BLAST

Школа биоинформатики, осень 2016

Basic Local Alignment Search Tool

На входе:

- биологическая последовательность (запрос, "query");
- банк биологических последовательностей.

На выходе:

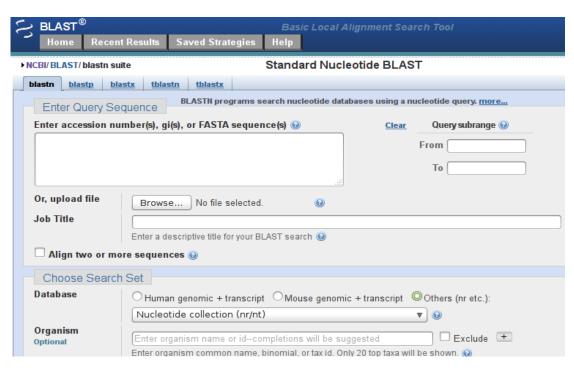
 – лучшие выравнивания запроса с последовательностями из банка.

Важно:

- выравнивания локальные (не глобальные);
- выравнивания не обязательно оптимальны по весу. Быстрый алгоритм позволяет за разумное время искать в большом банке, но может пропустить хорошее выравнивание.

Два воплощения

Online



Локальный (standalone)

Пять видов BLAST

| Программа | Запрос | Банк |
|-------------------------|--------|-------|
| BLASTN | НК | НК |
| BLASTP | Белок | Белки |
| BLASTX | НК | Белки |
| T BLAST N | Белок | НК |
| T BLAST X | НК | НК |

Standalone BLAST

(точнее, BLAST+)

Имеется файл с запросом (или многими запросами) и файл с «банком»

Банком может быть, например, чей-либо геном или протеом или даже «сырые» риды...

Оба файла – в формате Fasta!

Этапы работы:

- форматирование банка (программа makeblastdb)
- собственно поиск (blastp или blastn или blastx или tblastn)

Подробное описание: https://www.ncbi.nlm.nih.gov/books/NBK279690/

Форматирование банка

> makeblastdb -in db.fasta -dbtype prot -out mydb

Входной файл
Тип данных: Название
"prot" или "nucl" нового банка

В результате работы образуются три файла с расширениями .phr, .pin, .psq для типа prot .nhr, .nin, .nsq для типа nucl

BLASTP

> blastp -query adh4_human.fasta -db drome -out adh4.blastp
Файл с Название банка Выходной файл запросом

Выходной файл состоит из четырёх частей:

- справочный заголовок (версия программы, на что ссылаться, размер банка, длина запроса);
- список «хитов» (то есть находок, последовательностей банка, в которых нашлись лучшие выравнивания);
- выравнивания;
- использованные параметры.

Список находок

| Sequences producing significant alignments: | Score (Bits) | E Value |
|---|-----------------|------------|
| ADHX_DROME P46415 Alcohol dehydrogenase class-3 (1.1.1.1) (Alco | 430 | 1e-149 |
| MECR DROME Q9V6U9 Probable trans-2-enoyl-CoA reductase, mitocho | 33.5 | 0.032 |
| SMBT_DROME Q9VK33 Polycomb protein Sfmbt (Scm-like with four MB | 30.0 | 0.52 |
| PDI DROME P54399 Protein disulfide-isomerase (PDI) (dPDI) (5.3 | 28.9 | 0.95 |
| DHE3 DROME P54385 Glutamate dehydrogenase, mitochondrial (GDH) | 28.5 | 1.3 |
| YEMA DROME P25992 Yemanuclein-alpha | 26.6 | 5.3 |
| RM39 DROME Q9VUJ0 39S ribosomal protein L39, mitochondrial (MRP | 25.8 | 7.6 |

Связь E-value и веса в битах:

E-value =
$$mn2^{-B} = mnK \cdot \exp(-\lambda S)$$

где S – обычный вес, m – длина последовательности, n – суммарная длина банка, K и λ зависят от матрицы замен и штрафов за гэпы.

Упражнение: выразить B через S, K, λ .

Смысл: хочется, чтобы B имел примерно один смысл для разных матриц замен и штрафов за гэпы.

Выравнивание

```
> MECR DROME Q9V6U9 Probable trans-2-enoyl-CoA reductase, mitochondrial
(1.3.1.38) (Precursor)
Length=357
 Score = 33.5 bits (75), Expect = 0.032, Method: Compositional matrix adjust.
 Identities = 24/83 (29%), Positives = 40/83 (48%), Gaps = 1/83 (1%)
           EAGKPLCIEEVEVAPPKAHEVRIQIIATSLCHTDATVIDSKFE-GLAFPVIVGHEAAGIV
Query
           E + L + E ++ PK ++V ++I+A + D
                                                I K+
Sbjct 33
           EPOEVLOLVEDKLPDPKDNOVLVKILAAPINPADINTIOGKYPVKPKFPAVGGNECVAEV
                                                                         92
Query
      76
                                    98
           ESIGPGVTNVKPGDKVIPLYAPL
                     + G VIPL + L
Sbjct 93
            ICVGDKVKGFEAGQHVIPLASGL
                                    115
```

Query — это последовательность запроса, Subject — последовательность из банка (в данном случае MECR_DROME). Выравнивание локальное, сопоставляет участок 17—98 из запроса с участком 33—115 из находки. Знаком "+" обозначены сопоставления неодинаковых букв, вносящие тем не менее положительный вклад в вес выравнивания (то есть для которых значение из матрицы > 0).

«Хвост» выходного файла: использованные параметры

Lambda K H a alpha 0.320 0.138 0.418 0.792 4.96

Gapped

Lambda K H a alpha sigma 0.267 0.0410 0.140 1.90 42.6 43.6

Effective search space used: 484173366 <

Параметры подсчёта bit score и E-value

Вместо тп

Database: drome.fasta

Posted date: Jan 22, 2014 4:20 PM Number of letters in database: 1,948,893 Number of sequences in database: 3,203

Matrix: BLOSUM62

Gap Penalties: Existence: 11, Extension: 1

Neighboring words threshold: 11 Window for multiple hits: 40

Параметры поиска

Полезные параметры

–outfmt 6 или –outfmt 7табличная выдача (часто бывает удобна для автоматической обработки,

особенно если запросов много)

Кстати: если поместить во входной файл (query) не одну последовательность, а много, то будет выдан результат для всех, причём существенно быстрее, чем при поиске по каждой последовательности отдельно.

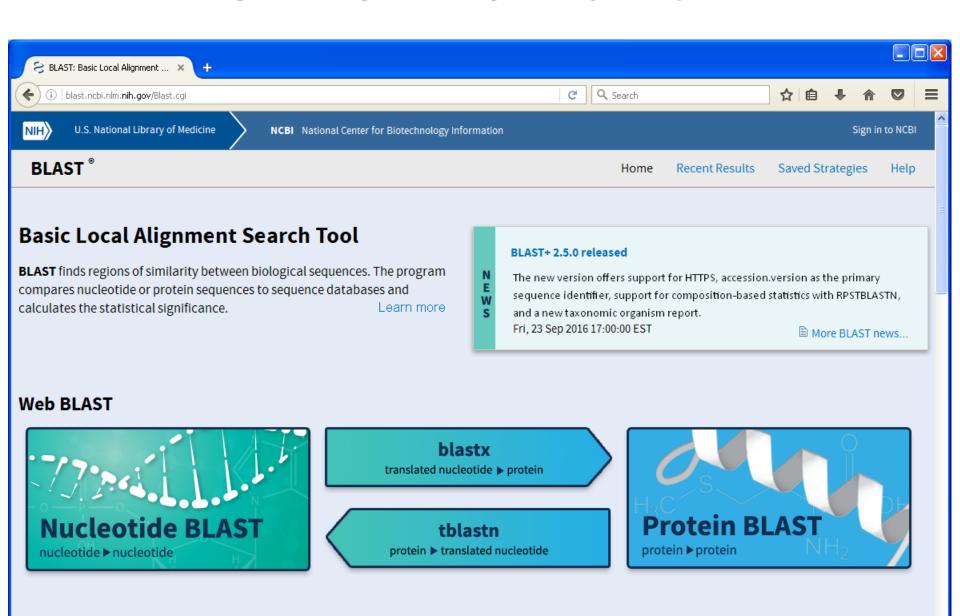
—outfmt 5
выдача в формате xml (для автоматической обработки любителями xml)

—evalue 0.001 ограничение не E-value, чтобы не перегружать выдачу слабыми находками

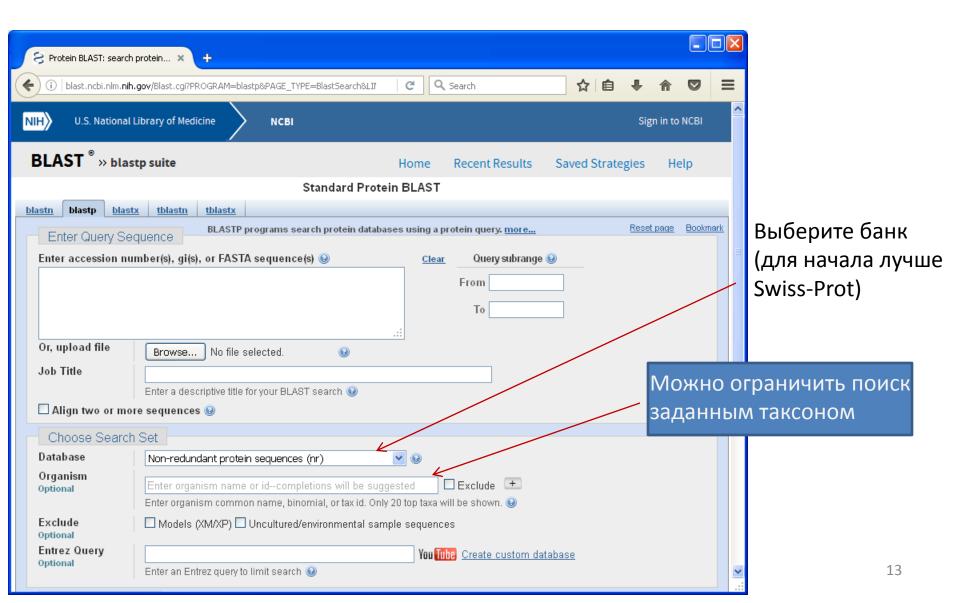
—num_alignments 1000 если в 250, выдаваемые по умолчанию, не помещаются интересные находки

Полный список параметров выдаётся по команде blastp —help, подробности — на https://www.ncbi.nlm.nih.gov/books/NBK279675/

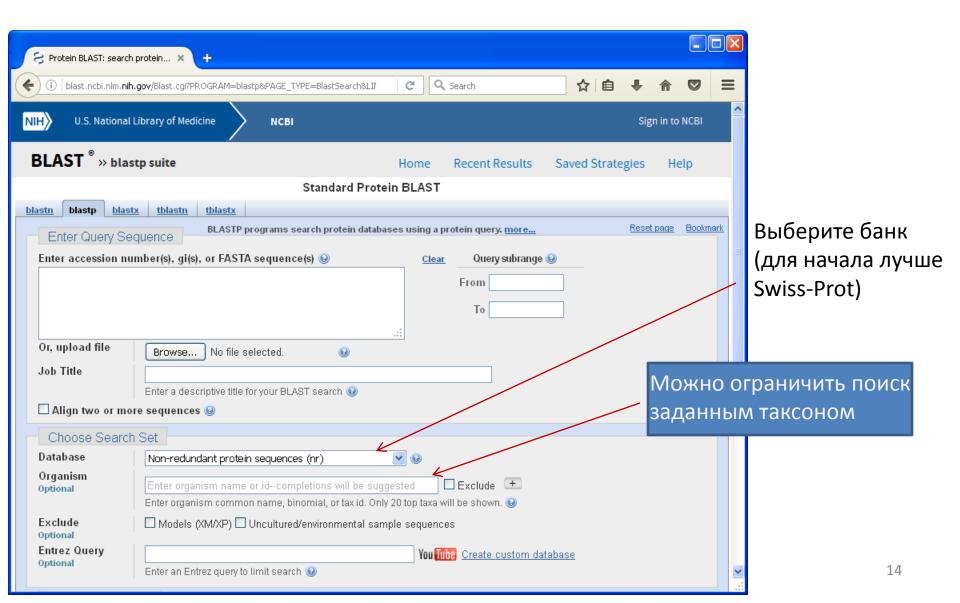
Online BLAST на NCBI



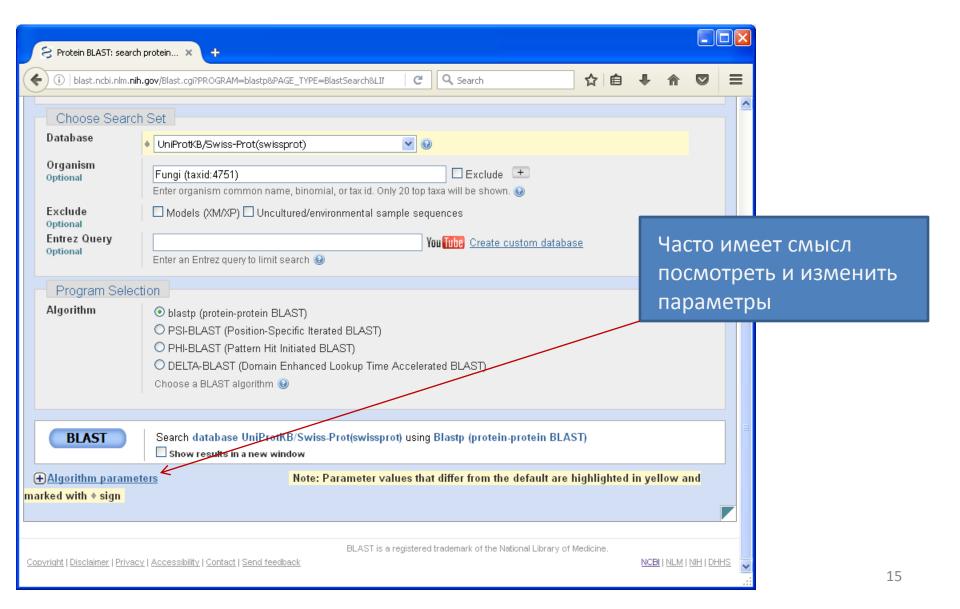
BLASTP Ha NCBI



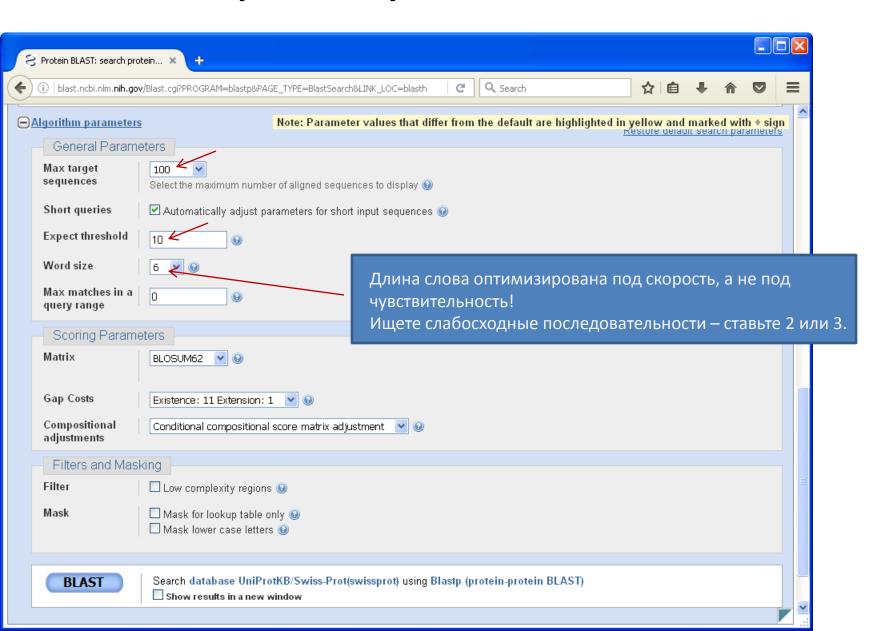
BLASTP Ha NCBI



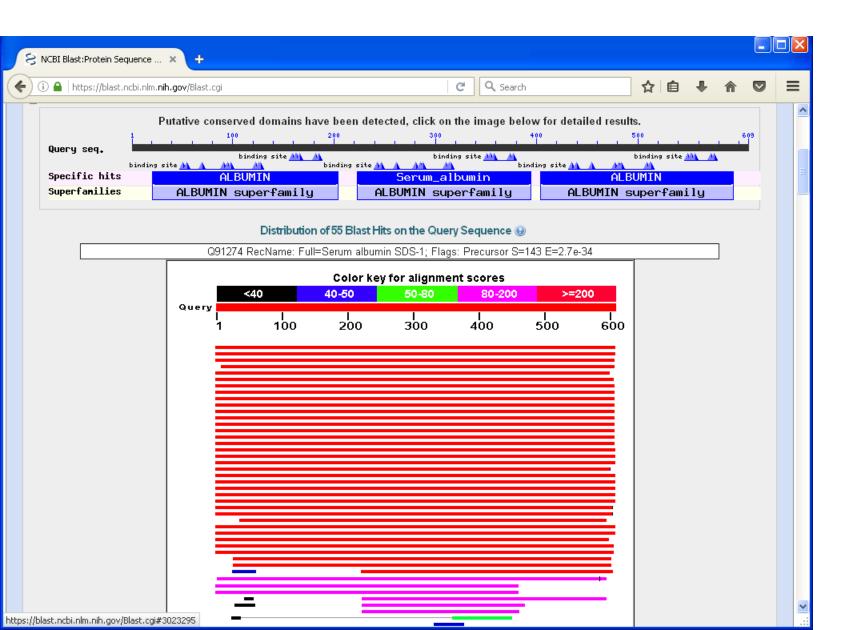
Нижняя часть страницы: параметры



Параметры Online-BLAST



Выдача Online BLAST



TBLASTN и BLASTX

Всё практически так же, как у BLASTP, только белковые последовательности запроса (у TBLASTN) или банка (у BLASTX) возникают как формальные трансляции в шести рамках соответственно нуклеотидного запроса или нуклеотидного банка.

BLASTX – способ черновой аннотации кодирующих участков в вашей нуклеотидной последовательности.

TBLASTN помогает идентифицировать гены гомологов интересующего вас белка в неаннотированном геноме, или даже в сырых ридах, или в коллекции EST.

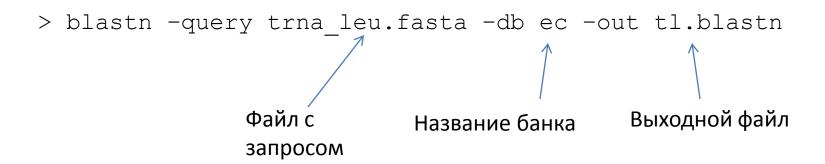
Online – те же проблемы с длиной слова и размером выдачи по умолчанию.

Три подвида BLASTN

| Программа | Затравка* | Цели |
|----------------------------|-----------|---------------------------|
| megablast | 28 нк | Найти запрос в банке |
| discontiguous megablast | 11 нк | Найти близкие гомологи |
| blastn | 11 нк | Найти гомологи |

^{*} по умолчанию

BLASTN



Казалось бы, всё как в blastp, но в данном примере будет запущен **megablast**: очень быстрый, но предназначенный для поиска в банке только последовательностей, совпадающих или почти совпадающих с запросом.

BLASTN

> blastn -task blastn query trna_leu.fasta -db ec -out tl.blastn

В таком варианте уже можно искать гомологи.

Но длина слова по умолчанию 11, что многовато.

Вес совпадения по умолчанию 2, вес несовпадения –3, что позволяет находить лишь гомологи с идентичностью > 70%.

BLASTN

> blastn -task blastn query trna_leu.fasta -db ec -out tl.blastn

В таком варианте уже можно искать гомологи.

Но длина слова по умолчанию 11, что многовато.

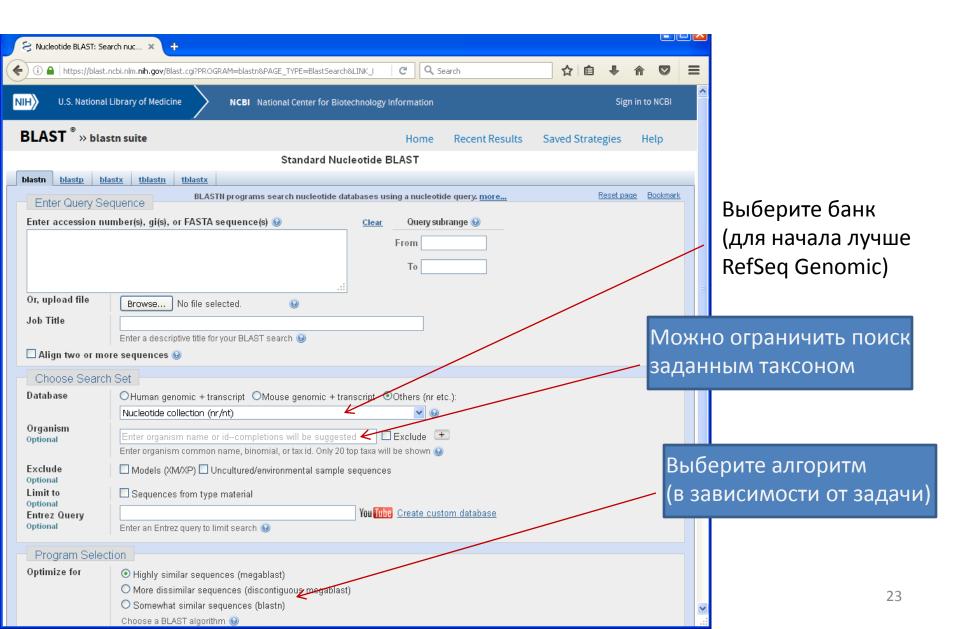
Вес совпадения по умолчанию 2, вес несовпадения –3, что позволяет находить лишь гомологи с идентичностью > 70%.

Поэтому лучше добавить опции: —word_size 4 —penalty —1 —reward 1

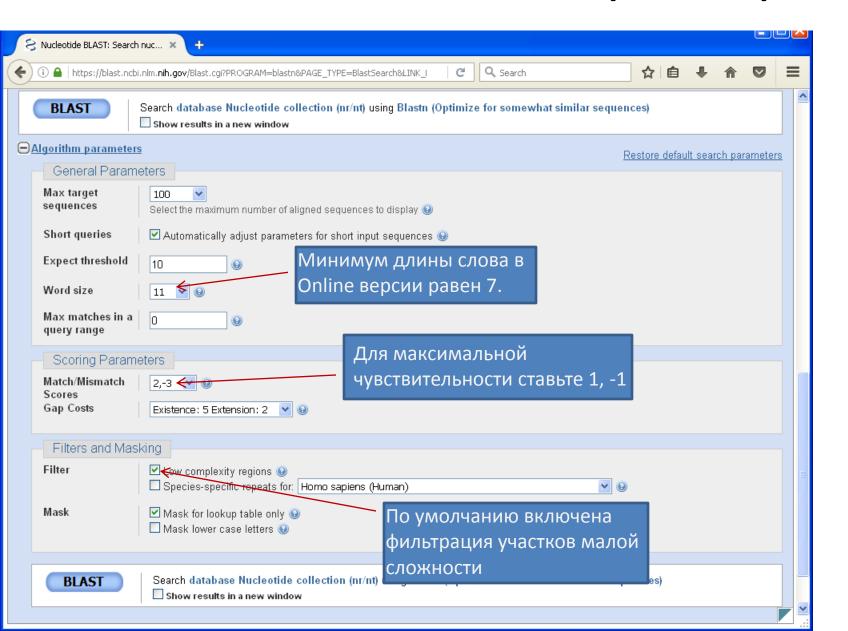
Длина слова 4 заметно замедляет работу. Разумный компромисс между чувствительностью и скоростью – длина слова 7.

Для гомологов с совсем низкой (порядка 50%) идентичностью разумно поставить —penalty —4 —reward 5, но при этом придётся поменять значения —gapopen и —gapextend (см. о том, какие сочетания четырёх параметров допустимы, в https://www.ncbi.nlm.nih.gov/books/NBK279678/).

BLASTN Ha NCBI



BLASTN на NCBI: параметры



Проблема малой сложности

... to appear soon

Удаленный бласт

Имя банка в NCBI

> blastn -task blastn -query query.fasta -db refseq_rna -out blast.out -evalue 0.001 -word_size 7 -outfmt 7 -remote -entrez_query 'arabidopsis[orgn]'

Ограничение по организму возможно только в сочетании с —remote

Опция —remote означает «запускать на сервере NCBI по банку NCBI»

Бласт двух последовательностей

> blastn -task blastn -query seq1.fasta -subject seq2.fasta

Фактически программа локального выравнивания.

В online-версии выдаёт, кроме самих выравниваний, ещё и удобную графическую карту локального сходства двух последовательностей.