

# BLAST: Basic Local Alignment Search Tool

С.А. Спирин, 11 мая 2021



# **BLAST** – алгоритм для нахождения участков локального сходства между последовательностями

Алгоритм сравнивает входную последовательность с последовательностями в базе данных, ищет сходные последовательности в базе данных и оценивает статистическую значимость находок.

# Напоминание: сходство и гомология

Гомология — общность происхождения

- У гомологичных белков можно говорить о парах гомологичных остатков
- В эволюционно правильном выравнивании все остатки в одной колонке гомологичны друг другу

Признак гомологии — сходство последовательностей

- Для выявления сходства последовательности надо выровнять
- Подбирают оптимальное выравнивание, то есть имеющее наибольший вес
- Оптимальное выравнивание существует для любых последовательностей, в том числе негомологичных
- Для двух последовательностей можно рассматривать или глобальное, или локальное выравнивание

# Идея поиска гомологов в банке последовательностей

На входе — последовательность, для которой хочется найти гомологичные («запрос»), и банк

Выравниваем запрос с каждой последовательностью банка, посчитаем веса этих парных выравниваний

Отберём те последовательности банка («находки»), для которых вес **существенно выше, чем мог бы быть по случайным причинам.**

# Почему локальное выравнивание?

Глобальное выравнивание следует применять только в случае заранее известной гомологии последовательностей по всей длине.

Часто у последовательностей гомологичны только отдельные части (примеры: гомеобелки, полипротеины, ...)

Если про белки заранее ничего не известно, то более информативным будет локальное выравнивание. Поэтому именно оно применяется при поиске в банках данных.

# Protein BLAST: поиск гомологов данного белка в банке аминокислотных последовательностей

## Алгоритмы

- blastp
- psi-blast
- phi-blast

Можно использовать:

- из командной строки
- через веб-интерфейс

# Что подаётся на вход программе BLAST?

- Последовательность запроса
- Банк последовательностей
- Параметры:
  - параметры выравнивания: матрица аминокислотных замен, штрафы за гэпы;
  - параметры поиска: длина слова и другие (см. далее);
  - параметры выдачи: максимальное число находок, пороги на качество выравнивания, форма выдачи (обычная, табличная, формат ASN, ...)

# Что выдает BLAST?

Выдача самой программы состоит из четырёх частей:

- заголовок с описанием программы, банка, запроса (query);
- список находок;
- выравнивания запроса с находками;
- несколько строк со статистическими показателями.

Веб-интерфейсы тем или иным способом перерабатывают выдачу программы. Раздел со статистикой обычно не показывается. Часто вставляется графическое изображение находок.

# Выравнивание, выданное BLAST

Sequence ID: Q51368.2 Length: 342

Range 1: 234 to 338

Участок найденного белка,  
попавший в выравнивание

Score:80.9 bits(198), Expect:1e-16,

Method:Compositional matrix adjust.,

Identities:46/115(40%), Positives:63/115(54%), Gaps:15/115(13%)

```
Query 123 SPFENTAPARLTSSTATAATSKPVTSVASGPRALSRNQYPARAQALRIEGQVKVKFDV 182
      +P + PA L S + + KP L + P YP AQA IEG+VKV F +
Sbjct 234 APSGSQGPAGLPSGSLNDS DIKP-----LRMDPPVYPRMAQARGIEGRVKVLFTI 283
```

```
Query 183 TPDGRVDNVQILSAKPANMFEREVKNAMRRWRYEPGKPGSGIVVN-----ILFKI 232
      T DGR+D++Q+L + P+ MF+REV+ AM +WR+EP G IV FKI
Sbjct 284 TSDGRIDDIQVLESVPSRMFDREVRQAMAKWRFEPRVSGGKIVARQATKMFFFKI 338
```

Отображение консервативности: между одинаковыми буквами ставится эта же буква, между сходными (positive) — знак +

Query = запрос, Sbjct (то есть Subject) — найденная последовательность

# Выравнивание, выданное BLAST

Sequence ID: Q51368.2 Length: 342 ← Длина найденного белка  
Range 1: 234 to 338

Score: 80.9 bits (198) ← Вес в битах, Expect: 1e-16, ← Вес  
Method: Compositional matrix adjust., ← E-value  
Identities: 46/115 (40%), Positives: 63/115 (54%), Gaps: 15/115 (13%)

Query	123	SPFENTAPARLTSSTATAATSKPVTSVASGPRALSRNQPOYPARAQALRIEGQVKVKFDV	182
		+P + PA L S + + KP L + P Y E AQA IEG+VKV F +	
Sbjct	234	APSGSQGPAGLPSGSLNDS DIKP-----LRMDPPVYPRMAQARGIEGRVKVLF TI	283
Query	183	TPDGRVDNVQILSAKPANMFEREVKNAMRRWRYEPGKPGSGIVVN-----ILFKI	232
		T DGR+D++Q+L + P+ MF+REV+ AM +WR+EP G IV FKI	
Sbjct	284	TSDGRIDDIQVLESVPSRMFDREVRQAMAKWRFEPRVSGGKIVARQATKMFFFKI	338

Число совпадений ↑  
Длина выравнивания ↑  
Число сходных букв ↑  
Число символов гэпа ↑

# Словарик BLAST

**Identities** — совпадения

**Positives** — сходные буквы, то есть те, для которых значение матрицы положительно

**Gaps** — знаки гэпа "-" (не индели!)

Для всех трёх приводится их число в виде числителя со знаменателем из длины выравнивания (не длины находки!) и процент от длины выравнивания

**Score** — вес выравнивания. Приводится в двух видах: сначала в битах (см. далее), затем в скобках обычный = сумма значений матрицы по сопоставлениям минус штраф за гэпы

**Expect** — E-value, то есть ожидаемое число выравниваний с тем же или большим весом. Запись вида  $9e-15$  означает  $9 \cdot 10^{-15}$ .

**E-value** – ожидаемое количество случайных находок с таким же и лучшим весом выравнивания, при поиске в той же базе данных, со случайным запросом той же длины и состава, с теми же параметрами на вычисление веса выравнивания.

В выдаче BLAST E-value называется “Expect”

**Чем меньше E-value, тем выше значимость находки.**

E-value зависит от:

- веса выравнивания (чем больше вес, тем меньше E-value);
- размера банка (чем больше банк, тем больше E-value);
- длины запроса (чем длиннее запрос, тем больше E-value);
- параметров, используемых для вычисления веса.

**E-value** – **ожидаемое** количество случайных находок с таким же и лучшим весом выравнивания, при поиске в той же базе данных, со случайным запросом той же длины и состава, с теми же параметрами на вычисление веса выравнивания.

Что означает слово «ожидаемое»?

**E-value** – **ожидаемое** количество **случайных** находок с таким же и лучшим весом выравнивания, при поиске в той же базе данных, со случайным запросом той же длины и состава, с теми же параметрами на вычисление веса выравнивания.

Что означает слово «ожидаемое»?

Формально это то, что называется «математическое ожидание случайной величины».

Случайной величиной в данном случае является **число находок** (*NB! Просьба запомнить!*)

На практике ожидаемое вычисляется как **среднее** по достаточно большому количеству испытаний.

Другое ключевое слово — «случайных». Нам нужно понять, сколько можно ожидать именно случайных, то есть бессмысленных, негомологичных находок, чтобы оценить, насколько надёжно утверждение, что данная находка — действительно гомолог.

# Как посчитать E-value

Прямой способ — вычислительный эксперимент: перемешать буквы в запросе очень много раз, каждый раз запуская BLAST, и посмотреть, сколько в среднем при одном запуске бывает находок с весом выше данного.

Такой способ, естественно, не применяется :)

*Стоит подумать: от чего и как может зависеть число случайных находок*

# Как посчитать E-value

Имеется замечательная теорема (С.Карлина):

$$E\text{-value} = Kmn \cdot e^{-\lambda S}$$

$S$  – Score (вес)

$m$  – длина исходной последовательности

$n$  – размер базы данных (суммарная длина всех последовательностей)

$K$  и  $\lambda$  – две константы

Коэффициенты  $K$  и  $\lambda$  зависят от параметров вычисления веса, то есть матрицы и штрафов за гэпы.

BLAST хранит значения  $K$  и  $\lambda$  для нескольких наборов параметров вычисления веса (их раз и навсегда нашли посредством вычислительного эксперимента).

# Вес в битах

Вес в битах  $B$  зависит от обычного веса  $S$  и параметров вычисления веса. Эта зависимость подобрана так, чтобы

$$E\text{-value} = mn \cdot 2^{-B}$$

$m$  – длина исходной последовательности

$n$  – размер базы данных

(констант  $K$  и  $\lambda$  теперь нет, они “загнаны внутрь  $B$ ”)

Нетрудно подсчитать, что  $B = (\lambda S - \ln K) / \ln 2$

Далее описан интерфейс, установленный на «родине» BLAST: National Center for Biotechnology Information (NCBI) в США, <http://blast.ncbi.nlm.nih.gov/>

### Standard Protein BLAST

[blastn](#) **[blastp](#)** [blastx](#) [tblastn](#) [tblastx](#)

BLASTP programs search protein databases using a protein query. [more...](#)

[Reset page](#) [Bookmark](#)

#### Enter Query Sequence

Enter accession number(s), gi(s), or FASTA sequence(s)

[Clear](#)

Query subrange

From

To

**BLAST results will be displayed in a new format by default**  
You can always switch back to the Traditional Results page.



Or, upload file

Выберите файл

Job Title

Enter a descriptive title for your BLAST search

Align two or more sequences

**ВВОДИМ  
ПОСЛЕДОВАТЕЛЬНОСТЬ**

#### Choose Search Set

Database

Non-redundant protein sequences (nr)

Organism  
Optional

Enter organism name or id—completions will be suggested  exclude  +  
Enter organism common name, binomial, or tax id. Only 20 top taxa will be shown.

Exclude  
Optional

Models (XM/XP)  Non-redundant RefSeq proteins (WP)  Uncultured/environmental samples

**банк**

- Non-redundant protein sequences (nr)
- Reference proteins (refseq\_protein)
- Model Organisms (landmark)
- UniProtKB/Swiss-Prot (swissprot)**
- Patented protein sequences (pataa)
- Protein Data Bank proteins (pdb)
- Metagenomic proteins (env\_nr)
- Transcriptome Shotgun Assembly proteins (tsa\_nr)

#### Program Selection

Algorithm

- Quick BLASTP (Accelerated protein-protein BLAST)
- blastp (protein-protein BLAST)**
- PSI-BLAST (Position-Specific Iterated BLAST)
- PHI-BLAST (Pattern Hit Initiated BLAST)
- DELTA-BLAST (Domain Enhanced Lookup Time Accelerated BLAST)

Choose a BLAST algorithm

**организм (если надо ограничить)**

**BLAST**

Search database nr using Blastp (protein-protein BLAST)

Show results in a new window

**дополнительные параметры**

[+ Algorithm parameters](#)

# Дополнительные параметры

максимальный  
размер выдачи

▼ Algorithm parameters

General Parameters

Max target sequences  Select the maximum number of aligned sequences to display ⓘ

Short queries  Automatically adjust parameters for short input sequences ⓘ

Expect threshold  ⓘ

Word size  ⓘ

Max matches in a query range  ⓘ

Scoring Parameters

Matrix  ⓘ

Gap Costs Existence: 11 Extension: 1 ⓘ

Compositional adjustments  ⓘ

Filters and Masking

Filter  Low complexity regions ⓘ

Mask  Mask for lookup table only ⓘ  
 Mask lower case letters ⓘ

порог на E-value

параметры  
выравнивания

борьба с «участками  
малой сложности»



# Участок малой сложности

Определяется как участок с смещенным составом (biased composition)

- Гомополимерные участки
- Короткие повторы
- Перепредставленность отдельных остатков

- ✓ Может мешать анализу последовательностей
- ✓ Вычисление E-value (параметры  $K$  и  $\lambda$ ) опирается на средние по всем белкам частоты аминокислотных остатков, поэтому на участках малой сложности оно становится некорректным
- ✓ Обычно ведет к ложным предсказаниям гомологии (false positives)
- ✓ Лучше использовать «Compositional adjustment» (по умолчанию включен)

# Выдача BLAST в интерфейсе NCBI

[< Edit Search](#)

[Save Search](#)

[Search Summary](#) ▾

[? How to read this report?](#)

[▶ BLAST Help Videos](#)

[↶ Back to Traditional Results Page](#)

<b>Job Title</b>	<b>P02929:RecName: Full=Protein TonB</b>
<b>RID</b>	<a href="#">9X5D3ACN014</a> Search expires on 04-22 14:41 pm <a href="#">Download All</a> ▾
<b>Program</b>	BLASTP <a href="#">?</a> <a href="#">Citation</a> ▾
<b>Database</b>	swissprot <a href="#">See details</a> ▾
<b>Query ID</b>	<a href="#">P02929.2</a>
<b>Description</b>	RecName: Full=Protein TonB [Escherichia coli K-12]
<b>Molecule type</b>	amino acid
<b>Query Length</b>	239
<b>Other reports</b>	<a href="#">Distance tree of results</a> <a href="#">Multiple alignment</a> <a href="#">MSA viewer</a> <a href="#">?</a>

## Filter Results

**Organism** *only top 20 will appear*  exclude

Type common name, binomial, taxid or group name

[+ Add organism](#)

**Percent Identity**

**E value**

**Query Coverage**

to

to

to

**Filter**

**Reset**

**Descriptions**

[Graphic Summary](#)

[Alignments](#)

[Taxonomy](#)

## Sequences producing significant alignments

[Download](#) ▾

[Manage Columns](#) ▾

Show  [?](#)

select all 10 sequences selected

[GenPept](#)

[Graphics](#)

[Distance tree of results](#)

[Multiple alignment](#)

	Description	Max Score	Total Score	Query Cover	E value	Per. Ident	Accession
<input checked="" type="checkbox"/>	<a href="#">RecName: Full=Protein TonB [Escherichia coli K-12]</a>	471	471	100%	4e-170	100.00%	<a href="#">P02929.2</a>
<input checked="" type="checkbox"/>	<a href="#">RecName: Full=Protein TonB [Salmonella enterica subsp. enterica serovar Typhimurium str. LT2]</a>	313	313	100%	5e-108	83.54%	<a href="#">P25945.2</a>
<input checked="" type="checkbox"/>	<a href="#">RecName: Full=Protein TonB [Klebsiella pneumoniae]</a>	270	270	97%	1e-90	67.08%	<a href="#">P45610.1</a>
<input checked="" type="checkbox"/>	<a href="#">RecName: Full=Protein TonB [Serratia marcescens]</a>	125	125	52%	5e-34	54.69%	<a href="#">P26185.1</a>
<input checked="" type="checkbox"/>	<a href="#">RecName: Full=Protein TonB [Klebsiella aerogenes KCTC 2190]</a>	116	116	25%	1e-30	87.10%	<a href="#">P46383.2</a>
<input checked="" type="checkbox"/>	<a href="#">RecName: Full=Protein TonB [Yersinia enterocolitica]</a>	110	110	48%	4e-28	47.06%	<a href="#">Q05740.1</a>

# Переход к текстовому виду

Чтобы скачать выдачу самой программы (а не её обработку интерфейсом), можно поступить так:

The screenshot shows a BLAST search results page. At the top, there are navigation links: '< Edit Search', 'Save Search', and 'Search Summary'. Below this is a metadata section for the job: Job Title (P02929:RecName: Full=Protein TonB), RID (9X5D3ACN014), Program (BLASTP), Database (swissprot), Query ID (P02929.2), Description (RecName: Full=Protein TonB [Escherichia coli K-12]), Molecule type (amino acid), and Query Length (239). There are also links for 'Distance tree of results', 'Multiple alignment', and 'MSA viewer'.

The 'Filter Results' panel is open, showing options to filter by organism, percent identity, E value, and query coverage. A red arrow points from the 'Filter Results' panel to the 'Download' menu in the 'Sequences producing significant alignments' section.

The 'Sequences producing significant alignments' section shows a list of sequences with checkboxes and a 'Description' column. The 'Download' menu is open, showing options: FASTA (complete sequence), FASTA (aligned sequences), GenBank (complete sequence), Hit Table (text), Hit Table (CSV), **Text**, XML, and ASN.1.

	Total Score	Query Cover	E value	Per. Ident	Accession
<input checked="" type="checkbox"/> RecName: Full=Protein TonB [Escherichia coli K-12]	171	100%	4e-170	100.00%	<a href="#">P02929.2</a>
<input checked="" type="checkbox"/> RecName: Full=Protein TonB [Salmonella enterica subsp. enterica serovar Typhimurium str. LT2]	113	100%	5e-108	83.54%	<a href="#">P25945.2</a>
<input checked="" type="checkbox"/> RecName: Full=Protein TonB [Klebsiella pneumoniae]	270	97%	1e-90	67.08%	<a href="#">P45610.1</a>
<input checked="" type="checkbox"/> RecName: Full=Protein TonB [Serratia marcescens]	125	52%	5e-34	54.69%	<a href="#">P26185.1</a>
<input checked="" type="checkbox"/> RecName: Full=Protein TonB [Klebsiella aerogenes KCTC 2190]	116	25%	1e-30	87.10%	<a href="#">P46383.2</a>
<input checked="" type="checkbox"/> RecName: Full=Protein TonB [Yersinia enterocolitica]	110	48%	4e-28	47.06%	<a href="#">Q05740.1</a>
<input checked="" type="checkbox"/> RecName: Full=Protein TonB [Pseudomonas aeruginosa PAO1]	109	46%	1e-16	40%	



# Текстовая выдача BLAST

RID: 9X6N7P7G016

Job Title:ORF1ab

Program: BLASTP

Query: ORF1ab ID: 1c1|Query\_23045(amino acid) Length: 7050

Database: swissprot Non-redundant UniProtKB/SwissProt sequences

Sequences producing significant alignments:

Description	Max Score	Total Score	Query cover	E Value	Ident	Per.	Accession
RecName: Full=Replicase polyprotein 1ab; Short=pp1ab; AltName:...	12928	12928	100%	0.0	85.90		P0C6X7.1
RecName: Full=Replicase polyprotein 1ab; Short=pp1ab; AltName:...	12882	12882	100%	0.0	85.76		P0C6W2.1
RecName: Full=Replicase polyprotein 1ab; Short=pp1ab; AltName:...	12867	12867	100%	0.0	85.52		P0C6V9.1
RecName: Full=Replicase polyprotein 1ab; Short=pp1ab; AltName:...	12861	12861	100%	0.0	85.50		P0C6W6.1
RecName: Full=Replicase polyprotein 1a; Short=pp1a; AltName:...	7476	7476	61%	0.0	80.46		P0C6U8.1
RecName: Full=Replicase polyprotein 1a; Short=pp1a; AltName:...	7461	7461	61%	0.0	80.20		P0C6F8.1
RecName: Full=Replicase polyprotein 1a; Short=pp1a; AltName:...	7436	7436	61%	0.0	79.80		P0C6F5.1
RecName: Full=Replicase polyprotein 1a; Short=pp1a; AltName:...	7431	7431	61%	0.0	79.71		P0C6T7.1
RecName: Full=Replicase polyprotein 1ab; Short=pp1ab; AltName:...	6323	6323	95%	0.0	48.41		P0C6W5.1
RecName: Full=Replicase polyprotein 1ab; Short=pp1ab; AltName:...	6135	6135	99%	0.0	45.73		P0C6W4.1
RecName: Full=Replicase polyprotein 1ab; Short=pp1ab; AltName:...	5990	6347	99%	0.0	50.39		K9N7C7.1
RecName: Full=Replicase polyprotein 1ab; Short=pp1ab; AltName:...	5608	6235	92%	0.0	55.76		P0C6W1.1
RecName: Full=Replicase polyprotein 1ab; Short=pp1ab; AltName:...	5599	6237	93%	0.0	55.66		P0C6W3.1
RecName: Full=Replicase polyprotein 1ab; Short=pp1ab; AltName:...	5435	5608	93%	0.0	49.37		P0C6W8.1
RecName: Full=Replicase polyprotein 1ab; Short=pp1ab; AltName:...	5434	5606	93%	0.0	49.39		P0C6W7.1
RecName: Full=Replicase polyprotein 1ab; Short=pp1ab; AltName:...	5418	5554	83%	0.0	48.88		P0C6X8.1
RecName: Full=Replicase polyprotein 1ab; Short=pp1ab; AltName:...	5404	5574	93%	0.0	49.26		P0C6X6.1
RecName: Full=Replicase polyprotein 1ab; Short=pp1ab; AltName:...	5397	5555	87%	0.0	48.81		P0C6X9.1
RecName: Full=Replicase polyprotein 1ab; Short=pp1ab; AltName:...	5395	5565	93%	0.0	49.23		P0C6W9.1
RecName: Full=Replicase polyprotein 1ab; Short=pp1ab; AltName:...	5346	5484	93%	0.0	48.37		P0C6Y0.1
RecName: Full=Replicase polyprotein 1ab; Short=pp1ab; AltName:...	5305	5483	91%	0.0	48.52		P0C6X3.1
RecName: Full=Replicase polyprotein 1ab; Short=pp1ab; AltName:...	5301	5477	91%	0.0	48.52		P0C6X4.1
RecName: Full=Replicase polyprotein 1ab; Short=pp1ab; AltName:...	5299	5474	91%	0.0	48.53		P0C6X2.1
RecName: Full=Replicase polyprotein 1ab; Short=pp1ab; AltName:...	5267	5435	87%	0.0	48.86		P0C6X0.1
RecName: Full=Replicase polyprotein 1ab; Short=pp1ab; AltName:...	4319	4397	72%	0.0	46.73		P0C6W0.1
RecName: Full=Replicase polyprotein 1ab; Short=pp1ab; AltName:...	4282	4358	70%	0.0	46.53		P0C6Y4.1
RecName: Full=Replicase polyprotein 1ab; Short=pp1ab; AltName:...	4266	4344	70%	0.0	46.72		P0C6X5.1
RecName: Full=Replicase polyprotein 1ab; Short=pp1ab; AltName:...	4227	4303	73%	0.0	45.46		P0C6X1.1
RecName: Full=Replicase polyprotein 1ab; Short=pp1ab; AltName:...	4179	4252	72%	0.0	45.66		P0C6Y5.1
RecName: Full=Replicase polyprotein 1ab; Short=pp1ab; AltName:...	4154	4217	71%	0.0	45.69		Q98VG9.2
RecName: Full=Replicase polyprotein 1ab; Short=pp1ab; AltName:...	4151	4216	70%	0.0	45.61		P0C6Y3.1

# Словарик (таблица находок BLAST)

**Max Score:** самый большой из весов (в битах) выравниваний запроса с данной находкой

**Total Score:** суммарный вес (в битах) всех выравниваний запроса с данной находкой

**Query cover:** процент длины запроса, покрытого выравниваниями

**E Value:** в таблице находок это E-value, посчитанное по особой формуле на основе **всех** выравниваний запроса с данной находкой

**Per. Ident:** процент идентичных букв в лучшем (по весу) из выравниваний запроса с данной находкой

# BLAST — эвристический алгоритм

Алгоритмы биоинформатики можно разделить на точные и эвристические.

**Точные** алгоритмы решают какую-либо точно сформулированную формализованную задачу. Пример: алгоритм Нидельмана – Вунша, который для данных последовательностей находит выравнивание с максимальным весом (реализован в программе needle).

**Эвристические** алгоритмы — те, для которых формальную задачу сформулировать нельзя.

BLAST **не гарантирует** нахождение оптимального локального выравнивания. За счёт этого достигается высокая скорость работы. Но теоретически возможно, что BLAST не найдёт в банке вполне достоверный (судя по выравниванию) гомолог.

# Дополнительные параметры

Algorithm parameters

General Parameters

Max target sequences: 100  
Select the maximum number of aligned sequences to display

Short queries:  Automatically adjust parameters for short input sequences

Expect threshold: 10

Word size: 3

Max matches in a query range: 0

Scoring Parameters

Matrix: BLOSUM62

Gap Costs: Existence: 11 Extension: 1

Compositional adjustments: Conditional compositional score matrix adjustment

Filters and Masking

Filter:  Low complexity regions

Mask:  Mask for lookup table only  
 Mask lower case letters

Длина слова

# Длина слова

Одним из параметров BLAST является длина слова (word size).

Чем больше длина слова, тем быстрее работает BLAST, но тем меньше его **чувствительность**. Это означает, что вероятность пропустить гомологи возрастает.

Сейчас на сайте NCBI значение длины слова по умолчанию равно 6, доступны значения 2 и 3.

# Идея алгоритма BLAST

Нам нужно найти в банке последовательности, хорошо выравнивающиеся с последовательностью запроса (хорошо — то есть с большим весом).

Можно было бы это делать алгоритмом Смита – Уотермена, последовательно выравнивая каждую банковскую последовательность с запросом (и такие сервисы существуют, например [ssearch](#) на сайте [ebi.ac.uk](#)). Но при нынешних объёмах банков это работает слишком медленно.

Идея состоит в том, чтобы заранее **проиндексировать** банк.

Индексы вы видели в конце почти любой научной книги, там имеется **алфавитный** список терминов (или, например, латинских названий растений) с указанием страниц, на которых упоминается этот термин.

В случае BLAST индексами служат **слова** заданной длины из букв, встречающихся в наших последовательностях. Например, для белков и при длине слова 3 это AAA, AAC, AAD, ..., YYY, всего  $20^3 = 8000$  слов.

Перед тем, как запускать собственно поиск, создаётся таблица, в которой для каждого слова указано, в какой последовательности банка и в каком месте это слово встретилось.

# Индекс — примерно то же, что алфавитный указатель в книге

## АЛФАВИТНЫЙ УКАЗАТЕЛЬ

*(цифры обозначают номера экспериментов или параграфов)*

- Агрегатное состояние 18, 19.  
Акустический указатель 169.  
Акция 128.  
Амплитуда колебания 162, 191, 196, 197, 211, 217.  
Апериодические колебания 205.
- Балансирование 65, 66, 70.  
Барометр чашечный § 1.  
Батавские слезки 61.  
Биение 217.  
Бифилярный подвес 150, 156, 162, 197, 207.  
Блок 84—86, § 2 — 1, 3, 4.  
Блок ступенчатый § 2—5.  
Болонская колбочка 61.
- Время, деление на равные промежутки 15, 16.  
Время, измерение 13—15, 113, § 3.  
Время падения 120.  
Высота падения 118, 120.  
Вытесняемость жидкости 8, 9, 21, 22.  
Вытесняемость твердых тел 20.
- Гармоническое колебание 191, 196, § 28.  
Градуирование шкалы динамометра 55.  
Грамм § 7.  
Графики 55, 147, 183, 193, 194, 199.  
Грузики с крючками § 2—10.
- Давления, сила 53, 135.  
Дальность полета 118, 122, 157.  
Движение волновое 201.

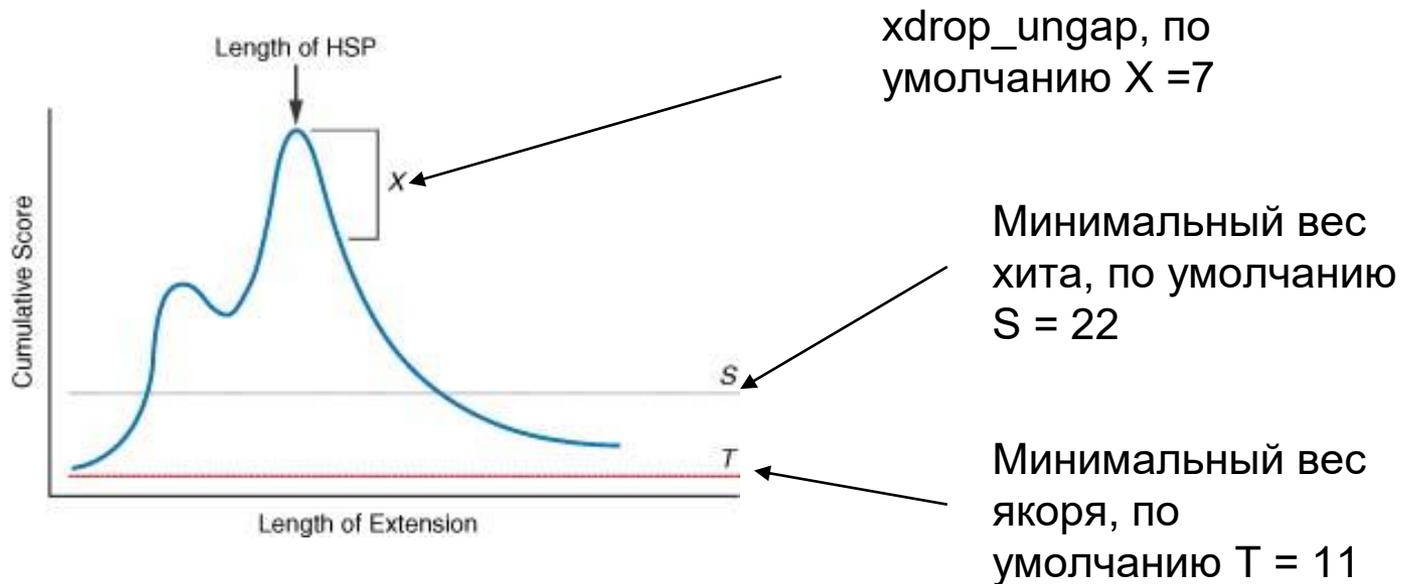
# BLAST: отбор слов

- Два параметра:
  - длина слова  
(word\_size,  $\geq 2$ , по умолчанию 3)
  - порог на сходство слов  
(threshold,  $\geq 0$ , по умолчанию 11)
- Берутся все слова из запроса (query)  
например, из aacddefg будут взяты (при длине слова 3):  
aac, acd, cdd, dde, def, dfg
- В индексах ищутся слова, имеющие сходство со словами из запроса на уровне не менее threshold

# BLAST: от якоря к выравниванию

- Выравнивание начинает строиться, если в запросе есть пара слов на расстоянии, меньшем параметра `window_size` (по умолчанию 40), для которых нашлась пара сходных слов в одной банковской последовательности на том же расстоянии. В результате получаем два якоря — выравнивания длины `word_size`.
- Второй якорь расширяется без гэпов в обе стороны, пока вес не упадёт на заданную величину от максимально достигнутого (по умолчанию этот параметр `xdrop_ungap` = 7 бит)
- Если максимально достигнутый вес больше 22 бит, то соответствующее выравнивание расширяется уже с гэпами (аналогично алгоритму Нидлмана – Вунша). Расширение продолжается, пока вес не упадёт ниже максимально достигнутого на величину, большую `xdrop_gap`, по умолчанию 15 бит

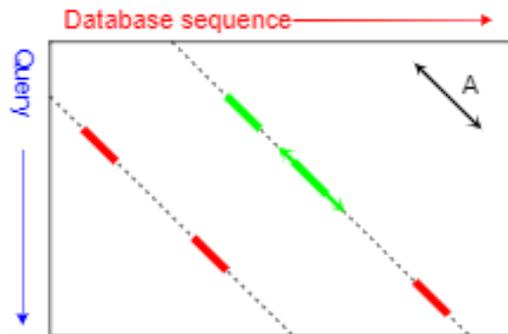
# BLAST: расширение якоря



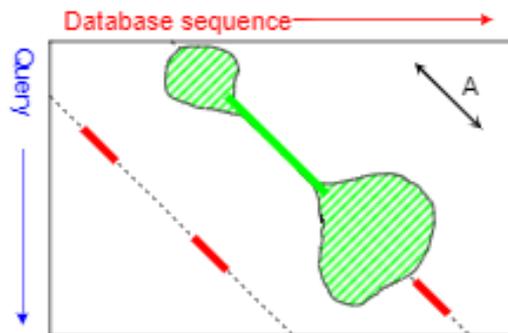
Это схема расширения в одну сторону; после того, как максимальное значение найдено, точно так же расширяем в другую.



## Indexing for Blast (3)



Ungapped extension if:  
2 "Hits" are on the same diagonal but  
at a distance less than A



Extension using **dynamic programming**  
limited to a restricted region  
limited through a **score drop-off**  
threshold

LR, Bezel October 2008



<https://docplayer.net/15013198-Databases-indexation.html>

Автор: Laurent Falquet, SIB

# BLAST: роль длины слова

## (мой эксперимент)

- Вход: последовательность из 466 остатков
- NCBI BLAST (<https://blast.ncbi.nlm.nih.gov/>)
- Область поиска: Swissprot, белки из бактерий
- Параметры, кроме "Word Size", по умолчанию.  
В частности, порог E-value = 10
- Word Size = 6
  - Найдено 16 последовательностей, в них 18 выравниваний
  - 8 выравниваний с  $E < 0,001$
  - Время работы сервиса – менее одной минуты
- Word Size = 2
  - Найдено 69 последовательностей, в них 75 выравниваний
  - 12 выравниваний с  $E < 0,001$
  - Время работы сервиса – около 35 мин

# Вопросы и ответы про BLAST

## **За счёт чего BLAST работает быстро?**

За счёт просмотра не всех возможных выравниваний, а только полученных расширением "затравок". Каждая "затравка" получается из слова длины  $k$  ( $k = 2, 3, \dots, 6$ ), встреченного в запросе, и очень сходного слова из какой-либо банковской последовательности.

"Затравки" находятся очень быстро благодаря предварительной индексации всех слов в банке. В результате индексации для каждого слова указано, в каких местах каких банковских последовательностей это слово встречается.

## **Что может поменяться при изменении параметра "Word size»?**

Чем длиннее слово, тем меньше машинного времени займёт поиск.

Чем короче слово, тем чувствительнее поиск (меньше опасность пропустить хорошее выравнивание).

# Standalone BLAST

BLAST можно установить на своём компьютере  
(а на kodomo он уже установлен)

Предположим, вам нужно найти гомологи белка, чья последовательность — в файле `myprot.fasta`, в протеоме, содержащемся в файле `proteom.fasta` (всё в `fasta`-формате, BLAST других не понимает).

Придётся сначала проиндексировать ваш банк программой `makeblastdb`, подав ей на вход протеом (читайте `makeblastdb -help`)

Эта программа создаст несколько файлов, необходимых для поиска, в том числе тот самый индекс якорей (сразу для всех допустимых длин слов)

После этого можно искать программой `blastp`, указав ей имя файла с запросом и название проиндексированного банка (читайте `blastp -help`, нужные опции: `-query`, `-db`, `-out`)

# Standalone BLAST

Впрочем, можно использовать BLAST и для обычного локального выравнивания двух последовательностей, безо всякой индексации:

```
blastp -query seq1.fasta -subject seq2.fasta -out result.blastp
```

Но имейте в виду, что BLAST и в таком варианте не гарантирует оптимального выравнивания (это **эвристический** алгоритм)! Зато можно быстро выровнять очень длинные последовательности (команде water может не хватить памяти) и получить не одно, а много локальных выравниваний.

(На самом деле в этом варианте BLAST «на ходу» индексирует вторую последовательность)

# BLAST: варианты формата выходного файла

```
-outfmt <String>  
alignment view options:  
  0 = Pairwise,  
  1 = Query-anchored showing identities,  
  2 = Query-anchored no identities,  
  3 = Flat query-anchored showing identities,  
  4 = Flat query-anchored no identities,  
  5 = BLAST XML,  
  6 = Tabular,  
  7 = Tabular with comment lines,  
  8 = Seqalign (Text ASN.1),  
  9 = Seqalign (Binary ASN.1),  
 10 = Comma-separated values,  
 11 = BLAST archive (ASN.1),  
 12 = Seqalign (JSON),  
 13 = Multiple-file BLAST JSON,  
 14 = Multiple-file BLAST XML2,  
 15 = Single-file BLAST JSON,  
 16 = Single-file BLAST XML2,  
 18 = Organism Report
```

0–4 — чтобы смотреть глазами

5–12 — чтобы парсить программами.

6, 7 и 10 можно импортировать в электронные таблицы