# Introduction to Sequence Alignment

Chris Overall

Department of Bioinformatics and Genomics
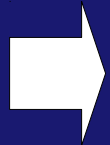
University of North Carolina – Charlotte

# Sequence alignment algorithms

An <u>algorithm</u> is a step-by-step procedure for solving a problem or accomplishing something

Examples ➡️

Finding prime numbers

Sorting numbers from lowest to highest
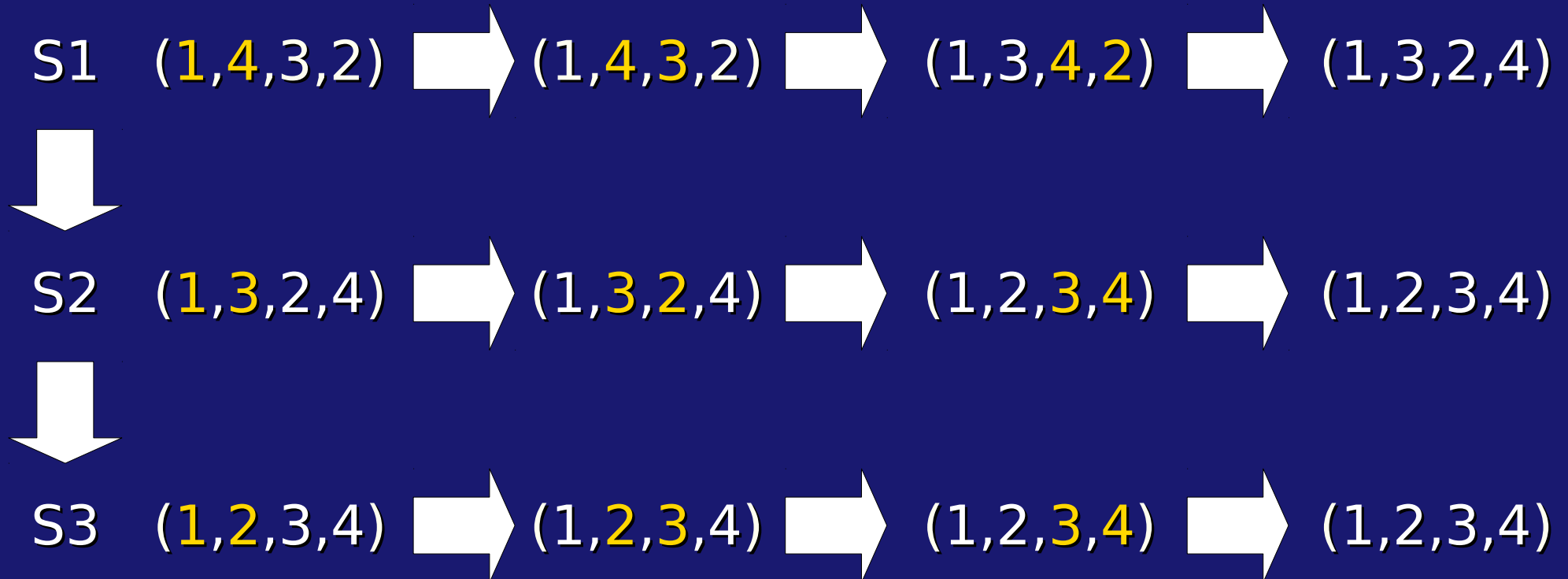
Calculating shortest delivery routes

A cooking recipe

A lab protocol

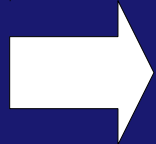The term algorithm usually refers to solving a problem (often mathematical) using a computer

# Bubble sort is a classic algorithm for sorting numbers

Bubble sort

Input: 1,4,3,2 → Output: 1,2,3,4

S1 (1,4,3,2) → (1,4,3,2) → (1,3,4,2) → (1,3,2,4)

S2 (1,3,2,4) → (1,3,2,4) → (1,2,3,4) → (1,2,3,4)

S3 (1,2,3,4) → (1,2,3,4) → (1,2,3,4) → (1,2,3,4)

_Sequence alignment algorithms_ are ways to arrange two or more "words" to find out how similar they are to each other
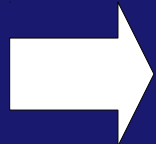
Word ➡ {

Sequence of letters (symbols) using a particular alphabet (set of letters)

Doesn't have to be the English alphabet

Examples ➡ {

FRAGILE

SUPERCALIFRAGILISTICEXPIALIDOCIOUS

1324647597

? 🎁 ✕ ⬜ ? ✈ ⓘ ✦ 🚓 ◀ ⏩ 🚓 ✔ ▢ 🚫

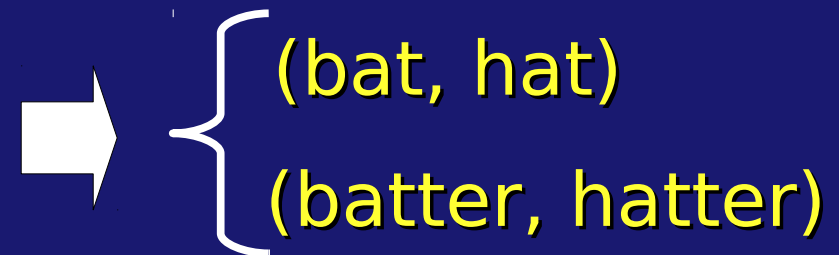**Let's look at a simple example using words from the English alphabet (size = 26)**

Consider these four words: hat, bat, hatter, batter

How similar are bat and hat? Why?

Which word is more similar to bat? Why? ➡ { hat
batter

Which pair of words is more similar to each other? Why? ➡ { (bat, hat)
(batter, hatter)

One way to figure out how similar two words (sequences) are is to find the <u>longest common subsequence</u>

What's the algorithm to do this?

<u>Step 1</u>: Left-align the two sequences

<u>Step 2</u>: Score the alignment

<u>Step 3</u>: Shift the shorter alignment to the right by one character and repeat steps 1-2 until you reach the end of the longer sequence

# Here's a simple example of finding the longest common subsequence

We'll evaluate the sub-alignments using this scoring scheme ➡

$$\begin{cases} \text{Match} = 1 \text{ point} \\ \text{Mismatch} = 0 \text{ point} \end{cases}$$

| BAT | BAT | BAT | BAT |
| :---: | :---: | :---: | :---: |
| \|\|\| | \|\|\| | \|\|\| | \|\|\| |
| BATTER | BATTER | BATTER | BATTER |
| 111000 | 000100 | 000000 | 000000 |
| ⬇ | ⬇ | ⬇ | ⬇ |
| 3 points | 1 point | 0 points | 0 points |

# Now we can score the alignments between our four words: hat, hatter, bat, batter

### Alignment 1

```
BAT
|||
HAT
011
```

**Score = 2**

### Alignment 2

```
HAT
|||
BATTER
011
```

**Score = 2**

### Alignment 3

```
BAT
|||
HATTER
011
```

**Score = 2**

### Alignment 4

```
BAT
|||
BATTER
111000
```

**Score = 3**

### Alignment 5

```
HAT
|||
HATTER
111000
```

**Score = 3**

### Alignment 6

```
BATTER
||||||
HATTER
011111
```

**Score = 5**

# Just when you thought you understood alignments, it gets a bit more complicated

Consider these three phrases ➡ {
THE CATS IN THE HAT

THE CAT IN THE HAT

THE CAT IS A HAT
}

Which two phrases are the most similar?

What algorithm did you use to figure that out?

Will our previous alignment method work?
Why or why not?

The current algorithm can't finding LCSs with additional or missing characters (insertions and deletions)

We'll use the same scoring scheme as before

➡

Match = 1 point

Mismatch = 0 points

```
THECATINTHEHAT
||||||||||||||
THECATSINTHEHAT
111111100000000
```

➡

```
THECATINTHEHAT
||||||||||||||||
THECATSINTHEHAT
00000011111111
```

⬇

⬇

6 points

8 points

How can we fix the algorithm?

# We can improve the algorithm by allowing gaps in the longest common subsequence

We'll use a modified scoring scheme ➡

$$
\begin{cases}
\text{Match} = 1 \\
\text{Mismatch} = 0 \\
\text{Gap open} = \text{-1} \\
\text{Gap extension} = 0
\end{cases}
$$

```
THECAT-INTHEHAT
||||||  ||||||||
THECATSINTHEHAT
111111-11111111
```

➡ 6 - 1 + 8 = 13

Better!

You might be asking yourselves, why do we need the gap penalty?

Which of these two alignments is better? Why?

➡

```
THECATINTHEHAT
||||||||||||||
THECATINTHEHAT

THECAT-INTHEHAT
||||||  |||||||
THECATSINTHEHAT
```

Without a gap penalty, both alignments have the same score (14)

We need gap penalties to reflect the intuition that, all things being equal, ungapped alignments are better than gapped alignments

# Now we can score the alignments between our three phrases (sequences)

## Alignment 1

```
THECAT-INTHEHAT
||||||  |||||||||
THECATSINTHEHAT
111111-11111111
```

**Score = 13**

## Alignment 2

```
THECAT---ISAHAT
||||||    ||||||
THECATSINTHEHAT
111111-00000111
```

**Score = 8**

## Alignment 3

```
THECATI--SAHAT
|||||||  |||||
THECATINTHEHAT
1111111-000111
```

**Score = 9**

# Alternative alignments for 2

## Alignment 2a

```
THECAT---ISAHAT
||||||   ||||||
THECATSINTHEHAT
111111-00000111
```

**Score = 8**

## Alignment 2b

```
THECAT-I--SAHAT
|||||| |   |||||
THECATSINTHEHAT
111111-1-000111
```

**Score = 8**

# Biological sequence alignments

<u>Biological sequence alignment algorithms</u> are ways of arranging two or more molecular sequences to identify regions of similarity between them

Types of molecules ➡️

- DNA
- RNA
- Protein

We'll focus on DNA, which consists of four nucleotides (alphabet size = 4) ➡️

- Adenine (A)
- Thymine (T)
- Guanine (G)
- Cytosine (C)

We align biological sequences in the same way as we did with English words and phrases

# Example → aligning two DNA sequences

ACTG vs. ACGTG

Match = 1, mismatch = 0, gap open = -1, gap extension = 0

```
ACTG            ACTG            A-CTG           AC-TG           ACT-G
||||            ||||            |  |||          ||  ||          ||| |
ACGTG           ACGTG           ACGTG           ACGTG           ACGTG
11000           00011           1-011           11-11           110-1
```

2 points     2 points     2 points     3 points     2 points

# Why are biological sequence alignments important?

The more similar two molecular sequences are, the more likely that the molecules are also similar in:

→

Structure

Function

Evolutionary history

```
ACCTG
| | | | |
ACGTG
      ↑
Point mutation
```

Deletion?

↓

```
AC-TG
| |   | |
ACGTG
      ↑
```

Insertion?

# We need computers and algorithms to find biological sequence alignments

Why not find all biological sequence alignments manually?

How many times can you find the <u>query sequence</u> ATCGGCCATTAC in the following <u>target sequence</u>? Is it there at all? If so, is it unique?

```
ATCACTGTAGTAGTAGCTGGAAAGAGAAATCTGTGACTCCAATTAGCCAGTTCCTGCAGACCTTGTGAGGACTAG
AGGAAGAATGCTCCTGGCTGTTTTGTACTGCCTGCTGTGGAGTTTCCAGACCTCCGCTGGCCATTTCCCTAGAGC
CTGTGTCTCCTCTAAGAACCTGATGGAGAAGGAATGCTGTCCACCGTGGAGCGGGGACAGGAGTCCCTGTGGCCA
GCTTTCAGGCAGAGGTTCCTGTCAGAATATCCTTCTGTCCAATGCACCACTTGGGCCTCAATTTCCCTTCACAGG
GGTGGATGACCGGGAGTCGTGGCCTTCCGTCTTTTATAATAGGACCTGCCAGTGCTCTGGCAACTTCATGGGATT
CAACTGTGGAAACTGCAAGTTTGGCTTTTGGGGACCAAACTGCACAGAGAGCGACTCTTGGTGAGAAGAAACAT
CTTCGATTTGAGTGCCCCAGAGAAGGACAAATTTTTTGCCTACCTCACTTTAGCAAAGCATACCATCAGCTCAGA
CTATGTCATCCCCATAGGGACCATTGGCCAAATGAAAAATGGATCAACACCCATGTTTAACGACATCAATATTTA
TGACCTCTTTGTCTGGATGCATTATTATGTGTCAATGGATGCACTGCTTGGGGGATCTGAAATCTGGAGAGACAT
TGATTTTGCCCATGAAGCACCAGCTTTTCTGCCTTGGCATAGACTCTTCTTGTTGCGGTGGGAACAAGAAATCCA
GAAGCTGACAGGAGATGAAAACTTCACTATTCCATATTGGGACTGGCGGGATGCAGAAAGTGTGACATTTGCAC
AGATGAGTACATGGG
```

# <u>B</u>asic <u>L</u>ocal <u>A</u>lignment <u>S</u>earch <u>T</u>ool
# (BLAST)

# There are two types of sequence alignment: local and global

Global alignment → find the single best alignment across the entire length of both sequences

Local alignment → find one or more highly similar local regions between both sequences

# BLAST is a very fast tool for finding local regions of similarity between biological sequences

## Basic Local Alignment Search Tool

```
ATCACTGTAGTAGTAG
CTGGAAAGAGAAATCT
GTGACTCCAATTAGCC
```

Your query sequence

Can be DNA, RNA or protein

Internet search or local search

BLAST your sequence: search target database for local alignments

Target DB

nt = non-redundant nucleotide sequence database

nr = non-redundant protein sequence database

Target databases are extremely large; millions of sequences

# There are six types of BLAST, depending on the type of query and target sequences

Nucleotide BLAST (blastn) → Search a nucleotide database using a nucleotide query

Protein BLAST (blastp) → Search a protein database using a protein query

blastx → Search a protein database using a translated nucleotide query

tblastn → Search a translated nucleotide database using a protein query

tblastx → Search translated nucleotide database using a translated nucleotide query

# Many animals use the Earth's magnetic field for orientation and navigation esp. during migration

Some examples: sea turtles, swallows, monarch butterflies and fruit flies (Drosophila melanogaster)

Cryptochome is a key protein for geomagnetic sensing; it seems to be a quantum compass

Humans produce cryptochrome in the retina, but we don't seem to have this geomagnetic perception

# Human cryptochrome, when it replaces fruit fly cryptochrome, works the same way

Researchers created cryptochrome-deficient flies, and they lost their ability to navigate

They then created transgenic flies with human cryptochrome instead of their normal version

The flies with the human cryptochrome could navigate just as well as the flies with the normal version

How similar are the protein sequences of human cryptochrome and fly cryptochrome?

# We can use BLAST to find protein sequences in humans that are similar to fly cryptochrome



**Distribution of 13 Blast Hits on the Query Sequence**

Mouse-over to show defline and scores, click to show alignments

**Color key for alignment scores**

| <40 | 40-50 | 50-80 | 80-200 | >=200 |

Query  1   100   200   300   400   500

**Descriptions**

Legend for links to other resources: **U** UniGene **E** GEO **G** Gene **S** Structure **M** Map Viewer **PubChem BioAssay**

**Sequences producing significant alignments:**

| Accession | Description | Max score | Total score | Query coverage | E value | Links |
|---|---|---|---|---|---|---|
| BAA31633.1 | KIAA0658 protein [Homo sapiens] | 375 | 375 | 94% | 3e-103 | G M |
| Q49AN0.2 | RecName: Full=Cryptochrome-2 >gb|AAH41814.1| Cryptochrome 2 (photolya | 375 | 375 | 94% | 3e-103 | G |
| BAF83949.1 | unnamed protein product [Homo sapiens] | 374 | 374 | 94% | 4e-103 | G M |
| NP_004066.1 | cryptochrome-1 [Homo sapiens] >sp|Q16526.1|CRY1_HUMAN RecName: Full | 374 | 374 | 94% | 4e-103 | U G M |
| NP_066940.2 | cryptochrome-2 isoform 1 [Homo sapiens] >dbj|BAG64048.1| unnamed prote | 374 | 374 | 94% | 4e-103 | U G M |
| NP_001120929.1 | cryptochrome-2 isoform 2 [Homo sapiens] | 332 | 332 | 85% | 2e-90 | U G M |
| BAG57993.1 | unnamed protein product [Homo sapiens] | 289 | 289 | 72% | 2e-77 | G M |
| AAH35161.1 | CRY2 protein [Homo sapiens] | 287 | 287 | 72% | 6e-77 | G M |
| EAW97796.1 | cryptochrome 1 (photolyase-like), isoform CRA_b [Homo sapiens] | 259 | 259 | 59% | 2e-68 | G |
| BAG58504.1 | unnamed protein product [Homo sapiens] | 241 | 241 | 47% | 5e-63 | G M |
| BAC05354.1 | unnamed protein product [Homo sapiens] | 65.9 | 65.9 | 14% | 3e-10 | G M |
| BAC86686.1 | unnamed protein product [Homo sapiens] | 31.2 | 31.2 | 14% | 8.6 | G M |
| NP_001091970.1 | sickle tail protein homolog isoform 2 [Homo sapiens] >emb|CAI12212.1| KIA | 31.2 | 31.2 | 9% | 9.4 | U G M |

# We get a list – and visual overview – of alignments of the query sequence to target sequences

# We can view the alignment between query and target sequences for each match

- Query coverage = 94%

- Score =  374 bits (960)

- Expect = 4e-103

- Identities = 214/521 (41%)

- Positives = 298/521 (57%)

- Gaps = 41/521 (8%)

# We can also view a detailed record for each matching target sequence



**cryptochrome-2 isoform 1 [Homo sapiens]**

NCBI Reference Sequence: NP_066940

FASTA    Graphics

Go to: ⌄

```
LOCUS       NP_066940                614 aa            linear   PRI 15-MAY-2011
DEFINITION  cryptochrome-2 isoform 1 [Homo sapiens].
ACCESSION   NP_066940
VERSION     NP_066940.2  GI:188536100
DBSOURCE    REFSEQ: accession NM_021117.3
KEYWORDS    .
SOURCE      Homo sapiens (human)
  ORGANISM  Homo sapiens
            Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;
            Mammalia; Eutheria; Euarchontoglires; Primates; Haplorrhini;
            Catarrhini; Hominidae; Homo.
REFERENCE   1  (residues 1 to 614)
  AUTHORS   Dai,H., Zhang,L., Cao,M., Song,F., Zheng,H., Zhu,X., Wei,Q.,
            Zhang,W. and Chen,K.
  TITLE     The role of polymorphisms in circadian pathway genes in breast
            tumorigenesis
  JOURNAL   Breast Cancer Res. Treat. 127 (2), 531-540 (2011)
   PUBMED   20978934
  REMARK    GeneRIF: Observational study of gene-disease association, gene-gene
            interaction, and gene-environment interaction. (HuGE Navigator)
REFERENCE   2  (residues 1 to 614)
  AUTHORS   Bailey,S.D., Xie,C., Do,R., Montpetit,A., Diaz,R., Mohan,V.,
            Keavney,B., Yusuf,S., Gerstein,H.C., Engert,J.C. and Anand,S.
  CONSRTM   DREAM investigators
  TITLE     Variation at the NFATC2 locus increases the risk of
            thiazolidinedione-induced edema in the Diabetes REduction
            Assessment with ramipril and rosiglitazone Medication (DREAM) study
  JOURNAL   Diabetes Care 33 (10), 2250-2253 (2010)
   PUBMED   20628086
  REMARK    GeneRIF: Observational study of gene-disease association,
            gene-environment interaction, and pharmacogenomic / toxicogenomic.
            (HuGE Navigator)
REFERENCE   3  (residues 1 to 614)
  AUTHORS   Fontaine-Bisson,B., Renstrom,F., Rolandsson,O., Payne,F.,
            Hallmans,G., Barroso,I. and Franks,P.W.
  CONSRTM   MAGIC
  TITLE     Evaluating the discriminative power of multi-trait genetic risk
            scores for type 2 diabetes in a northern Swedish population
```

# BLAST Tutorial

http://www.digitalworldbiology.com/BLAST/slide1.html

# BLAST Tutorial → Slide 1

# BLAST Tutorial → Slide 2

# BLAST Tutorial → Slide 3

# BLAST Tutorial → Slide 4

# BLAST Tutorial → Slide 5

# BLAST Tutorial → Slide 6

# BLAST Tutorial → Slide 7



The total score includes scores from non-contiguous portions of the subject sequence that match the query.

The Max % ident corresponds to the match to a subject sequence with the highest percentage of identical bases.

A description of the sequence

Legend for links to other resources: **U** UniGene  **E** GEO  **G** Gene  **S** Structure  **M** Map Viewer

Sequences producing significant alignments:
(Click headers to sort columns)

| Accession | Description | Max score | Total score | Query coverage | E value | Max ident | Links |
|-----------|-------------|-----------|-------------|----------------|---------|-----------|-------|
| X16893.1 | Tarantula mRNA for hemocyanin subunit a | 4057 | 4057 | 100% | 0.0 | 100% | |
| AJ547807.1 | Nephila inaurata madagascariensis mRNA for hemocyani | 662 | 662 | 79% | 0.0 | 73% | |
| AJ547811.1 | Nephila inaurata madagascariensis mRNA for hemocyani | 202 | 319 | 43% | 5e-48 | 88% | |
| AJ547809.1 | Nephila inaurata madagascariensis mRNA for hemocyani | 185 | 241 | 14% | 7e-43 | 83% | |
| AJ307908.1 | Cupiennius salei mRNA for hemocyanin subunit 5' (hc-5' ; | 175 | 298 | 19% | 6e-40 | 84% | |
| AJ277492.1 | Eurypelma californicum mRNA for hemocyanin subunit g | 171 | 278 | 22% | 8e-39 | 79% | |

The Accession number is linked to the GenBank record.

A score that indicates how well the sequences match. For nucleotide sequences, this is approximately equal to twice the length of the matching region.

The Max score is linked to data that show where the sequences match.

The query coverage corresponds to the fraction of the query sequence that matches a subject sequence.

See the next page to learn more about the E value.

The E value is equal to the number of sequences that you would expect to find in a database composed entirely of random sequences.

Two important parameters that influence the E value are:

- The number of sequences in the database (database size).
- The length of the query sequence.

There is a greater chance of finding a match in a larger database. And the chance of finding a match for a short sequence is greater than the chance of finding a match to a longer sequence.

E Value

0.0
1e-21
7e-20
4e-18
2e-14
1e-09
4e-09
1e-08
1e-08
1e-08
0.014
0.014
0.014
0.014
0.014
0.014
0.22
0.22
0.86
0.86
0.86
0.86
3.4
3.4

In this example, the E value equals

$$1 \times 10^{-21}$$

The letter "e" is used to show that -21 is the exponent. You would "expect" to find very few random sequences in this database that match the query sequence this well.

This sequence has an E value of 3.4. A database of random sequences would be likely to contain 3.4 sequences that matched the query sequence equally well.

# BLAST Tutorial → Slide 9



Legend for links to other resources: U UniGene  E GEO  G Gene  S Structure  M Map Viewer

Sequences producing significant alignments:
(Click headers to sort columns)

| Accession | Description | Max score | Total score | Query coverage | E value | Max ident | Links |
|---|---|---|---|---|---|---|---|
| X16893.1 | Tarantula mRNA for hemocyanin subunit a | 4057 | 4057 | 100% | 0.0 | 100% | |
| AJ547807.1 | Nephila inaurata madagascariensis mRNA for hemocyanin | 662 | 662 | 79% | 0.0 | 73% | |
| AJ547811.1 | Nephila inaurata madagascariensis mRNA for hemocyanin | 202 | 319 | 43% | 5e-48 | 88% | |
| AJ547809.1 | Nephila inaurata madagascariensis mRNA for hemocyanin | 185 | 241 | 14% | 7e-43 | 83% | |
| AJ307908.1 | Cupiennius salei mRNA for hemocyanin subunit 5' (hc-5' ; | 175 | 298 | 19% | 6e-40 | 84% | |
| AJ277492.1 | Eurypelma californicum mRNA for hemocyanin subunit g | 171 | 278 | 22% | 8e-39 | 79% | |

The Max score is linked to data that show where the sequences match.

See where the sequences align.

# BLAST Tutorial → Slide 10

# BLAST Tutorial → Slide 10a

# BLAST Tutorial → Slide 10b



Legend for links to other resources: [U] UniGene  [E] GEO  [G] Gene  [S] Structure  [M] Map Viewer

Sequences producing significant alignments:
(Click headers to sort columns)

| Accession | Description | Max score | Total score | Query coverage | E value | Max ident | Links |
|---|---|---|---|---|---|---|---|
| X16893.1 | Tarantula mRNA for hemocyanin subunit a | 4057 | 4057 | 100% | 0.0 | 100% | |
| AJ__807.1 | Nephila inaurata madagascariensis mRNA for hemocyani | 662 | 662 | 79% | 0.0 | 73% | |
| A____1.1 | Nephila inaurata madagascariensis mRNA for hemocyani | 202 | 319 | 43% | 5e-48 | 88% | |
| AJ__809.1 | Nephila inaurata madagascariensis mRNA for hemocyani | 185 | 241 | 14% | 7e-43 | 83% | |
| AJ__808.1 | Cupiennius salei mRNA for hemocyanin subunit 5' (hc-5' | 175 | 298 | 19% | 6e-40 | 84% | |
| AJ__892.1 | Eurypelma californicum mRNA for hemocyanin subunit g | 171 | 278 | 22% | 8e-39 | 79% | |

Look at the GenBank record.

Copyright © Digital World Biology All rights reserved.

# BLAST Tutorial → Slide 11

# BLAST Tutorial → Slide 12

# Hand-on BLAST Tutorial

(1) Open the BLAST web application:

http://blast.ncbi.nlm.nih.gov/Blast.cgi

(2) In another tab, open this web page:

http://www.digitalworldbiology.com/BLAST/62000sequences.html

(3) Copy and paste the example sequence into the text box on the BLAST page

# Questions?

# Questions?