Выравнивание последовательностей

17 сентября 2025

Так выглядит выравнивание

```
82 KFKVFDKEFTPQQISAFILQKIKKDA-EAFLGEPVNEAVITVPAYFNDNQR 131
DNAK THEAC
          82 KYKIFGKEYTPQQISAFILQKIKRDA-EAFLGEPVTDAVITVPAYFNDNQR 131
DNAK PICTO
         116 RLRTVAGEKSPVEVSAEILRVLKERAVETLGGEP - EGAVITVPAYFDEAQR 165
HSCA_ACIF2
         116 RLRTVAGEKSPVEVSAEILRVLKERAVETLGGEP - EGAVITVPAYFDEAQR 165
HSCA ACIF5
         132 QATKDAGTIAGFDVKRIINEPTAAALAYGVDKSGKSEKILVFDLGGGTLDV 182
DNAK THEAC
DNAK PICTO
         132 QATKDAGAIAGLNVRRIINEPTAACLAYGIDKLNQTLKIVIYDLGGGTLDV 182
HSCA ACIF2
         166 QATKDAARLAGLNVLRLLAEPTAAAVAYGLDKGSEGI-FAIYDLGGGTFDI 215
         166 QATKDAARLAGLNVLRLLAEPTAAAVAYGLDKGSEGI-FAIYDLGGGTFDI 215
HSCA ACIF5
         183 TIMDFGDGVFQVLSTSGDTRLGGTDMDEAIVNYIADDFQKKEGIDLRKDRS 233
DNAK THEAC
         183 TIMDFGQGVFQVLSTSGDTHLGGTDMDEAIVNFLADNFQRENGIDLRKDHS 233
DNAK PICTO
HSCA ACIF2
         216 SILRLQAGVFEVLATAGDSALGGDDMDHALAEWLMQE - - - - EGGDASDPLW 262
HSCA ACIF5
         216 SILRLQAGVFEVLATAGDSALGGDDMDHALAEWLMQE----EGGDASDPLW 262
         234 AYIRLRDAAEKAKIELSTTLSTDIDLPYITVTNSGPKHIKMTLTRAKLEEL 284
DNAK THEAC
         234 AYIRLRDAAEKAKIELSTVLETEINLPYITATQDGPKHLQYTLTRAKFEEL 284
DNAK PICTO
         263 RRQVLQQ-ARTAKEALSAVAET----MIVLTPSGRAAREIKLSRGRLESL 307
HSCA ACIF2
         263 RRQVLQQ-ARTAKEALSAVAET----MIVLTPSGRAAREIKLSRGRLESL 307
HSCA ACIF5
         285 ISPIVERVKGPIDKALEGAKLKKTEITKLLFVGGPTRIPYVRKYVEDYLG I 335
DNAK THEAC
         285 IAPIVDRSKVPLDTALEGAKLKKGDIDKIILIGGPTRIPYVRKYVEDYFGR 335
DNAK PICTO
         308 IQPVIQRSLPACRRALRDAGLKLDEIEGVVLVGGATRVPAVRAMVEEFFRQ 358
HSCA ACIF2
         308 IQPVIQRSLPACRRALRDAGLKLDEIEGVVLVGGATRVPAVRAMVEEFFRQ 358
HSCA ACIF5
```

Секвенирование миллиардов последовательностей делается главным образом ради выравниваний

Пример из интернета: выравнивание в жизни (спецслужб?)

```
Text1: The caller ident--ified the bomber Text2: The caller ---n-am---ed the bomber Text1: as ------ Yussef Attala, 20, Text2: as 20-year old Yussef Attala----
```

Text1: from the Balata refugee camp near Text2: from the Balata refugee camp near

Text1: Nablus

Text2: Nablus

Бессмысленное

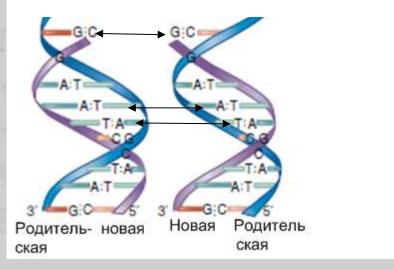
выравнивание

В выравнивании есть содержательные колонки и

бессмысленные

3

Гомология нуклеотидов ДНК



При репликации ДНК большинство нуклеотидов потомка "знают" своего предка в геноме родителя

Иногда из-за ошибок репликации ДНК (или после деления клетки) в ДНК новой клетки появляются небольшие локальные изменения по сравнению с ДНК родительской клетки.

При замене нуклеотида на другой можно указать какой именно нуклеотид родителя изменился!

Гомологичные нуклеотиды двух потомков — те, которые произошли от того же самого нуклеотида общего предка

Гомологичные нуклеотиды разных последовательностей располагают в одной колонке выравнивания

ЛОКАЛЬНУЮ ЭВОЛЮЦИЮ ПОСЛЕДОВАТЕЛЬНОСТИ МОЖНО ОТОБРАЗИТЬ ВЫРАВНИВАНИЕМ

Правильное выравнивание последовательностей ДНК живущих сегодня организмов

Гомологичные нуклеотиды ставим друг под другом

```
ПРЕДОК 1. TAT--GCGAAT-GCCCTGAA
            2. TAT--GCAAAT-GCCCTGAA замена
             3. TAT--GCAAAT-GCTCTGAA замена
               TAT - - GCAAATCGCTCGGAA вставка и замена
    правнук
            5. TAT--GCAAAGCGCTCGGAA замена
праправнук-1
праправнук-2 6. ТАТ - - GCAAA - CGCTCGGAA делеция
живет сейчас а. ТАТ--GCATA-CGC---GAA дел. 3, зам.1.
             b. TATATGCAAAGCGCTCGGAA вставка 2 п.н.
живет сейчас
живет сейчас С. ТАТ--GCAAA--GCGCTGAA дел. 1, зам. 2
```

Задача построения правильного выравнивания сложна и неоднозначно решается

Вот правильное выравнивание с пред. слайда

- a. TAT--GCATA-CGC---GAA
- b. TATATGCAAAGCGCTCGGAA
- C. TAT--GCAAA--GCGCTGAA

Программа выравнивания ориентируется на сходство. Сдвиг даст больше совпадений b. и с.

ПРОГРАММЫ ВЫРАВНИВАНИЯ МОГУТ ОШИБАТЬСЯ!

Эволюция последовательности белка – следствие эволюции кодирующей последовательности

Последовательности большинства белков находится под стабилизирующим отбором – против изменений

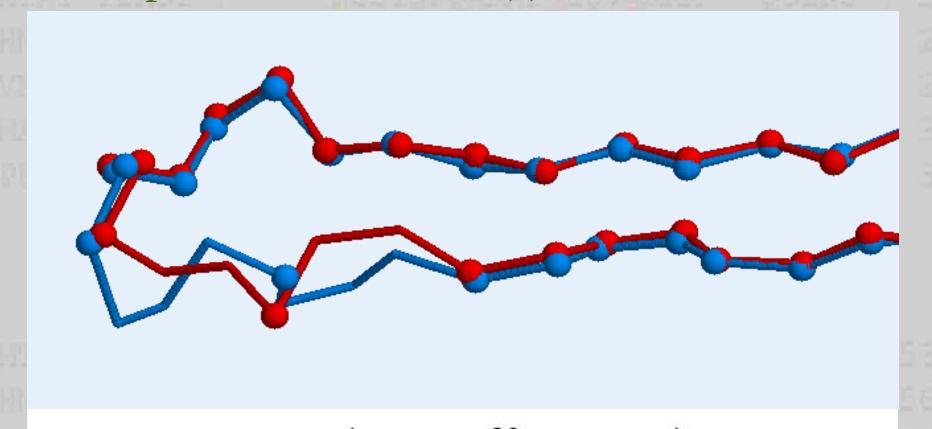
Выравнивание белков

- <u>Теория</u>: аминокислотные остатки белков гомологичны, если их кодоны гомологичны. В выравнивании гомологичные остатки располагают друг под другом
- <u>Практика</u>: гомологичность белков, участков и отдельных остатков выводят из выравнивания, построенного программой на основе сходства букв
- Выравнивать последовательности генов или белков? БЕЛКОВ:
 - в нуклеотидном выравнивании программа может нарушить кодонную структуру
 - в нуклеотидном выравнивании больше мутаций, чем в белковом, из-за синонимичных замен
 - белковое выравнивание учитывает сходство свойств аминокислот

Можно ли проверить правильность выравнивания?

- Программа может построить неправильное выравнивание. Есть ли способы независимой проверки?
- ДНК: нет способов (не считая генноинженерного мутагенеза в лаборатории)
- <u>Белки</u>: можно! если известна 3D структура белка.
 - Если при совмещении полипептидных цепей хорошо совмещаются Сα атомы, то такие остатки можно считать гомологичными
 - Ограничение: Структуры известны для 80 000 белков, а последовательности для 250 000 000
 - Есть базы эталонных выравниваний, построенных по совмещению структур (Balibase и др.). Их используют для сравнения программ выравнивания
- <u>РНК</u>: частично да. Тоже можно учитывать 3D структуру

Белки: совмещение структур и выравнивание последовательностей



Seq_A : LTGYGRWEAEFagnkae--sdtaqqKTrlAFAGLK : 33 Seq_B : LTGYGQWEYNFqgnnsegadaqtgnKTrlAFAGLK : 35 Aligned : AAAAAAAAAAAAAAAAAAAA : 27 Все программы ищут сходство, а сходство — не то же, что гомология!

ПРОГРАММЫ ВЫРАВНИВАНИЯ

Что нужно знать об алгоритме выравнивания двух последовательностей?

- Матрица весов замен
- Штраф за открытие "гэпа"
- Штраф за продолжение "гэпа"
- Вес выравнивания
- Оптимальное выравнивание выравнивание с максимальным весом

Программа успешно выровняет любые две последовательности, даже не гомологичные!

Упражнение: вычислите вес выравнивания

PDB:1osm : KIDGLHYFSD--Dkd : 28

PDB:1hxx : KAVGLHYFSKgnGen : 30

ARNDCQEGHI... A 4-1-2-2 0-1-1 0-2-1-1 R-150-2-310-20-3-2 N -2 0 6 1 -3 0 0 0 1 -3 -3 D -2 -2 1 6 -3 0 2 -1 -1 -3 -4 C 0-3-3-39-3-4-3-3-1-1 Q-1100-352-20-3-2 E-1002-425-20-3-3 G 0-20-1-3-2-26-2-4-4 H-201-1-300-28-3-3

1 2 2 2 1 2 2 1 2 1

gap open —6

gap extension -2

Штраф за продолжение маленький, т.к. делеция нескольких кодонов может произойти как одно событие

Матрица BLOSUM62

```
Ε
                     G
                        Н
                            LKMF
                                       P S T W Y
                          Ι
               -1 -1 0 -2 -1 -1 -1 -2 -1 1 0 -3 -2 0 -2 -1
                  0 -2 0 -3 -2 2 -1 -3 -2 -1 -3 -2 -3 -1
       6 1 -3
                     0
                        1 -3 -3
                                0 -2 -3 -2
                0
                  0
                                           1 0 -4 -2 -3
                  2 -1 -1 -3 -4 -1 -3 -3
                                           0 -1 -4 -3 -3
   -3 -3 -3 -3 -4 -3 -3 -1 -1 -3 -1 -2 -3 -1 -1 -2 -2 -1 -3 -3 -2 -4
0 -1
                5 2 -2
                                   0 - 3 - 1
          0 -3
                        0 - 3 - 2
                                1
                                           0 -1
                  5 -2
                        0 - 3 - 3
                                1 -2 -3 -1
                                           0 -1 -3 -2 -2
                                         0 -2 -2 -3 -3 -1
       0 -1 -3 -2 -2 6 -2 -4 -4 -2 -3 -3 -2
       1 -1 -3
                  0 -2 8 -3 -3 -1 -2 -1 -2 -1 -2 -2 2
I -1 -3 -3 -3 -1 -3 -4 -3
                          4 2 -3
                                  1
                                     0 -3 -2 -1 -3 -1
L -1 -2 -3 -4 -1 -2 -3 -4 -3
                          2 4 - 2
                                  2
                                    0 -3 -2 -1 -2 -1
                  1 -2 -1 -3 -2 5 -1 -3 -1
       0 -1 -3
                                           0 -1 -3 -2
                0 -2 -3 -2
                          1
                             2 -1
                                  5
                                     0 -2 -1 -1 -1 -1
                             0 -3
F -2 -3 -3 -3 -2 -3 -3 -1
                          0
                                  0
                                    6 -4 -2 -2 1
                                                  3 -1 -3 -3 -1 -4
P -1 -2 -2 -1 -3 -1 -1 -2 -2 -3 -3 -1 -2 -4 7 -1 -1 -4 -3 -2 -2 -1 -2 -4
                     0 -1 -2 -2 0 -1 -2 -1 4
       0 -1 -1 -1 -1 -2 -2 -1 -1 -1 -1 -2 -1
                                         1
                                             5 -2 -2 0
W -3 -3 -4 -4 -2 -2 -3 -2 -3 -2 -3 -1
                                     1 -4 -3 -2
                                               11
                      2 -1 -1 -2 -1 3 -3 -2 -2
Y -2 -2 -2 -3 -2 -1 -2 -3
V 0 -3 -3 -3 -1 -2 -2 -3 -3
                         3
                            1 -2 1 -1 -2 -2 0 -3
                                                  -1 4 -3 -2 -1 -4
                  1 -1
                        0 -3 -4 0 -3 -3 -2 0 -1
          4 -3
                0
                                               -4 -3 -3
                  4 - 2
                        0 -3 -3
                               1 -1 -3 -1
                                          0 -1 -3 -2 -2
X 0 -1 -1 -1 -2 -1 -1 -1 -1 -1 -1 -1 -1 -2 0 0 -2 -1 -1 -1 -1 -4
```

Парное выравнивание

- Локальное (алгоритм *Smith Waterman*) находит наиболее сходные участки двух последовательностей.
- Глобальное (Needleman Wunsch) по всей длине последовательностей: пригодно только для последовательностей, гомологичных по всей длине!

Выравнивание последовательностей

Номер столбца MTA1 YEAST: ----KSSISPOARAFIEQVERK---QSINS KPYRGHRF<mark>T</mark>KENVRI<mark>LESWEAK</mark>NIENPYLDT MAT2 YEAST : L13T.F. 40 60 KEKEEVAKKCGITPLQVRVWFINKRMRSK-MTA1 YEAST : 53 KGLENLMKNTSLSRIQIKNWVSNRRRKEKT MAT2 YEAST N4R 63 6Q64 W Название последовательности Функционально консервативная позиция Консервативный остаток Номер последнего в строке остатка ИЗ ЭТОЙ ПОСЛЕДОВАТЕЛЬНОСТИ

Множественное выравнивание

- Как правило, глобальное; программы
 - Muscle
 - MAFFT
 - Clustal Omega
 - T-coffee
 - •
- Локальное находит лучшие т.н. "мотивы" короткие похожие участки
 - MEME
 - •

Примеры выравниваний

РНК-зависимые РНК полимеразы пикорнавирусов

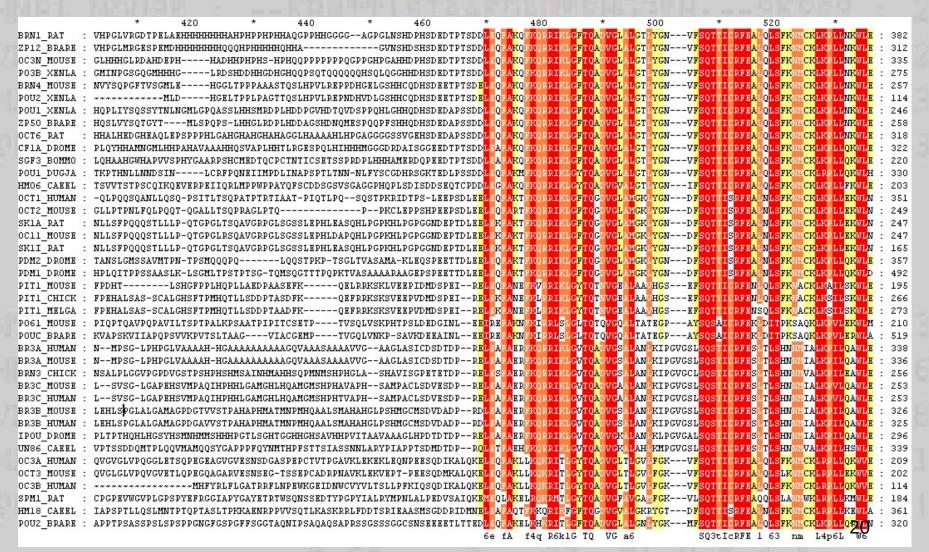
```
20
                                 340
                                                      360
             LIYDEILNTTICLANGMVIRKNVGNNSGOPSTVVDNTLVLMTAFLYAYIHKTG
                                                                       331
POL1 BAYMG
POLN SOUV3
             VVAQDLLAPSEMDVGDYVIRVKEGLPSGFPCTSQVNSINHWLITLCALSEVTG
                                                                       321
             IVAEDLLSPSVMDVGDFKISINEGLPSGVPCTSQWNSIAHWLLTLCALSEVTN
                                                                       321
POLN LORDV
             SACATLKSNPIGIFNGVAFKVAGGLPSGMPLTSIINSLNHCLMVGSAVVKALE
                                                                       318
POLN SMSV1
             TALINTIIYSKHLLYNCCYHVCGSMPSGSPCTALLNSIINNVNLYYVFSKIFG
                                                                       329
POLG HPAV2
                                                                       319
             AEYLRSLAVSRHAYEDRRVLIRGGLPSGCAATSMLNTTMNNVIIRAALYLTYS
POLG TMEVB
                                                                       315
POLG CXA9/
             TNYIDYLCNSHHLYRDKHYFVRGGMPSGCSGTSIFNSMINNIIIRTLMLKVYK
                                                                       314
POLG CXA21
             VDYIDYLNHSHHLYKNKTYCVKGGMPSGCSGTSIFNSMINNLIIRTLLLRTYK
POL1 GFLV/
             KNLLLAICGRLSICGNQVYATEAGIPSGCALTVVLNSIFKELLMRYCFKKIVP
                                                                       346
                                     q pSG
```

Фрагменты геномов двух видов бруцелл

```
* 20 * 40
AE008917_1899844_1899896 : tttagaaa<mark>ttc</mark>cagagcggttccggttaaaacggaatcgttg : 42
AE008917_686288_686337 : tttag---ttctggagcggttcctgttttaacagaatcgttg : 39
```

AE008917_1899844_1899896 : <mark>ga</mark>g<mark>ccgctcta</mark> : 53 AE008917_686288_686337 : ga<mark>ccgctcta</mark> : 50

Пример: выравнивание POU-белков. Блоки достоверного выравнивания. Домены.



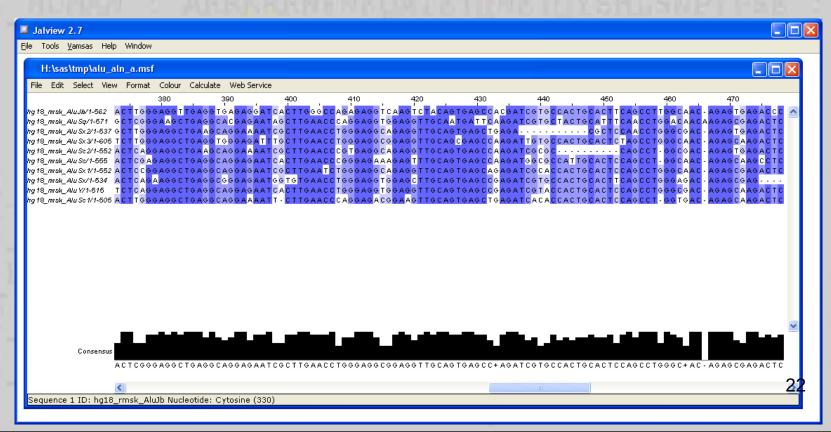
Продолжение



Биологически осмысленное выравнивание может быть в одной части выданного программой выравнивания и не быть в другой! 21

Множественное выравнивание

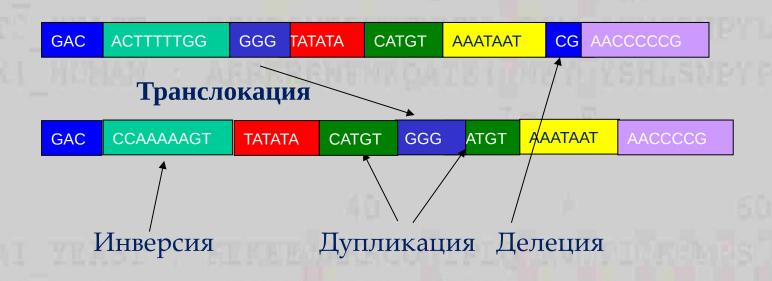
• Программа JalView https://www.jalview.org/ позволяет визуализировать выравнивание. Она же может послать ваши последовательности на один из серверов, делающих выравнивания.



Кроме точечных изменений редко, но происходят крупные перестройки ДНК (или РНК – для РНК вирусов)

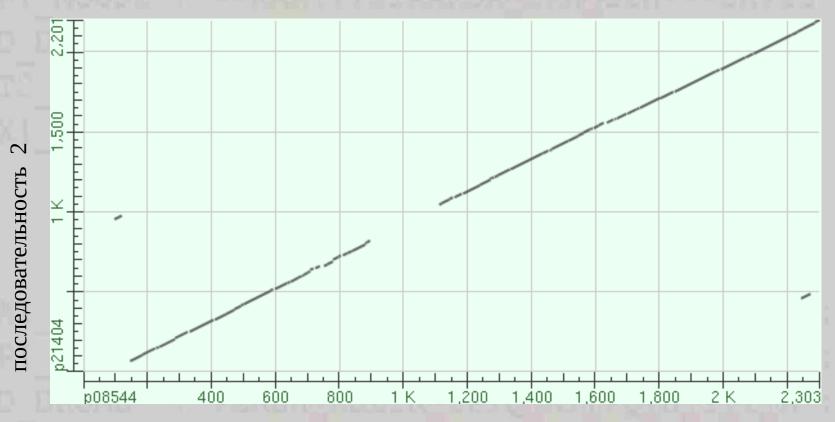
ГЛОБАЛЬНЫЕ ПЕРЕСТРОЙКИ ГЕНОМОВ

Крупные перестройки ДНК

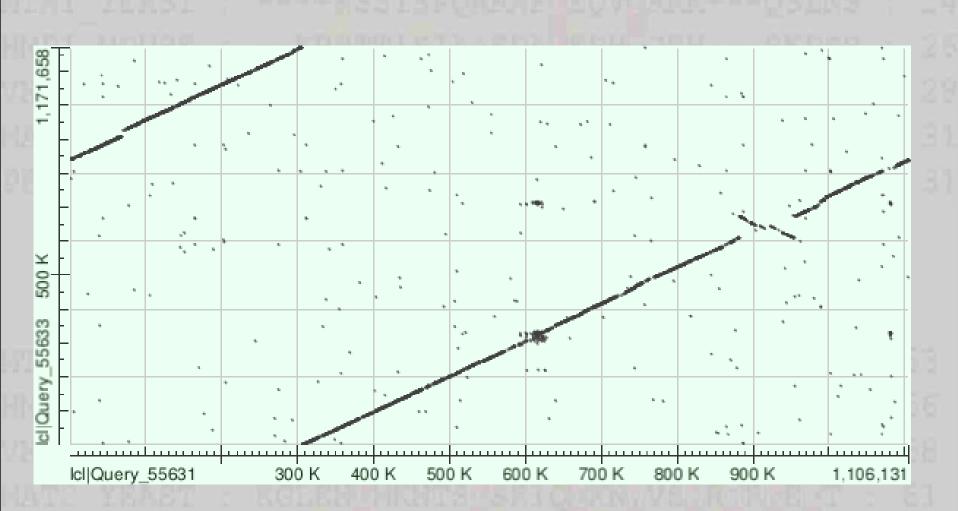


Участки могут состоять из сотен, тысяч и миллионов пар нуклеотидов

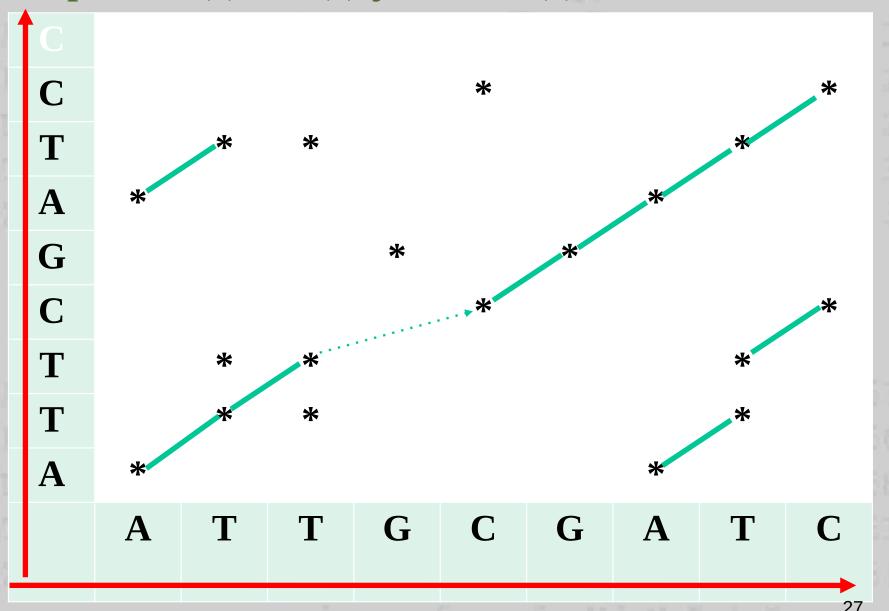
Карта локального сходства позволяет описывать крупные перестройки геномов



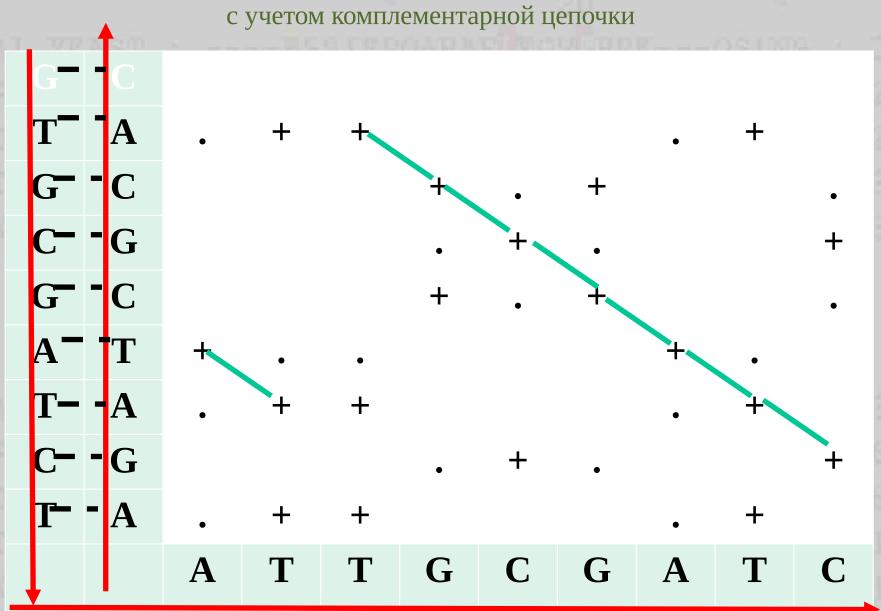
последовательность 1



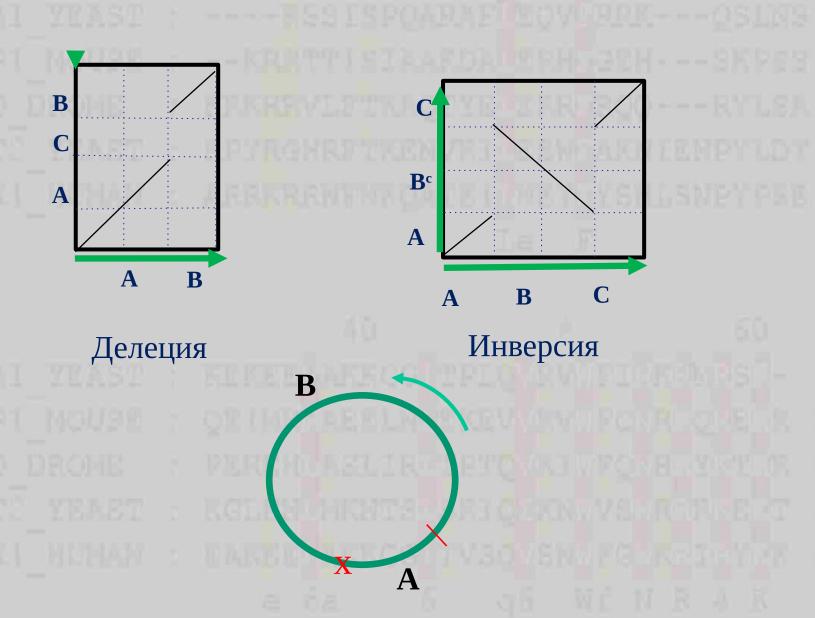
Карта сходства двух последовательностей



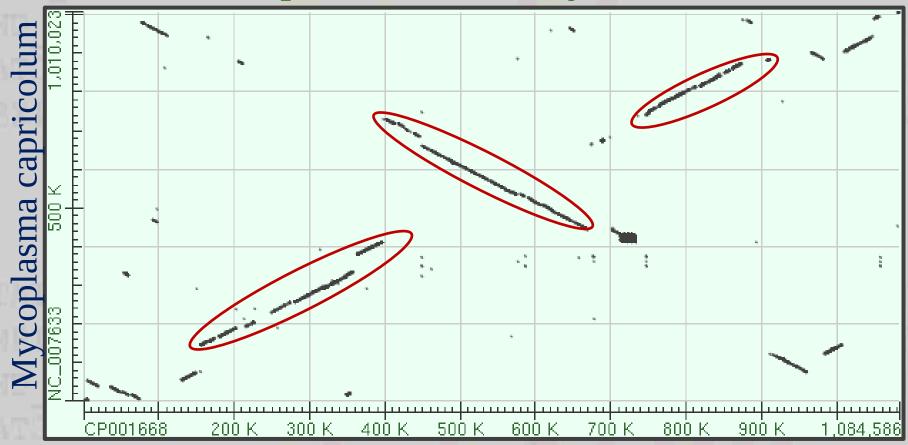
Карта сходства двух последовательностей с учетом комплементарной цепочки



Крупные эволюционные события на карте локального сходства. На примере двух бактериальных хромосом



Карта локального сходства геномов M.capricolum и M.mycoides



Mycoplasma mycoides

Pfam –одна из популярных БД семейств доменов белков и их выравниваний



SEARCH HOME













architectur sequences interaction

species

🏠 🕶 🔝 🕶 🚔 🕶 Page 🕶 Safety 🕶 Tools 🕶 🕡

Family: *Pico_P2A* (PF00947)

Summary

Domain organisation

Alignments

HMM logo

Trees

Curation & models

Species

Interactions

Structures

Summary

Picornavirus core protein 2A (Add annotation

This protein is a protease, involved in cleavage of the polyprotein.

Literature references

1. Petersen JF, Cherney MM, Liebig HD, Skern T, Kuechler E, James MN; , EMBO J 1999;18:5463-5475.: The structure of the 2A proteinase from a common cold virus: a proteinase responsible for the shut-off of host-cell protein synthesis. PUBMED: 10523291 [™]

InterPro entry IPR000081

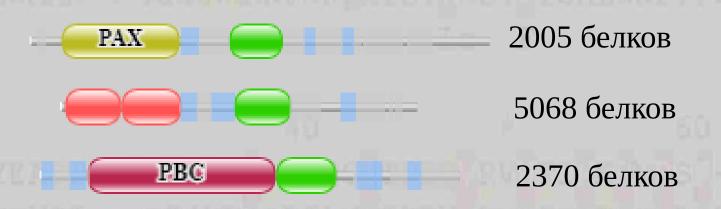


Example structure

PDB entry 2hrv: 2A CYSTEINE PROTEINASE FROM HUMAN RHINOVIRUS 2 31 View a different structure:

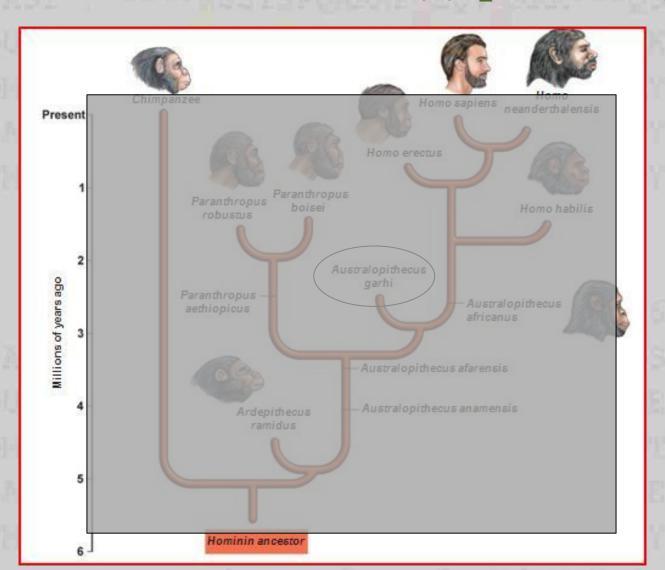
Крупные перестройки встречаются и в последовательностях белков

Примеры доменных архитектур с гомеодоменом (зеленый)

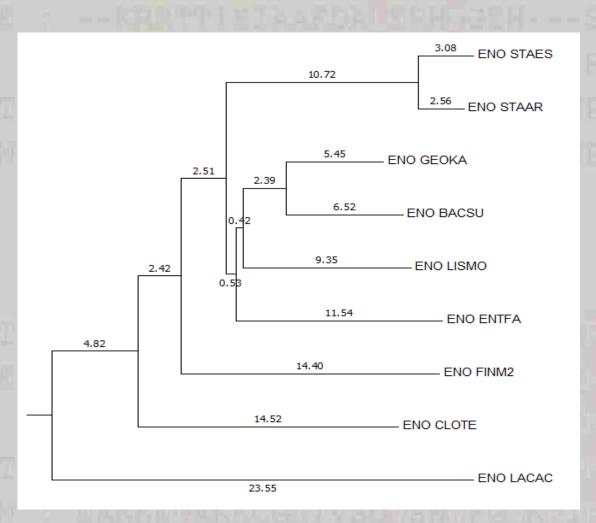


Молекулярная филогения

Филогенетическое дерево видов

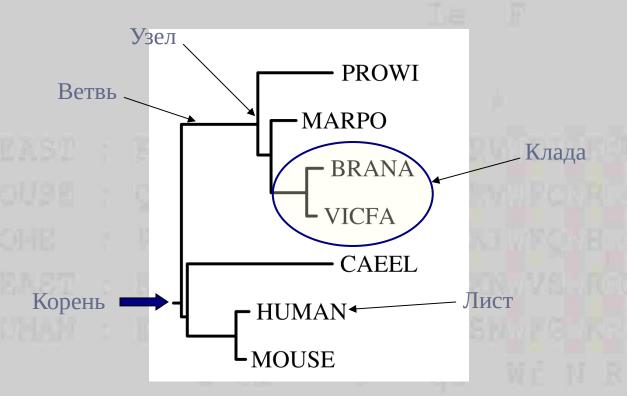


Филогенетическое дерево последовательностей энолаз из фирмикут



Описание структуры дерева (терминология)

- <u>Узел (node)</u> точка разделения предковой последовательности (вида, популяции) на две независимо эволюционирующие. Соответствует внутренней вершине графа, изображающего эволюцию.
- <u>Лист (leaf)</u> реальный (современный) объект; внешняя вершина графа.
- Ветвь (branch) связь между узлами или между узлом и листом; ребро графа.
- **Корень (root)** гипотетический общий предок.
- **Кла́да** группа организмов, которые являются потомками единственного общего предка и всех потомков этого предка.



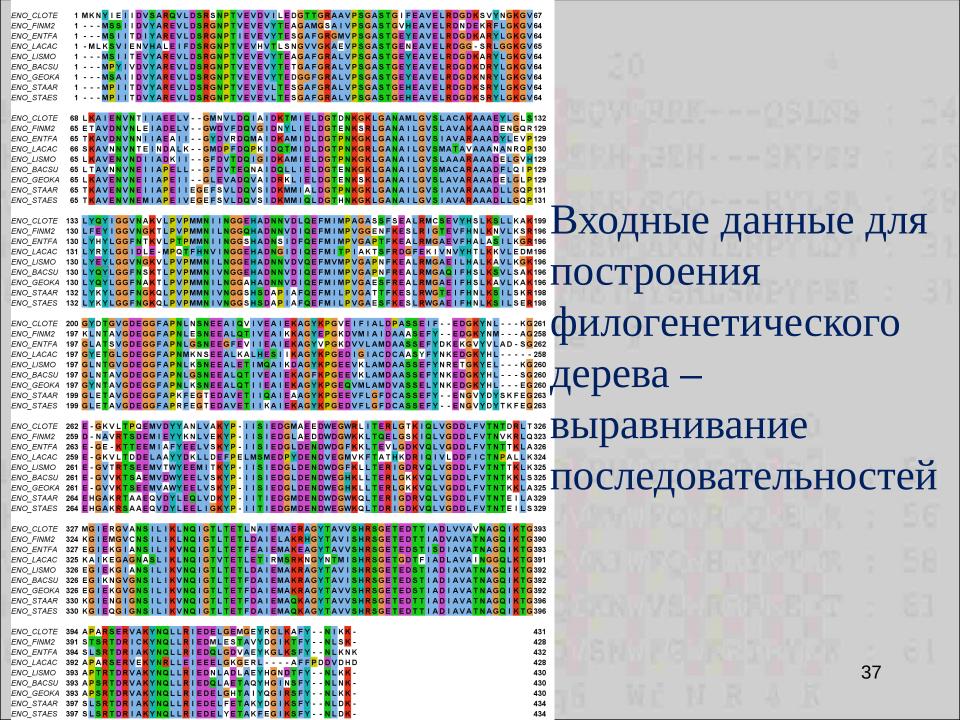
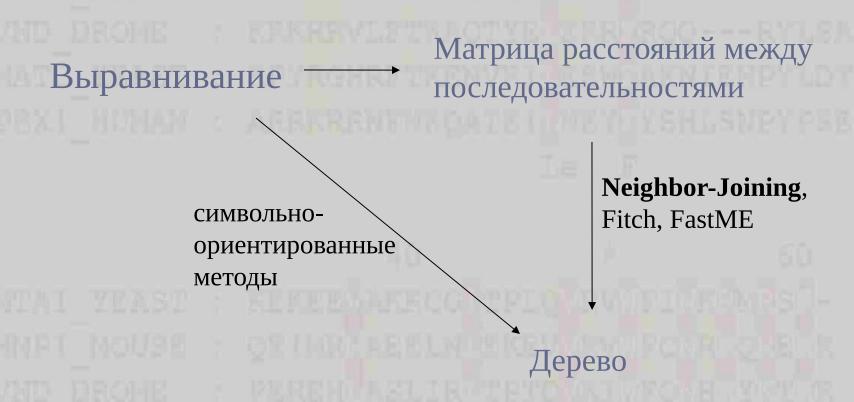


Схема алгоримов построения деревьев



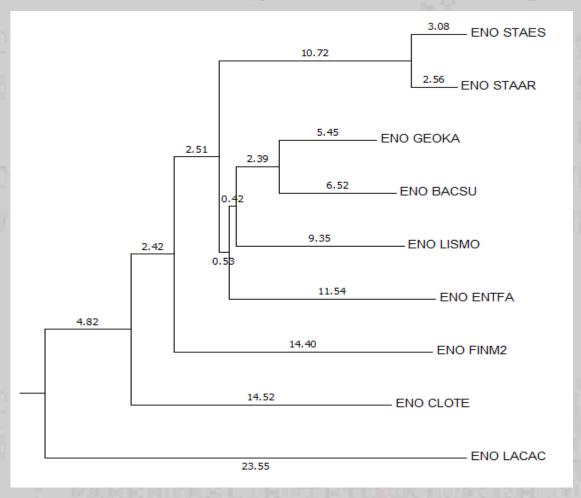
Матрица расстояний

[0]	MUSDO	CHICK	BOVIN	HUMAN
MUSDO	0	9.5	8.9	9.2
CHICK	9.5	0	3.4	2.8
BOVIN	8.9	3.4	0	1.7
HUMAN	9.2	2.8	1.7	0

Матрица расстояний

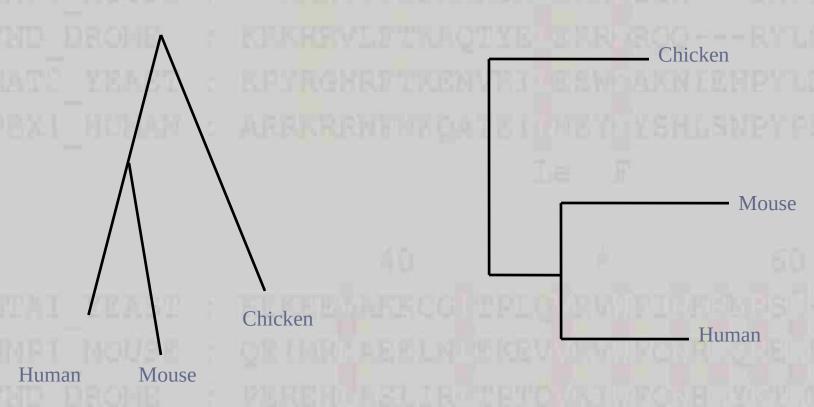
```
ENO CLOTE/
                       0.422314
                                  0.412416
                                             0.703889
                                                       0.374250
                                                                  0.356587
            0.000000
                                                       0.345432
ENO FINM2/
            0.422314
                       0.00000
                                  0.397118
                                             0.748527
                                                                  0.328530
ENO_ENTFA/
            0.412416
                       0.397118
                                  0.00000
                                             0.756866
                                                       0.240851
                                                                  0.271009
ENO LACAC/
            0.703889
                       0.748527
                                  0.756866
                                             0.00000
                                                       0.732893
                                                                  0.710658
ENO_LISMO/
            0.374250
                       0.345432
                                  0.240851
                                             0.732893
                                                       0.00000
                                                                  0.212414
                                  0.271009
                                                       0.212414
                                                                  0.00000
ENO BACSU/
            0.356587
                       0.328530
                                             0.710658
```

Расстояние между листьями на дереве



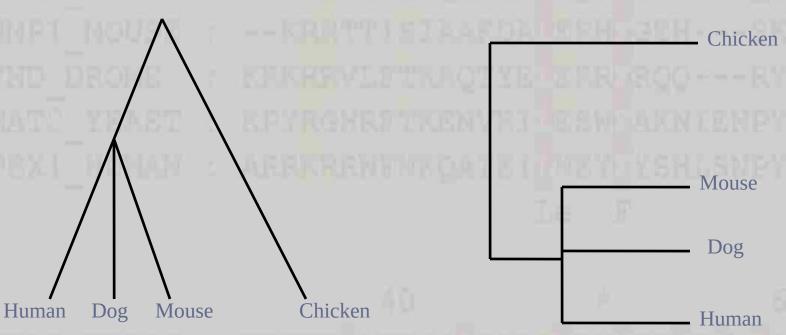
Задача алгоритма – построить такое дерево, чтобы расстояния между листьями на дереве было примерно таким же, как в матрице расстояний

«Молекулярные часы»: всегда идут, но иногда неточно

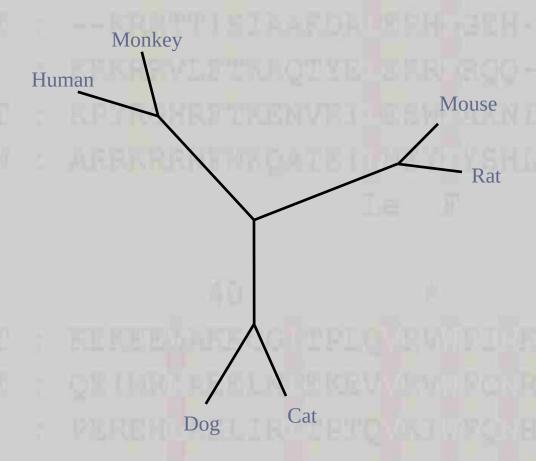


Когда хотят отразить разное число мутаций, произошедших на пути от общего предка, получается что-то вроде такого.

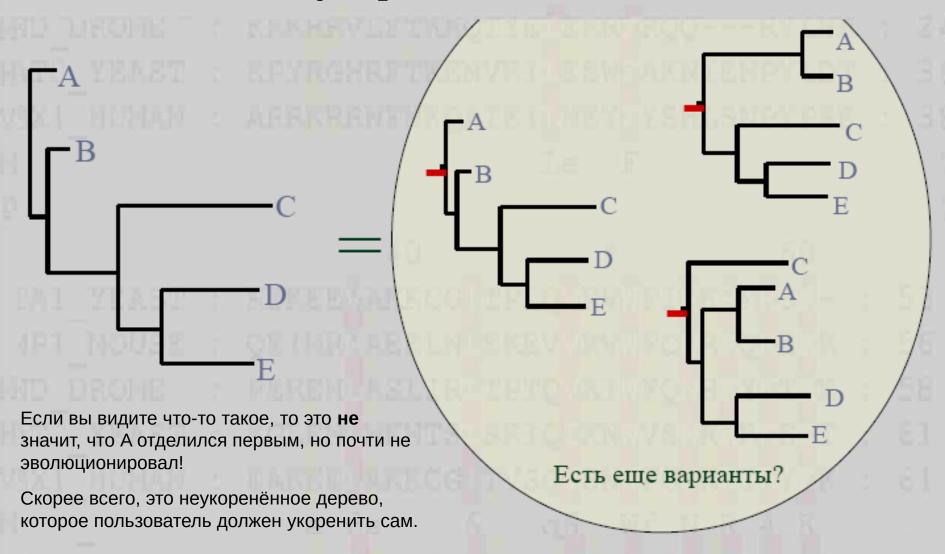
Небинарное дерево



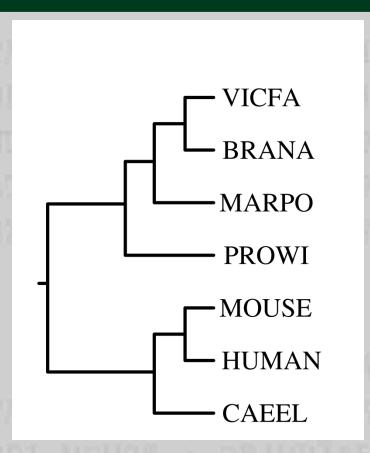
Неукоренённое дерево



Может быть укоренено многими способами



Скобочная формула



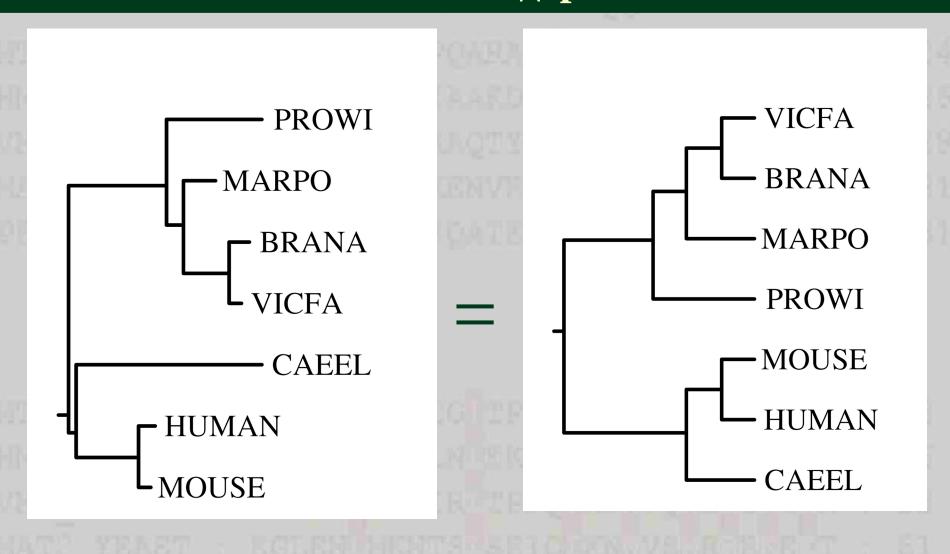
Newick Standard:

((((VICFA:3, BRANA:3):3, MARPO:6):2, PROWI:8):7, ((MOUSE:3, HUMAN:3):3, CAEEL:6):9);

«The reason for the name is that the second and final session of the committee met at Newick's restaurant in Dover, and we enjoyed the meal of lobsters.»

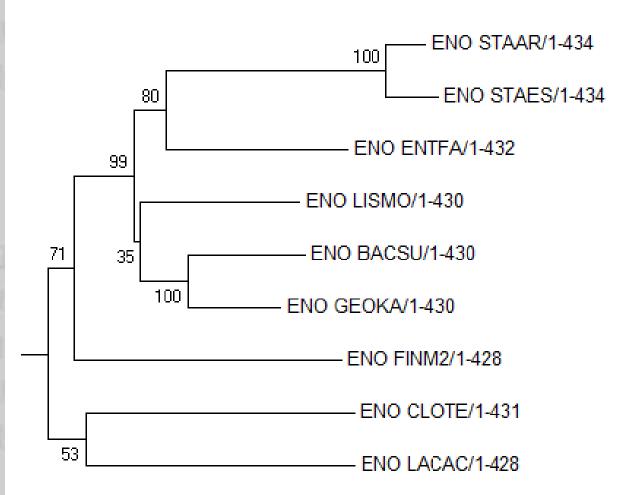
Joseph Felsenstein, http://evolution.genetics.washington.edu/phylip/newicktree.ht/

Топология дерева



Можно ли проверить достоверность дерева? Бутстрэп-анализ

(пример результата)



Эволюция видов и эволюция белков

Когда виды разделяются, то разделяются пути эволюции всех их белков...

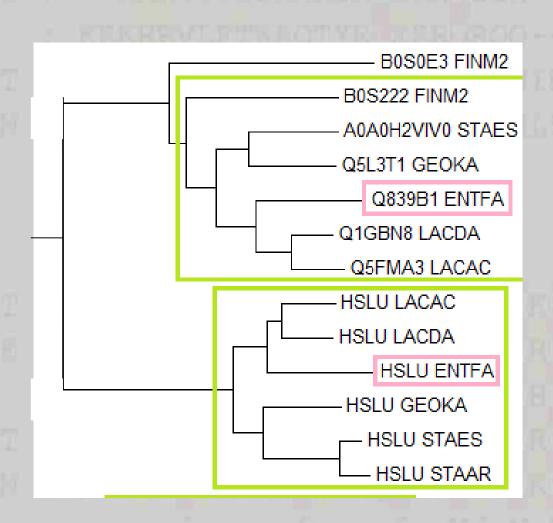
В результате большинству белков одного вида соответствует ортолог в другом виде.

Ho:

- 1)Бывают дупликации белков без разделения видов: два родственных белка существуют в одном геноме и эволюционируют (почти) независимо такие белки называются паралогами.
- 2)Бывают потери генов.
 - Если в двух видах потерялись по одному белку из пары паралогов, то может получиться, что общий предок белков, которые выглядят как ортологи, «жил» существенно раньше, чем общий предок видов.
- 3)Бывает, что два белка объединяются в один многодоменный, и наоборот.

Поэтому правильнее говорить об эволюции белковых доменов.

Ортологи и паралоги. Пример.



Программы

Веб-интерфейс к нескольким программам:

https://ngphylogeny.fr/

Cepвис iTOL https://itol.embl.de/ можно использовать для изображения деревьев