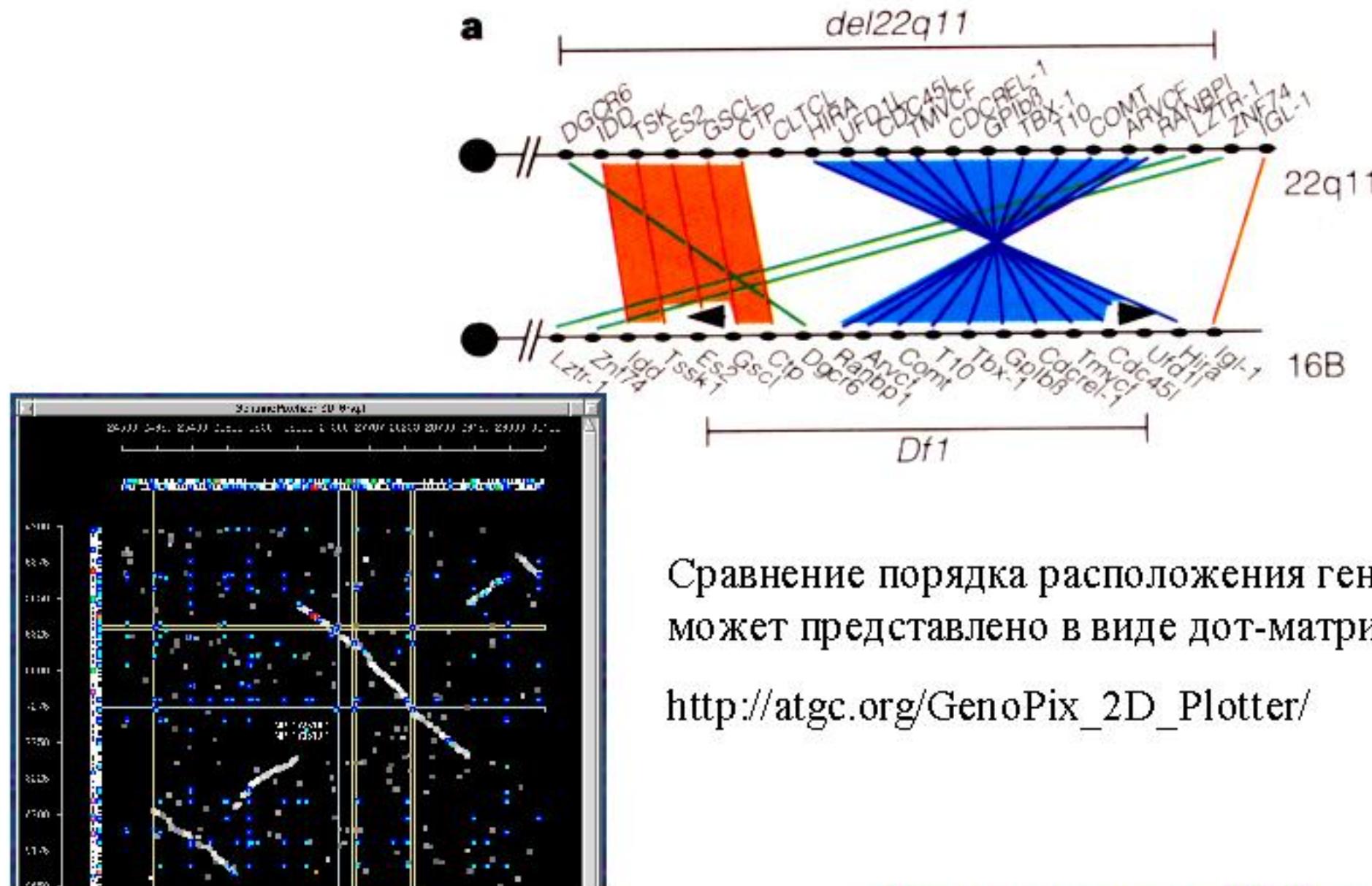
Равновероятная модель: $P(\text{ATGG})=(0.25)^4$

Модель Бернули: $P(\text{ATGG})=P(\text{A}) * P(\text{T}) * P(\text{G}) * P(\text{G})$

Марковская модель 2-го порядка:

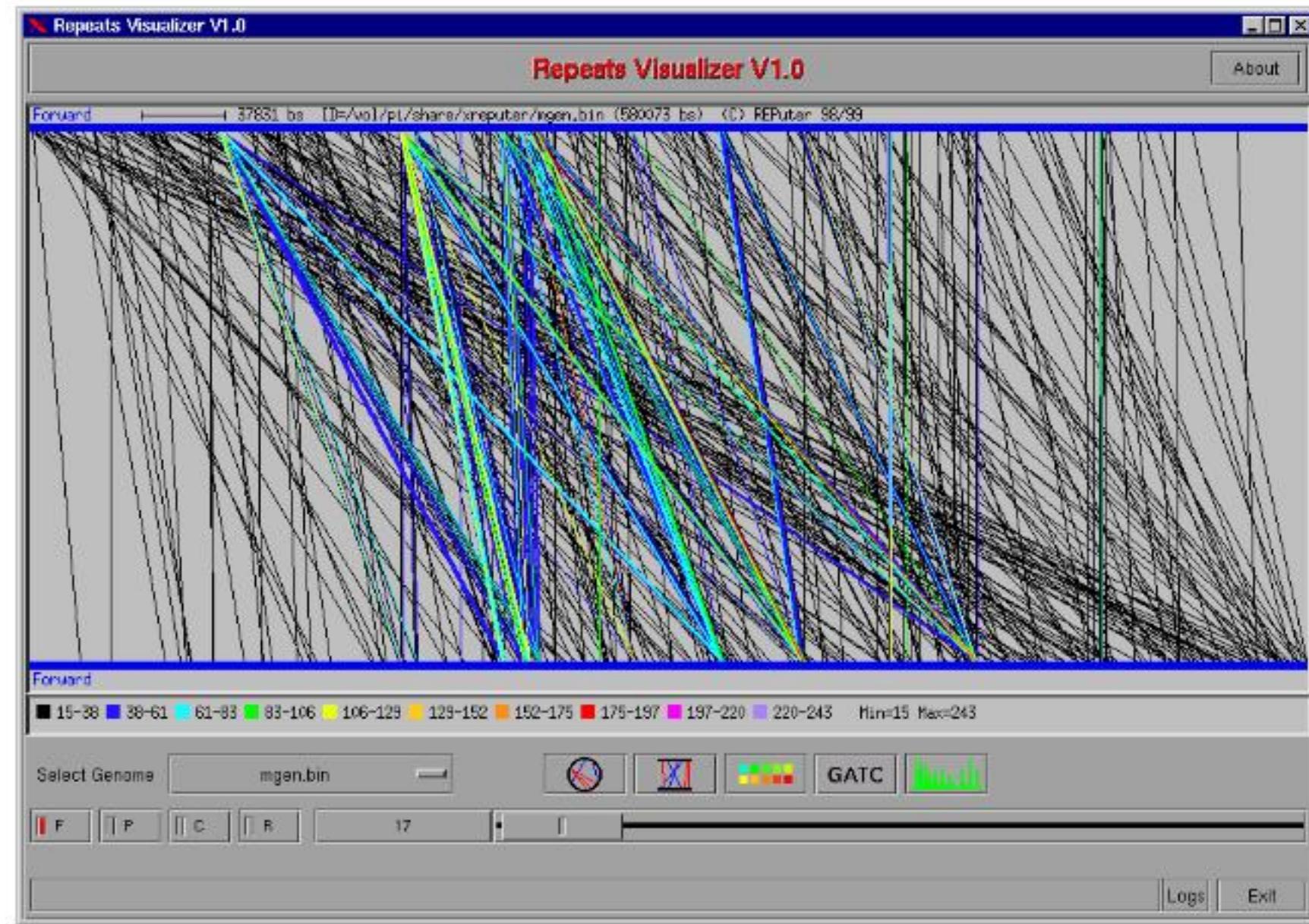
$P(\text{ATGG})=P(\text{A}) * P(\text{T}|\text{A}) * P(\text{G}|\text{T}) * P(\text{G}|\text{G})$

Сравнительный анализ района 22 хромосомы человека и 16 хромосомы мыши



Статистический анализ.

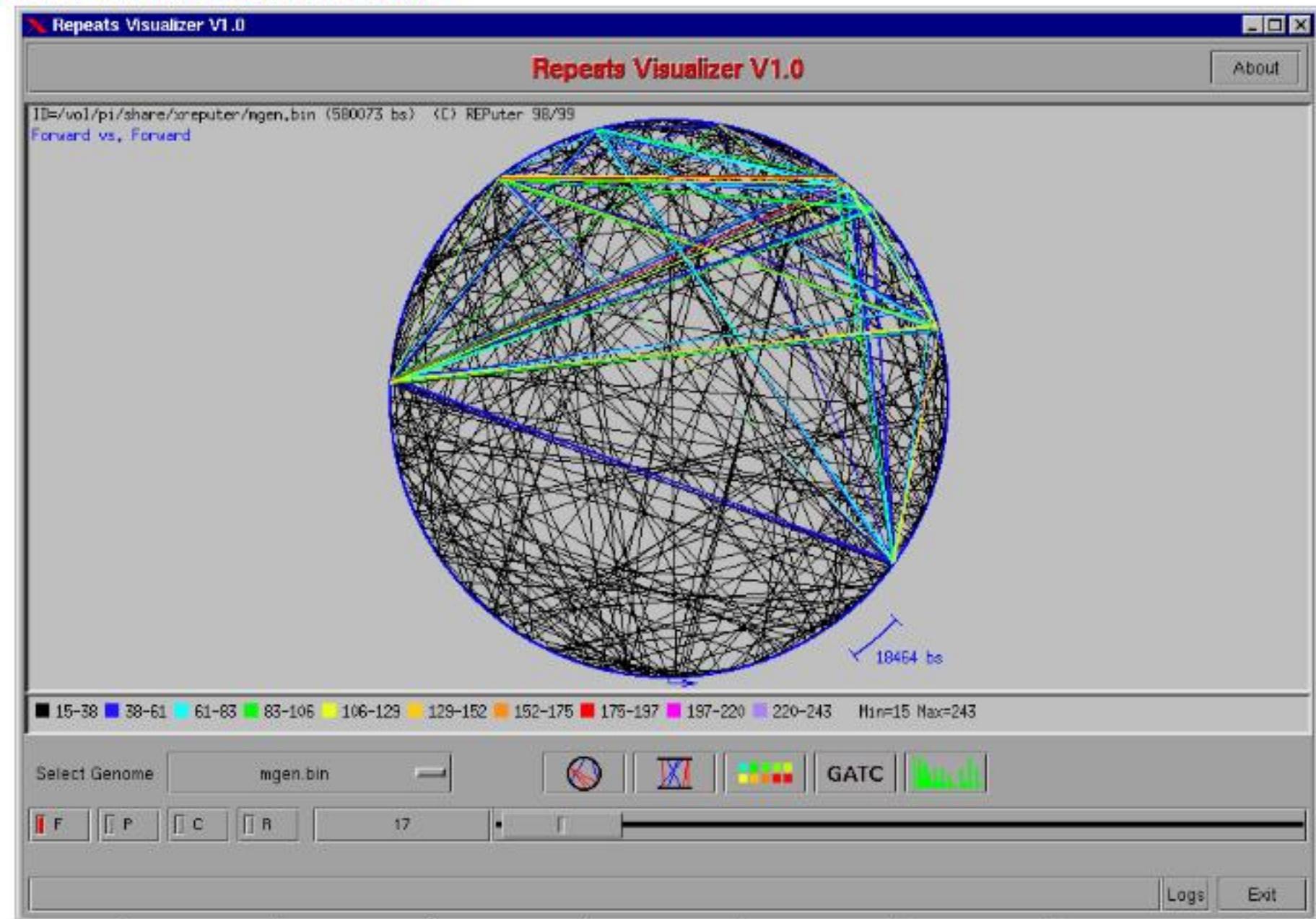
Сравнение и сегментация геномов



Визуализация сравнения полных геномов с помощью программы REPuter
(<http://bibiserv.techfak.uni-bielefeld.de/reputer/>)

Статистический анализ.

Сравнение и сегментация геномов



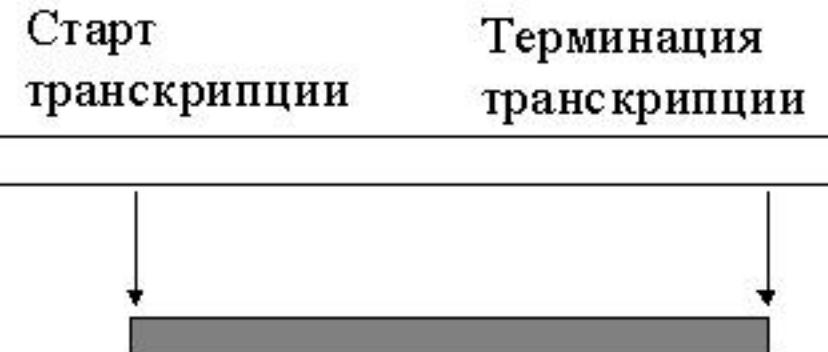
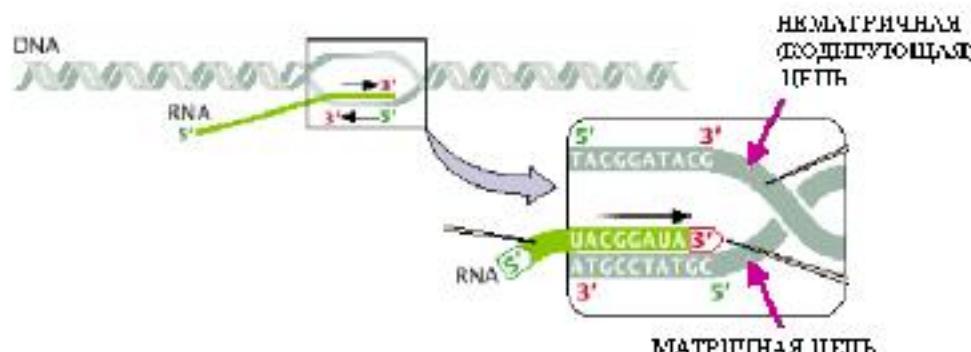
Визуализация сравнения полных геномов с помощью программы REPuter
(<http://bibiserv.techfak.uni-bielefeld.de/reputer/>)

Основные задачи компьютерного анализа генетических текстов:

3) предсказание кодирующих участков генов и открытых рамок считываания;

Бактериальный ген (прокариоты)

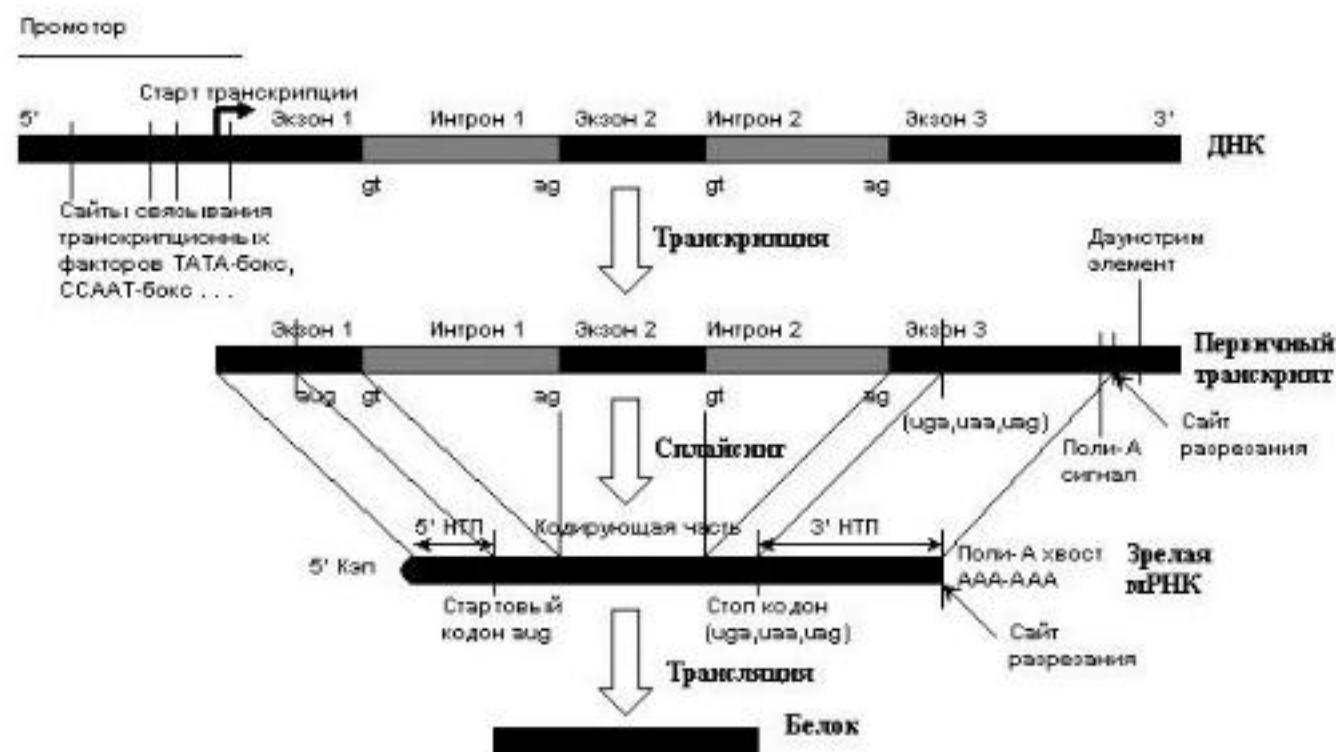
Задача предсказания структуры гена остается одной из важнейших задач биоинформатики, несмотря на полное секвенирование геномов и новые экспериментальные методики



Перекодирование последовательности в аминокислотную

Кодирующий потенциал – частоты гексануклеотидов

Предсказание кодирующих участков. Гены эукариот



ПРИНЦИПЫ работы программы распознавания структуры гена:

- Кодирующий потенциал
- Информация о сайтах сплайсинга
- Скрытые марковские модели
- Сравнение с известными генами по базам данных
- Экспертные системы

Предсказание кодирующих участков. Гены эукариот

Понятие рамки считываания
ORF.

Сдвиг рамки считываания

Статистическая задача
предсказания экзонов



W1
ATCGATGGATGAATAGAGC

W2
CTAACTGATAACGGATGCGGA

S_{1,1} ATC GAT GGA TGA ATA GAG C

S_{2,1} CT AAC TGA TAC GGA TGC GGA

S_{1,2} AT CGA TGG ATG AAT AGA GC

S_{2,2} C TAA CTG ATA CGG ATG CGG A

S_{1,3} A TCG ATG GAT GAA TAG AGC

S_{2,3} CTA ACT GAT ACG GAT GCG GA

...

ATG AATAGAGC - CTA ACT GATAACGGATGCGGA

Кодирующий потенциал
(частоты использования
кодонов)

Определение маршрута, начинающегося с AUG и
заканчивающегося стоп-кодоном или концом
последовательности

Предсказание кодирующих участков. Гены эукариот

Предсказание кодирующих частей генов с помощью программ ORFScan, ESTScan, GENSCAN, GRAIL, GENFIND на примере DT_101886.

Ген human ferrochelatase. Ошибочно определенные позиции маркированы цветом

true : mrslganmaaalraagvllrdplassswrvccqpwrlwksgaaaaaavttetaqhhaqqakpqvqpqkrkpktgilmnmggpetlgdvhdfllrlflqdmlmtlpqknklapf
ORFScan : MRSLGAMMAAALRAAGVLLRDPLASSSWRVCCQPWRLWKSGAAAAAAVTTETAQHAAQGAKPQVQPQKRKPKTGILMLNMGGPETLGDVHDPLLRLFLDQDLMTLPIQNKLAPF

ESTScan : NAFIGANMAAALRAAGVLLRXSAGIQGLEGLSAMEEVEUGAAAAAVTTETAQHAAQGAKPQVQPQKRKPKTGILMLNMGGPETLGDVHDPLLRLFLDQDLMTLPIQNKLAPF

GENSCAN : _____ MAAALRAAGVLLRDPLASSSWRVCCQPWRLWKSGAAAAAAVTTETAQHAAQGAKPQVQPQKRKPKTGILMLNMGGPETLGDVHDPLLRLFLDQDLMTLPIQNKLAPF

GRAIL : MRSLGAMMAAALRAAGVLLRDPLASSSWRVCCQPWRLWKSGAAAAAAVTTETAQHAAQGAKPQVQPQKRKPKTGILMLNMGGPETLGDVHDPLLRLFLDQDLMTLPIQNKLAPF

GenFind : _____ MAAALRAAGVLLRDPLASSSWRVCCQPWRLWKSGAAAAAAVTTETAQHAAQGAKPQVQPQKRKPKTGILMLNMGGPETLGDVHDPLLRLFLDQDLMTLPIQNKLAPF

true : iakrrtpkiqeinqrrriggspikwtksqgegmvk11delspntaphkyyigfrvhplteaeemerdgeraiatftqypqyscsttgslnaiyryynqvgrkptmk
ORFScan : IAKRRTPKIQEQQYRRLEADPP3RY??SKQGEGMVKLLDELSPNTAPHKYYIGIRYVHPLTEE&IEEMERDGLERAIAFTQYPQYSCTTGSSLNAIYRYYNQUGRKPTMK

ESTScan : IAKRRTPKIQEQQYRRXWRRIPHQIMDFQAGRHH_____APHKYYIGIRYVHPLTEE&IEEMERDGLERAIAFTQYPQYSCTTGSSLNAIYRYYNQUGRKPTMK

GENSCAN : IAKRRTPKIQEQQYRRXWRRIPHQIMDFQAGRHH_____APHKYYIGIRYVHPLTEE&IEEMERDGLERAIAFTQYPQYSCTTGSSLNAIYRYYNQUGRKPTMK

GRAIL : IAKRRTPKIQEQQYRRTWRRIPHQIMDFQAGRHH_____

GenFind : IAKRRTPKIQEQQYRRTWRRIPHQIMDFQAGRHH_____EEAIEEMERDGLERAIAFTQYPQYSCTTGSSLNAIYRYYNQUGRKPTMK

true : wstidrwpthhlliqcfadhlkeldhfplekrsevlfshslpmssvvnrqdpypqevsatvqkvmerleycnpyrlvwqskvgpmpwlgpqtdesikglcergrknii
ORFScan : WSTIDRWPTHHLLIQCFADHLKELDHFPLEKRSEVJILF3AHSLPMSSUUNRGDPYPQEVSATUQKIMERLEYCNPYRLUWQSKUGPMPWLGPQTDESIKGLCERGRKNI

ESTScan : WSTIDRWPTHHLLIQCFADHLKELDHFPLEKRSEVJILF3AHSLPMSSUUNRGDPYPQEVSATUQKIMERLEYCNPYRLUWQSKUGPMPWLGPQTDESIKGLCERGRKNI

GENSCAN : WSTIDRWPTHHLLIQCFADHLKELDHFPLEKRSEVJILF3AHSLPMSSUUNRGDPYPQEVSATUQKIMERLEYCNPYRLUWQSKUGPMPWLGPQTDESIKGLCERGRKNI

GRAIL : _____GEAAGQEVSATUQKIMERLEYCNPYRLUWQSKUGPMPWLGP_____

GenFind : WSTIDRWPTHHLLIQCFADHLKELDHFPLEKRSEVJILF3AHSLPMSSUUNRGDPYPQEVSATUQKIMERLEYCNPYRLUWQSKUGPMPWLGPQTDESIKGLCERGRKNI

true : llvpiaftsdhietyldeysqlakecgvenirraeslgnmplfskaladlvhshiqsnelskqlt1scplcvnpvcrtksfftsqql

ORFScan : LLVPIAFTSDHIETLYEYDIEYSQULAKECGVENSIRRAESLNGIHCS?KALADLV?FTHPUKRAUFQAADPDCPLCUNPVCRETKSFPT?PAAUTPAGGPRGVSKCPT3RYLRCGEGUI

ESTScan : LLVPIAFTSDHIETLYEYDIEYSQULAKECGVENSIRKELSLFWGKFRCSLKAQGRLGAFTHPSQTSCUPKQLTPELUNPVCRETKSFHMPA33CEPPPUUDFUALANQPPDTSDVER

GENSCAN : LLVPIAFTSDHIETLYEYDIEYSQULAKECGS_____

GRAIL :

Пример программы предсказания структуры гена GenScan

The New GENSCAN Web Server at MIT

Identification of complete gene structures in genomic DNA

2017



[For information about Genscan, click here]

This server provides access to the program Genscan for predicting the locations and exon-intron structures of genes in genomic DNA.

<http://genes.mit.edu/GENSCAN.html>

Предсказание кодирующих участков. Гены эукариот

Результаты предсказания
кластера глобиновых генов
человека (73308 по) с
помощью программы
GraileXP v3.31 [March, 2002]

<http://grail.lsdornl.gov/grailexp/>

Предсказание кодирующих участков. Гены эукариот

Название программы, ссылка Интернет-адрес, краткое описание

GENEID (Wiehe et al., 2001)	http://www1.inim.es/geneid.html (R. Guigo, Spain Institut Municipal de Investigacio Medica, Испания)
SLAM(Pachter et al, 2002)	http://bio.math.berkeley.edu/slam/ (Марковские модели и парное выравнивание)
GENIE (Reese et al., 2000)	http://www.cse.ucsc.edu/~dkulp/cgi-bin/genie (Скрытые марковские модели)
SELFID (Audic and Claverie, 1998)	http://igs-server.crrs-mrs.fr/~audic/selfid.html (Франция, поиск генов в микробной ДНК)
MZEF(Zhang 1997)	http://argus.cshl.edu/genefinder/ (Cold Spring Harbor Lab, США, Квадратичный дискриминантный анализ)
WEBGENE (Milanesi et al, 1999)	http://www.itba.mi.cnr.it/webgene/ (ITBA, CNR, Milan, Italy)
GeneMark (Lukashin and Borodovsky, 1998)	http://opal.biology.gatech.edu/GeneMark/ (GIT, Borodovsky's lab, School of Biology, США, Скрытые марковские модели)
FrameD(Schiex et al, 2000)	http://www.toulouse.inra.fr/FrameD/cgi-bin/FrameD (INRA, Toulouse, Франция) поиск генов и рамок считываания в G+C богатых прокариотических последовательностях
EuGene(Schiex et al., 2001)	http://www-bia.inra.fr/T/EuGene/ поиск генов <i>Arabidopsis thaliana</i>
GLIMMER(Delcher et al, 1999)	http://www.tigr.org/~salzberg/glimmer.html (TIGR, Salzberg's lab) поиск генов в ДНК микробов
VEIL(Henderson et al., 1997)	http://www.tigr.org/~salzberg/veil.html (VEIL - the Viterbi Exon-Intron Locator. Скрытая марковская модель)
MORGAN(Salzberg et al, 1998)	http://www.tigr.org/~salzberg/morgan.html (Решающие деревья для поиска генов в ДНК позвоночных)
GENSCAN(Tiwari et al., 1997)	http://202.41.10.146/ (Jawaharlal Nehru Univ., Индия. Поиск генов с использованием преобразований Фурье)
GENSCAN	http://genes.mit.edu/GENSCAN.html (C. Burge, Massachusetts Institute of Technology, США)
Diogenes	http://www.cbc.umn.edu/diogenes/index.html (США. Предсказание для коротких последовательностей)
GRAIL(Uberbacher et al., 1996)	http://compbio.ornl.gov/Grail-bin/EmptyGrailForm (Oak Ridge National Lab, США)
FGENES (Solovyev and Salamov, 1997)	http://genomic.sanger.ac.uk/gf/gf.html (Sanger Centre, UK) http://www.softberry.com/berry.phml (новая версия FGENES)
HMMGENE(Krogh, 1997)	http://www.cbs.dtu.dk/services/HMMgene/ (Technical Univ. of Denmark, Дания. Скрытые марковские модели)
YEASTGENE(Zhang and Wang, 2000)	http://ubic.tju.edu.cn/cgi-bin/Yeastgene.cgi (TianJin University, Китай, техника распознавания Z-curve)
GENEPARSER(Snyder and Stormo, 1995)	http://beagle.colorado.edu/~eesnyder/GeneParser.html (Динамическое программирование и нейронные сети)

Предсказание кодирующих участков. Гены эукариот

Вопрос компьютерного предсказания генов, экзон-инtronной структуры по нуклеотидной последовательности остается открытым.

В настоящее время даже число генов в полностью секвенированном геноме человека остается неизвестным – разные компьютерные методики дают оценку от минимум 24500 до 45000 генов. Минимальную оценку дают сравнения с кодирующими частями родственных геномов (например мыши), максимальную – предсказание по свойствам самой последовательности ДНК (GENSCAN).

Эти цифры сильно уменьшены по сравнению с первоначальными оценками в 100 тысяч генов, затем 50 тысяч генов.

Есть проблемы в определении понятия гена – как считать варианты альтернативного сплайсинга, низкоэкспрессирующиеся гены, короткие гены (один экзон) и гены, кодирующие только РНК, но не белок.

Проблемой остается предсказание генов на так называемой «темной стороне» (dark side) генома человека – районов с низкой плотностью генов. Существующие программы предсказания генов не могут адекватно находить гены в этих частях генома.

Elizabeth Pennisi. (2003) Gene Counters Struggle to Get the Right Answer. *Science* Aug 22 : 1040-1041.

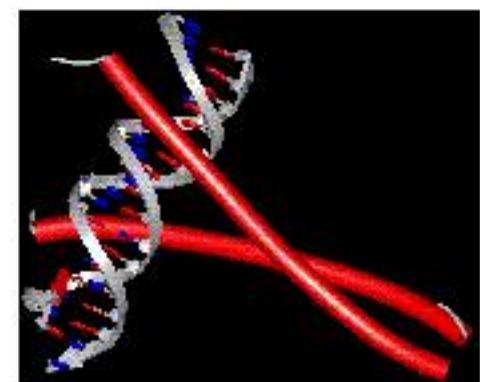


Основные задачи компьютерного анализа генетических текстов:

4) предсказание функциональных сигналов (функциональных сайтов и регуляторных районов)

Понятие консенсуса, весовой матрицы, паттерна

консенсус ТАТА-бокс связывающего белка (ТВР) имеет вид ТАТААААА, расширенный аналог в вырожденном 15-буквенном коде имеет вид
STWTAWADRSSSSSS



Соответствие представлений функц.сайта в 4-буквенном и 15-буквенном алфавитах

---**ТАТАААА**---

STWTAWADRSSSSSS

Более точным способом представления и анализа выборок выровненных последовательностей длины L являются весовые матрицы размерности $L \times 4$. Элемент $f(i,j)$ весовой матрицы $F = |f(i, j)|$ определяет частоту встречаемости нуклеотида i ($i = 1, 2, 3, 4$ соответствует символам A, T, G и C) в позиции j ($j = 1, \dots, L$), подсчитанную по выборке выровненных нуклеотидных последовательностей. Оптимизированная весовая матрица $W = |w(i, j)|$ может быть вычислена в логарифмической форме с учетом ожидаемых частот

$$w(i, j) = \ln\left(\frac{f(i, j)}{e(i, j)} + \frac{s}{100}\right) + c(i)$$

Понятие весовой матрицы

A	5	0	0	3	18	0	3	4	2
T	4	1	1	0	9	2	1	1	5
G	7	3	2	1	13	22	19	15	17
C	27	39	40	39	3	19	20	23	19
	C	C	C	C	A/G	G/C	G/C	G/C	
	C	C	C	C	R	S	S	S	S

R	G/A
Y	T/C
M	A/C
K	G/T
W	A/T
S	G/C
B	-A
V	-T
H	-G
D	-C
N	A/T/G/C

Предсказание функциональных сигналов

Пример весовой матрицы и консенсуса ТАТА-бокса в промоторах эукариот

- TATA-box HMM trained from 900 unrelated general promoter sequences:

Position	1	2	3	4	5	6	7	8	9	10	11	12
% A	21.4	15.9	3.7	91.1	0.0	94.5	67.3	97.3	52.1	40.7	16.5	23.6
% C	22.7	39.3	9.8	0.0	0.0	0.0	0.0	0.0	0.0	9.1	34.8	37.1
% G	28.2	35.2	2.9	0.0	0.0	0.0	0.0	2.7	12.0	40.2	38.0	30.4
% T	27.7	9.6	83.6	8.9	100.0	5.5	32.7	0.0	35.9	10.0	10.7	8.9
Consensus			T	A	T	A	W	A	W	R		

ТАТА-бокс в промоторах растений:

- TATA-box HMM trained from 134 unrelated plant promoter sequences:

Position	1	2	3	4	5	6	7	8	9	10	11	12
% A	31.6	16.3	2.0	90.8	0.0	94.9	57.1	100.0	27.6	69.4	11.2	24.5
% C	24.5	60.2	3.0	2.1	0.0	0.0	0.0	0.0	0.0	3.1	39.8	52.0
% G	15.3	10.2	0.0	2.0	1.0	0.0	0.0	0.0	2.0	13.3	37.8	21.4
% T	28.6	13.3	94.9	5.1	99.0	5.1	42.9	0.0	70.4	14.3	11.2	2.1
Consensus			T	A	T	A	W	A	W	A		

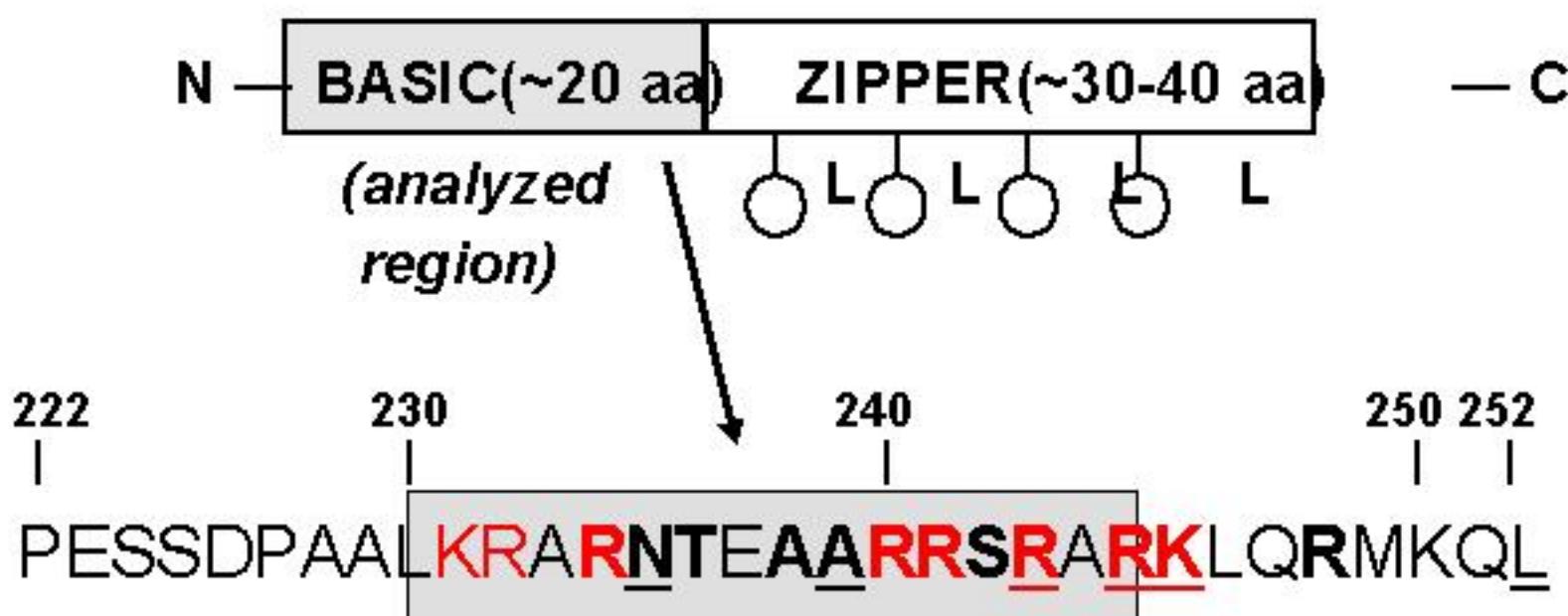
Понятие паттерна как вида записи функционального сайта:

TAT(A)₂₋₅

R[Y]SA[G][T]G₀₋₃C₀₋₆

структура домена bZIP

Функция: связывание с ДНК формирование димера



5) анализ вторичной структуры РНК и сигналов трансляции;

Предсказание вторичной структуры РНК

Шпилька РНК –

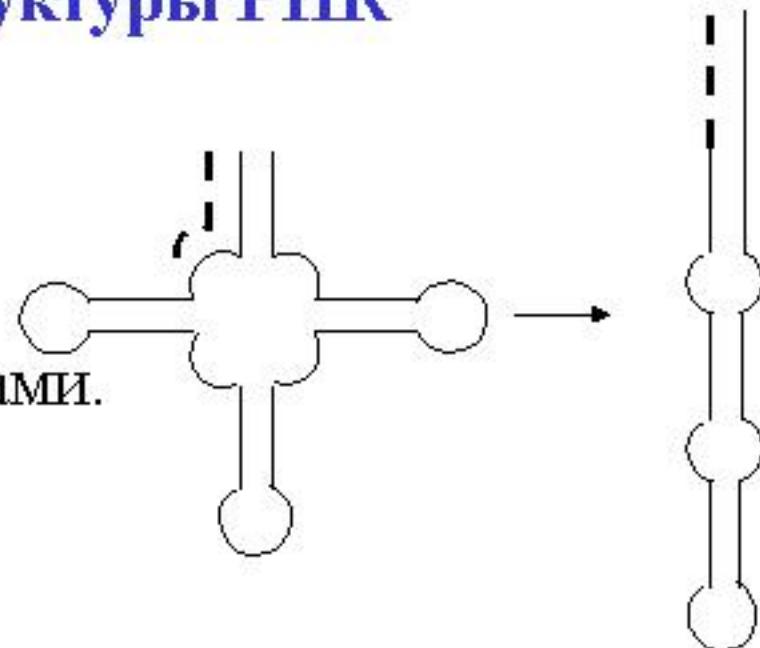
Оценка энергии по числу

водородных связей между нуклеотидами.

G-C (3H)

A-U (2H)

G-U (1H)



Задача – выбор оптимальной, наиболее энергетической вторичной структуры по анализу нуклеотидной последовательности

Nussinov and Jacobson, 1980; Zuker and Stiegler, 1981

ribosomal RNA database

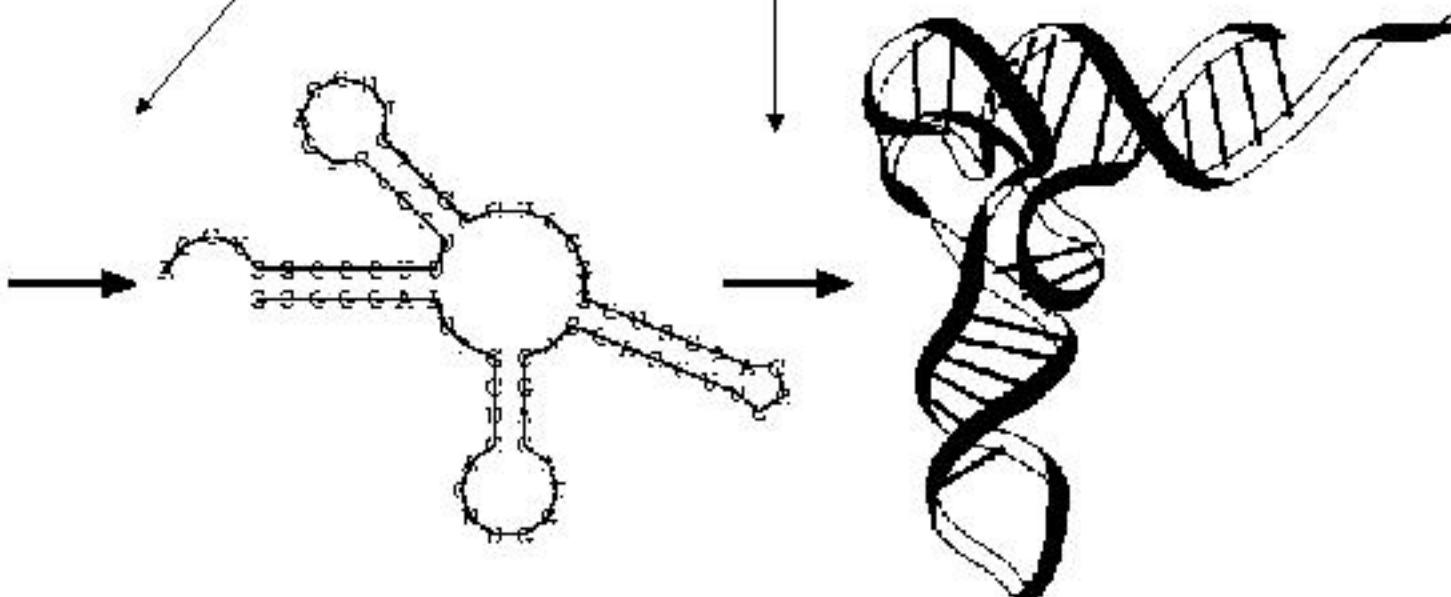
<http://www.cme.msu.edu/RDPhtml/index.html>

Анализ вторичной структуры РНК

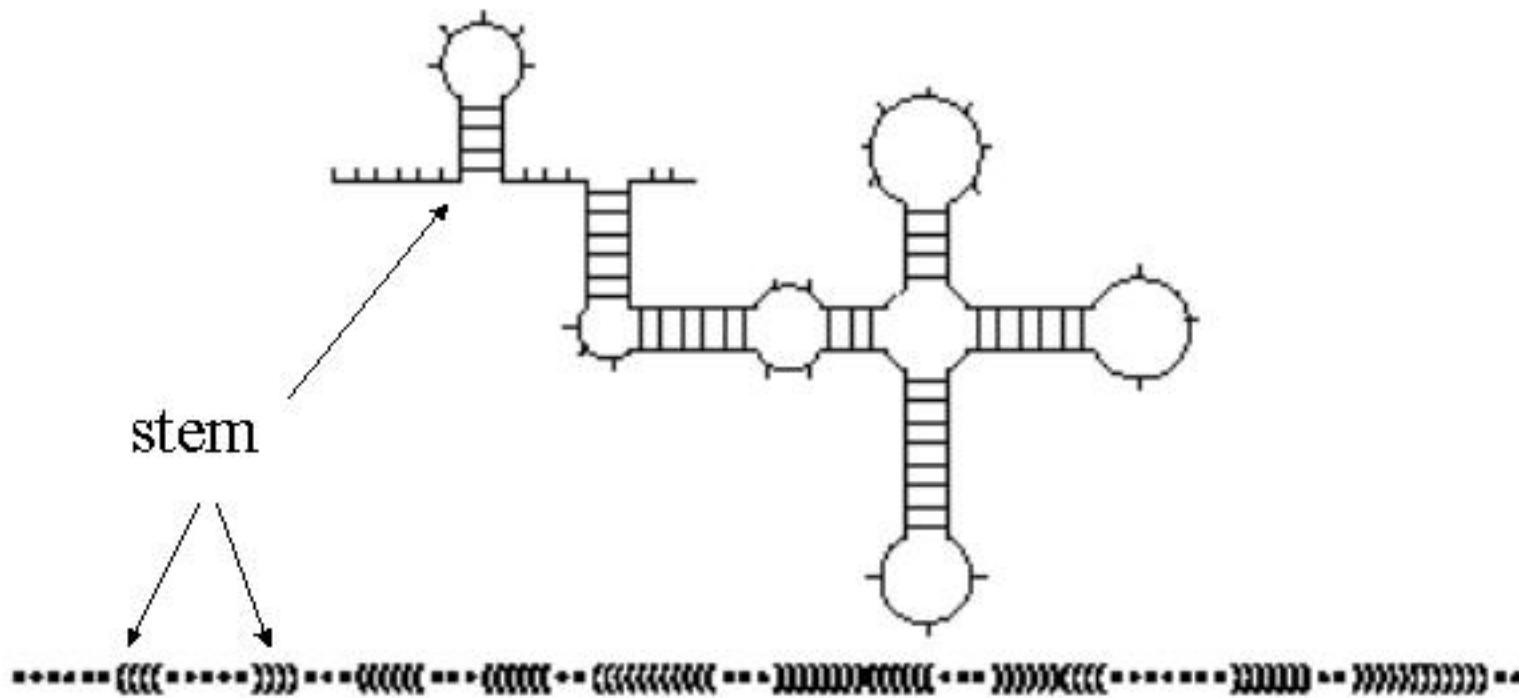
Предсказание
пространственной
структуры

tRNA^{Phe}

© С. В. Баранов, Университетская школа геномики и биотехнологии

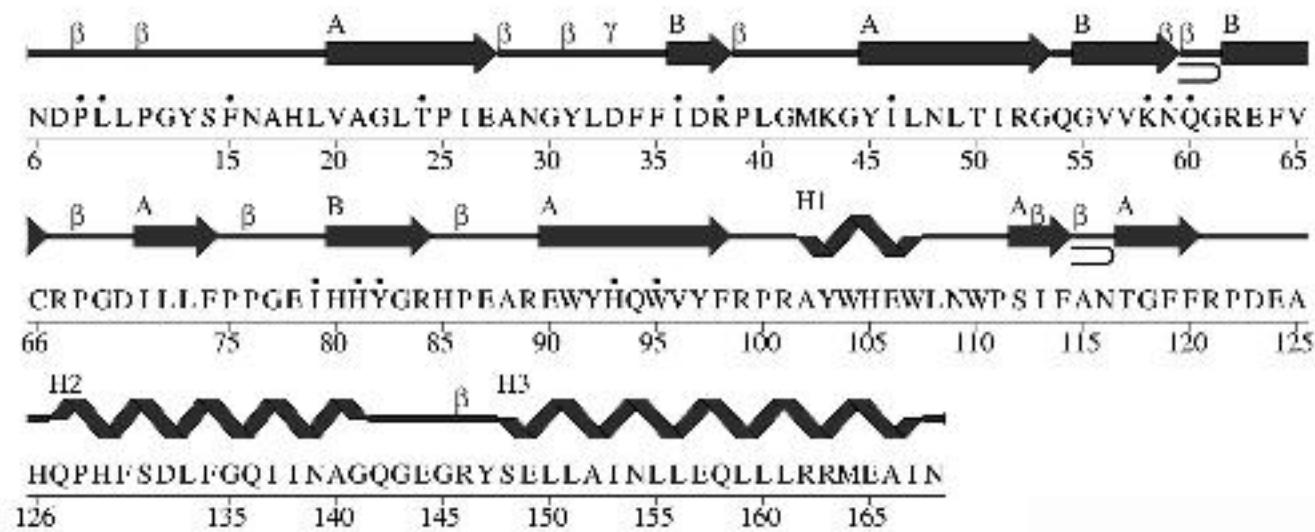


Представление (визуализация) вторичной структуры – матричное, линейное, графическое.



Линейное представление структуры – запись с помощью скобок (бракет)

б) анализ аминокислотных последовательностей белков, предсказание вторичной структуры, функциональных сайтов и доменов глобулярных белков по их аминокислотным последовательностям;



Уровни структурной организации белка на примере регуляторного белка lacR из *E. coli* (идентификатор PDB 2aac). На рисунке показаны: а) первичная структура белка – аминокислотная последовательность; и разметка вторичной структуры белка, выполненная программой PDBSum (Laskowski, 2001; <http://www.biochem.ucl.ac.uk/bstm/pdbsum/2aac/main.html>); б) пространственная структура димера, мономеры показаны разными оттенками серого цвета, элементы вторичной структуры показаны схематически в виде цилиндров (α -спиралей) и стрелок (α -нитей).



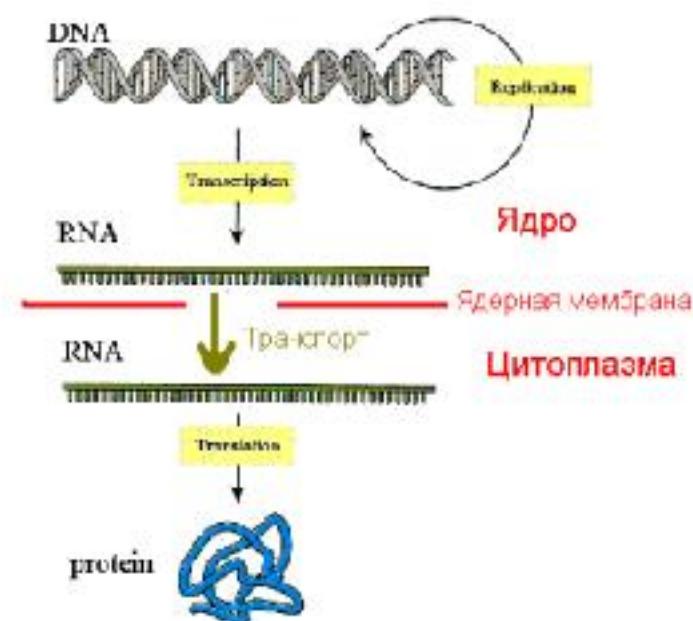
Построение последовательности белка по ДНК.

Предсказание вторичной структуры белка. Постановка задачи.

Полипептидная цепь –

Вторичная структура –

Трехмерная структура



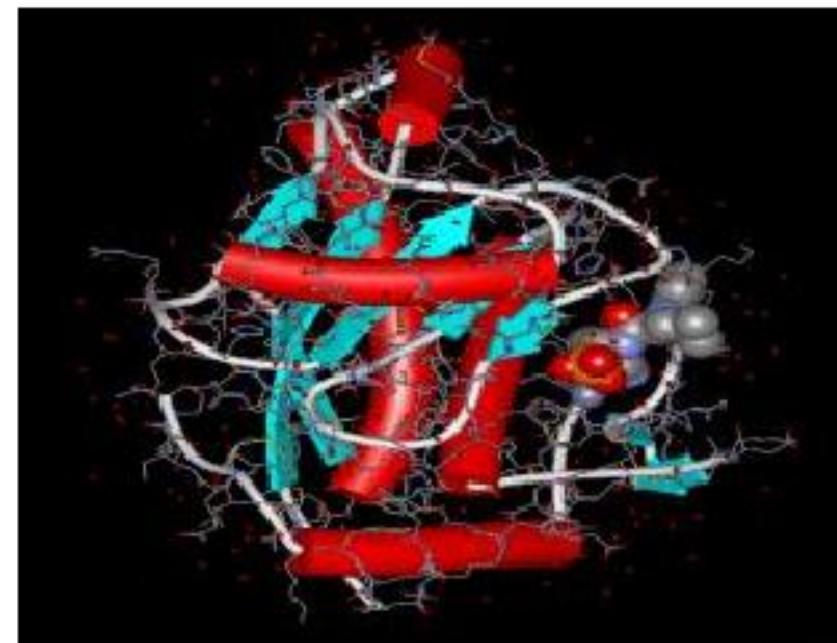
Адаптировано из:

URL Address: <http://www.cytochemistry.net/Cell-biology/ribosome.htm>
Gwen V. Childs, Ph.D.

Предсказание структуры белка

Расшифрованы
аминокислотные
последовательности
сотен тысяч белков

№	330	340	350	360	370	380
1ATP_A	● LGTMCDPKISIGIVQDESPIINILNGWTNAEIGRNLLGMEEDGWWEDCLRGWVWWGICLHRPGL					
query	● VEGMCRA-SGGGVSTDESELPIGAATNAEIGBSLGISED					
462302	● MGGMCNAKNSVGVKIDESSSNVFMVAUTHTHEIGRNLLGMEEDGWWEDCLRGWVWWGICLHRPGL					
462095	● VGGMCQKLNSTGVIQDESAINILVALTNAEIGRNLLGMEEDGWWEDCLRGWVWWGICLHRPGL					
123525	● TCGMCUPRNSVGIVQDESPLKTLIAUTNAEIGRNLLGMEEDGWWEDCLRGWVWWGICLHRPGL					
2231613	● MASMCDPKRSVGIQIDESSTINIMMAUTNAEIGRNLLGMEEDGWWEDCLRGWVWWGICLHRPGL					
1916617	● LEGMCTA-NSGGGWSHHSSENALIGAAATNAEIGRNLLGMEEDGWWEDCLRGWVWWGICLHRPGL					
7993727	● IHSMKCTadQSGGIVHIDESDNPIGAATLAEIGRNLLGMEEDGWWEDCLRGWVWWGICLHRPGL					
1709103	● VSALCSR-HSGAWNQDESNSIGVASTNAEIGRNLLGMEEDGWWEDCLRGWVWWGICLHRPGL					
2499914	● VSAMCSH-SGCAWNQDESNSKRPVGAVSTNAEIGRNLLGMEEDGWWEDCLRGWVWWGICLHRPGL					

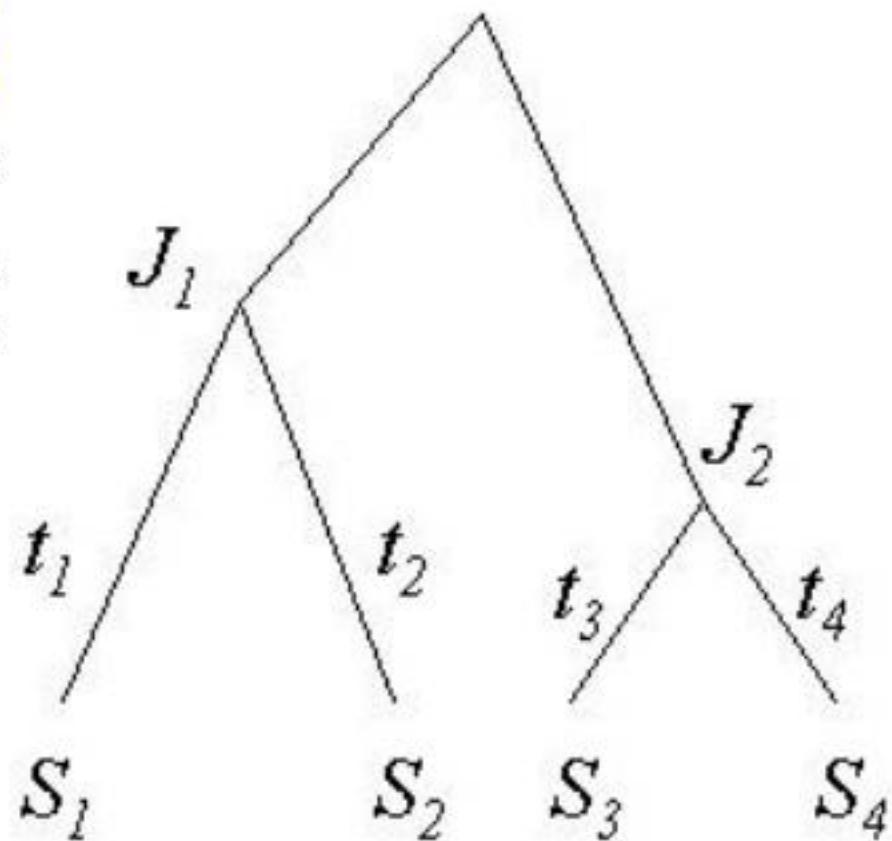


Задачи анализа белков будут рассмотрены в отдельном курсе
к.б.н. Афонникова Д.А.

7) филогенетические сравнения

Рассмотрим задачу восстановления филогенетического построения дерева по набору последовательностей (ДНК или белков). В качестве исходных данных выступают современные последовательности S (S_1, S_2, \dots, S_n). Как правило, последовательности S рассматриваются гомологичны и имеется множественное выравнивание.

Результат построения – дерево T , висячие вершины которого представляют современные последовательности S (S_1, S_2, \dots, S_n), а внутренние J (J_1, J_2, \dots, J_n) – гипотетические предковые.



Ребра (ветви) дерева имеют длину t (t_1, t_2, \dots, t_n). При этом дерево удовлетворяет ряду требований, зависящих от конкретной методики, с целью получение топологии дерева, максимально близкой к истинной.

Заметим, что даже верно восстановленное дерево, построенное для генов, не всегда соответствует видовому дереву, построенному по фенотипическим характеристикам.

Установление эволюционных взаимосвязей по последовательностям

Рибосомальные РНК, белки

Метод Фитча и Марголиаша

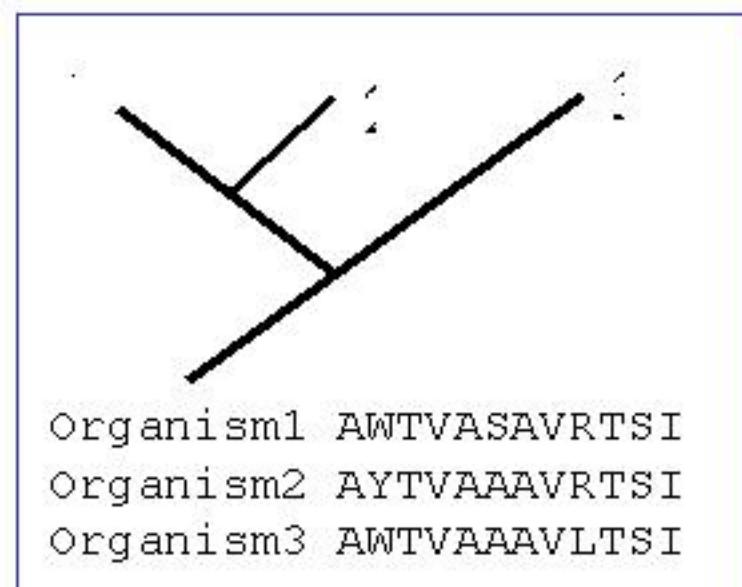
Fitch and Margoliash, 1987

Метод объединения соседей

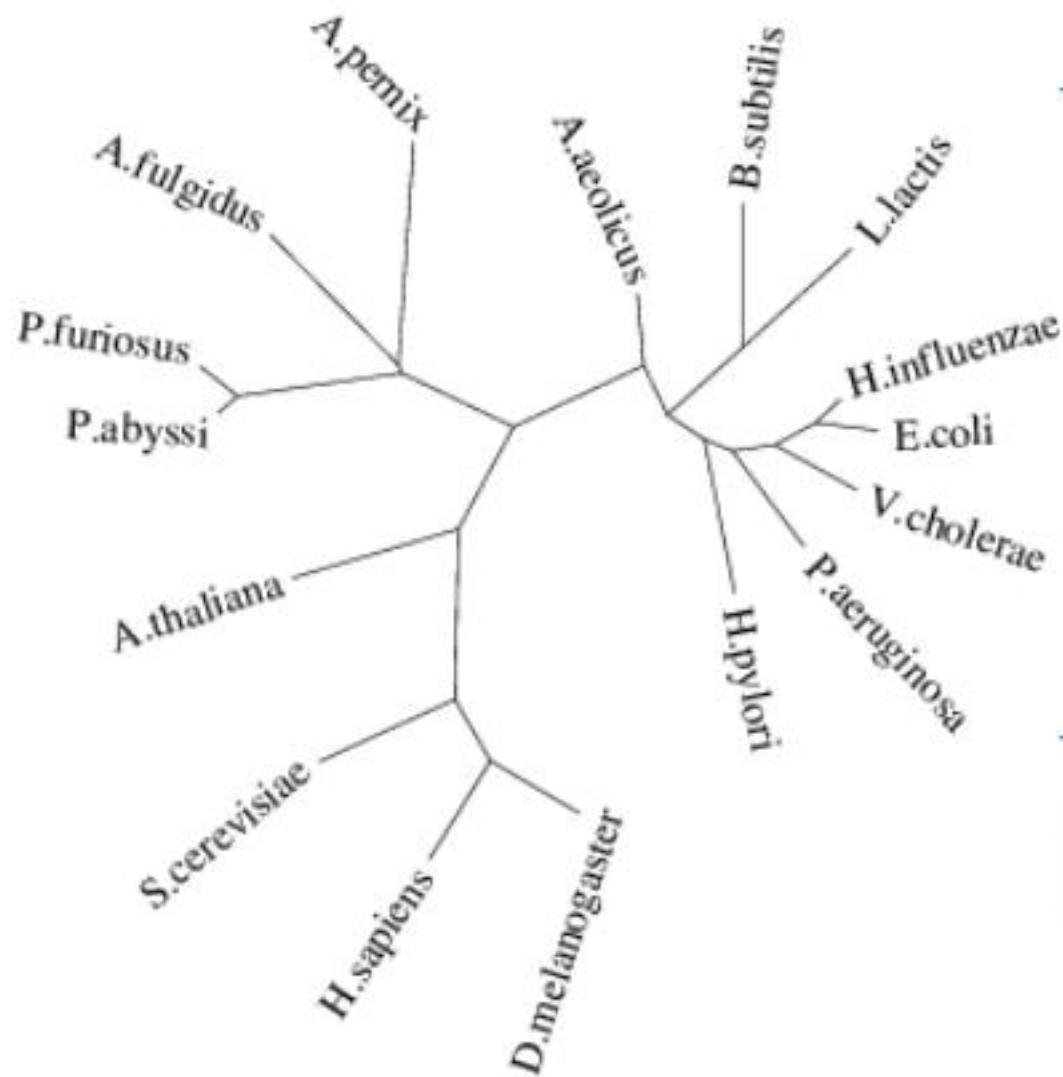
Saitou and Nei, 1987

Метод максимальной парсimonии

Felsenstein, 1988



Дерево построено на основе сходства групп ортологичных генов



SHOT - Shared
Ortholog and Gene
Order Tree
Reconstruction Tool

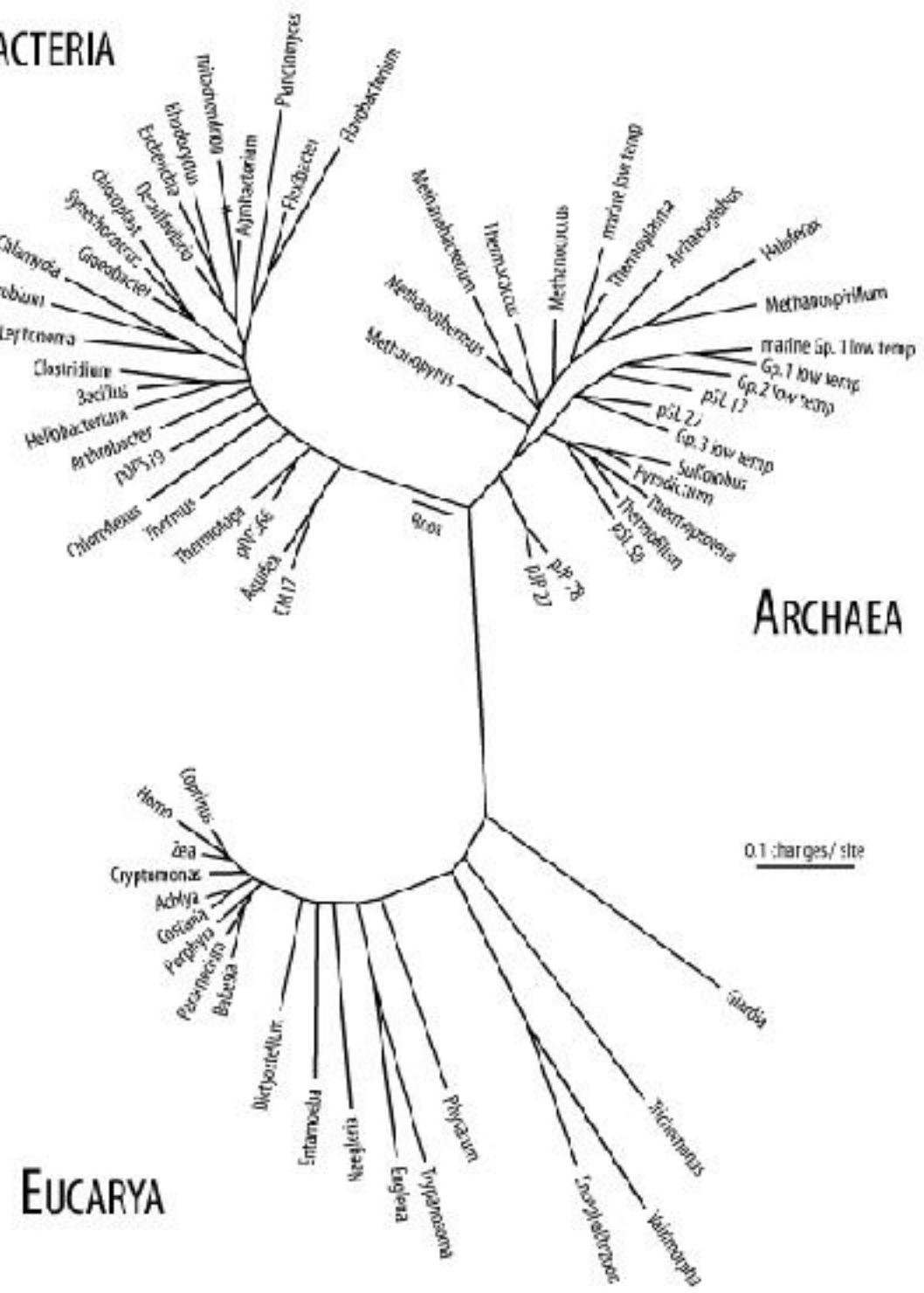
<http://www.bork.embl-heidelberg.de/~korbel/SHOT/>

филогенетические сравнения

BACTERIA

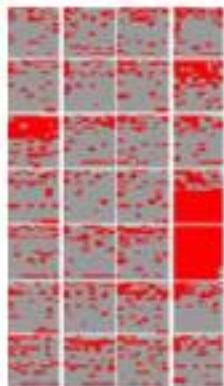
Дерево Жизни построено по последовательностям рРНК

(Diagram courtesy of Norman Pace)

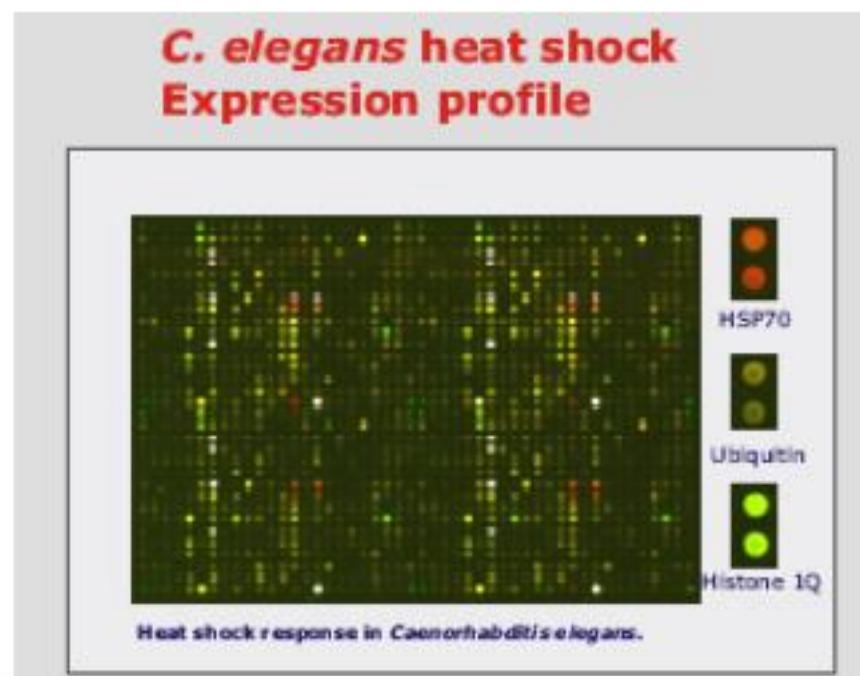


8) ДНК-чибы – экспрессионные кривые

Выявление и анализ закодированных в последовательностях ДНК функциональных сигналов требует применения современных методов распознавания образов, статистических подходов и применения специальных оптимизированных алгоритмов для преодоления вычислительных трудностей, связанных с обработкой огромных массивов информации.



Экспериментальная технология – закрепление олигонуклеотидов (30 нт) на стекле, флюоресцентное мечение гибридизации.

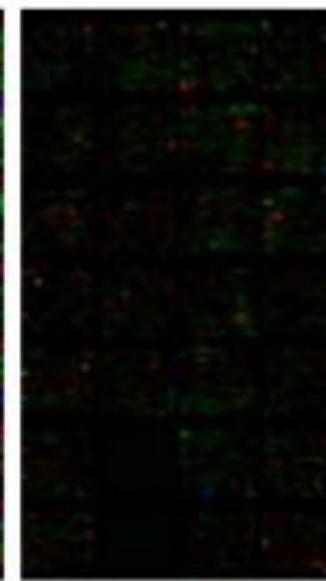
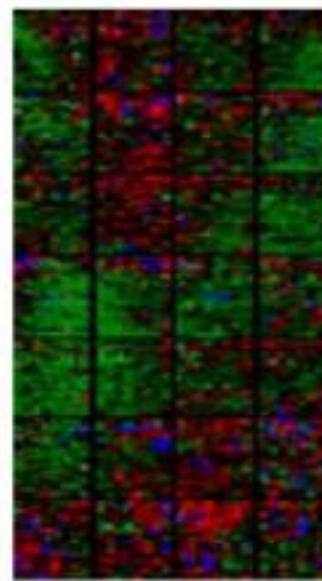


Компьютерные задачи – выбор оптимальных олигонуклеотидов для проб. Восстановление последовательности по олигонуклеотидам и, Актуальное направление сегодня – Анализ кривых генной экспрессии в зависимости от условий эксперимента (ткани, времени, стадии клеточного деления, и т.д.)

<http://oligo.lnatoools.com/expression/>

Микропробы ДНК (микроэррэй – місгоаггау) широко используются в биологических исследованиях. По анализу различной гибридизации на одной пластине с точечно нанесенными пробами можно определить изменения в уровнях экспрессии мРНК, вариации в числе копий ДНК и расположение сайтов связывания транскрипционных факторов в геномной шкале.

При выполнении экспериментов наибольшая проблема – обработка больших, зашумленных данных с целью. Определения специфического набора элементов, которые действительно гибридизуются различным образом. Такая обработка требует объединения различных методов и программ в единую технологическую линию.



<http://array.mbb.yale.edu/analysis/>

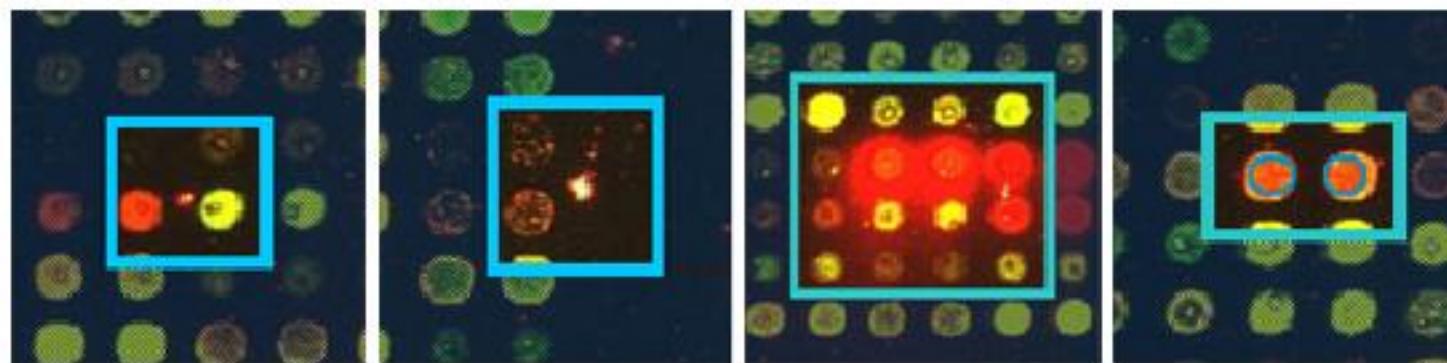
<http://www.cbs.dtu.dk/services/GenePublisher/>



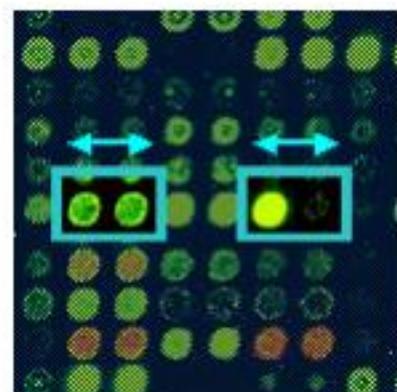
Хотя технология микроЭРРЫ относительно нова, многие аспекты анализа данных после стадии эксперимента хорошо определены.

Это измерение интенсивности флюоресценции, оцифровка изображения микрочипа с помощью компьютерных алгоритмов, кластеризация сходно экспрессирующихся генов и интеграция данных эксперимента с геномной информацией (базами данных).

Научная проблема – как трактовать численные данные, полученные сразу после сканирования и оцифровки изображения. Этой цели служит обработка данных (процессинг): (i) определение и минимизация уровня шума, связанного с экспериментом, (ii) оценка качества данных, полученных в эксперименте и (iii) идентификация элементов чипа, которые действительно по-разному гибридизуются.



Примеры точек (проб)
гибридизации



9) задачи оперирования с большими массивами информации и управления (Интернет-навигации) разрозненными специализированными базами данных

Gene Express 2.1

INSTITUTE OF CYTOLOGY AND GENETICS

Databases

Transcription Regulatory Response Database (TRRD Rel.4.2.x)
[TRRDGENES4](#) [TRRDEIB4](#) [TRRDEITES4](#) [TRRDEXP4](#) [TRRCFACTORS4](#) [TRRDUNITS4](#)

GeneNet Database
[GN](#) [CELL](#) [GN](#) [GENE](#) [GN](#) [RNA](#) [GN](#) [LITERATURE](#) [PROCESS](#) [GN](#) [SUBSTANCE](#)
[GN](#) [PROTEIN](#) [GN](#) [SCHEME](#) [ON](#) [RELATION](#) [ON](#) [CROAN](#) [SM](#) [ON](#) [COMPARTMENT](#)
[GN](#) [EXPERIMENT](#)

Research with Leader mRNA
[LEADER](#) [3'Q](#) [LEADER](#) [REF](#) [LEADER](#) [SCI](#) [LEADER](#) [KX](#) [LEADER](#) [WHY](#)

Samples Database
[SAMPLES](#) [CONSENSUS](#) [MATRIX](#) [ALIGNED](#) [FEATURES](#) [PROFILE3](#) [PROFILE](#) [LIST](#)
[SAMPLE](#) [PHO](#)

Activity Research Database
[ACTIVITY](#) [REFERENCE](#) [SCIENTIFIC](#) [WEIGHT](#) [KNOWLEDGE](#) [PROPERTY](#)

ProteinEDSrel1
[ENPDD](#)

Selex Databases
[SELEX_DB](#) [SELEX_EX](#) [BB](#) [SELEX_TOO_8](#)

Artificial selected protein/protein databases (AEPD)
[AEPD](#) [ALIGN](#) [ASED](#) [RFF](#)

Список вопросов к лекции:

- 2) поиск гомологии и выравнивание генетических текстов**
(Понятие дот-матрицы и ее элементов, основные методы поиска гомологий)

- 2) статистический анализ генетических текстов (структура повторов, классификация повторов)**
- 3) предсказание кодирующих участков генов и открытых рамок считываания (постановка задачи выбора оптимального варианта белок-кодирующей последовательности)**
- 4) предсказание функциональных сигналов (функциональных сайтов и регуляторных районов) (понятие консенсуса и паттерна)**
- 5) анализ вторичной структуры РНК**
- 6) анализ аминокислотных последовательностей**
- 7) филогенетические сравнения**
- 8) ДНК-чипы**
- 9) Интернет-навигация**