# BIMM 143

## Introduction to Bioinformatics
### Lecture 2

**Barry Grant**

UC San Diego

http://thegrantlab.org/bimm143

---

## Recap From Last Time:

- Bioinformatics is computer aided biology.
  - Deals with the collection, archiving, organization, and interpretation of a wide range of biological data.

- The NCBI and EBI are major online bioinformatics service providers.

- Introduced via **hands-on session** the BLAST, Entrez, GENE, OMIM, UniProt, Muscle and PDB bioinformatics tools and databases.

    - Muddy point assessment (see results)

- There are a large number of bioinformatics databases (see handout!).

- Also covered: Course structure; Supporting course website, Ethics code, and Introductions…

---

# Today's Menu

| | |
|---|---|
| **Classifying Databases** | Primary, secondary and composite Bioinformatics databases |
| **Using Databases** | **Vignette** demonstrating how major Bioinformatics databases intersect |
| **Major Biomolecular Formats** | How nucleotide and protein sequence and structure data are represented |
| **Alignment Foundations** | Introducing the *why* and *how* of comparing sequences |
| **Alignment Algorithms** | **Hands-on** exploration of alignment algorithms and applications |

---

## Primary, secondary & composite databases

Bioinformatics databases can be usefully classified into *primary*, *secondary* and *composite* according to their data source.

- **Primary databases** (or *archival databases*) consist of data derived experimentally.
  - **GenBank**: NCBI's primary nucleotide sequence database.
  - **PDB:** Protein X-ray crystal and NMR structures.

- **Secondary databases** (or *derived databases*) contain information derived from a primary database.
  - **RefSeq**: non redundant set of curated reference sequences primarily from GenBank
  - **PFAM**: protein sequence families primarily from UniProt and PDB

- **Composite databases** (or *metadatabases*) join a variety of different primary and secondary database sources.
  - **OMIM**: catalog of human genes, genetic disorders and related literature
  - **GENE**: molecular data and literature related to genes with extensive links to other databases.

## Slide 1

# DATABASE VIGNETTE

You have just come out a seminar about gastric cancer and one of your co-workers asks:

> "*What do you know about that 'Kras' gene the speaker kept taking about?*"

You have some recollection about hearing of 'Ras' before. How would you find out more?

- Google?
- Library?
- **Bioinformatics databases at NCBI and EBI!**

http://www.ncbi.nlm.nih.gov/

## Slide 2

http://www.ncbi.nlm.nih.gov/



Hands on demo (or see following slides)

## Slide 3

# Example Vignette Questions:

- What chromosome location and what genes are in the vicinity of a given query gene? NCBI **GENE**

- What can you find out about molecular functions, biological processes, and prominent cellular locations? EBI **GO**

- What amino acid positions in the protein are responsible for ligand binding? EBI **UniProt**

- What variants of this gene are associated with gastric cancer and other human diseases? NCBI **OMIN**

- What is known about the protein family, its species distribution, number in humans and residue-wise conservation? EBI **PFAM**

- Are high resolution protein structures available to examine the details of these mutations? How might we explain their potential molecular effects? RCSB **PDB**

## Slide 4

1 AND 2     ras AND disease (1185 results)

1 OR 2     ras OR disease (134,872 results)

1 NOT 2     ras NOT disease (84,448 results)

## GO: Gene Ontology

GO provides a controlled vocabulary of terms for describing gene product characteristics and gene product annotation data

# Why do we need Ontologies?

- Annotation is essential for capturing the understanding and knowledge associated with a sequence or other molecular entity

- Annotation is traditionally recorded as "free text", which is easy to read by humans, but has a number of disadvantages, including:
  - ‣ Difficult for computers to parse
  - ‣ Quality varies from database to database
  - ‣ Terminology used varies from annotator to annotator

- Ontologies are annotations using standard vocabularies that try to address these issues

- GO is integrated with UniProt and many other databases including a number at NCBI

# GO Ontologies

- There are three ontologies in GO:
  - ‣ **Biological Process**
    A commonly recognized series of events
    e.g. cell division, mitosis,

  - ‣ **Molecular Function**
    An elemental activity, task or job
    e.g. kinase activity, insulin binding

  - ‣ **Cellular Component**
    Where a gene product is located
    e.g. mitochondrion, mitochondrial membrane

## Slide 1

**Recap:** Major NCBI and EBI databases

- What chromosome location and what genes are in the vicinity of a given query gene?  NCBI **GENE**

- What can you find out about molecular functions, biological processes, and prominent cellular locations?  EBI **GO**

- What amino acid positions in the protein are responsible for ligand binding?  EBI **UniProt**

- What variants of this gene are associated with gastric cancer and other human diseases?  NCBI **OMIN**

- What is known about the protein family, its species distribution, number in humans and residue-wise conservation?  EBI **PFAM**

- Are high resolution protein structures available to examine the details of these mutations? How might we explain their potential molecular effects?  RCSB **PDB**

## Slide 2

# Today's Menu

| Classifying Databases | Primary, secondary and composite Bioinformatics databases |
|---|---|
| Using Databases | **Vignette** demonstrating how major Bioinformatics databases intersect |
| Major Biomolecular Formats | How nucleotide and protein sequence and structure data are represented |
| Alignment Foundations | **Introducing the *why* and *how* of comparing sequences** |
| Alignment Algorithms | **Hands-on** exploration of alignment algorithms and applications |

## Slide 3

ALIGNMENT FOUNDATIONS

- **Why…**
  - ‣ Why compare biological sequences?

- **What…**
  - ‣ Alignment view of sequence changes during evolution (matches, mismatches and gaps)

- **How…**
  - ‣ Dot matrices
  - ‣ Dynamic programing
    - Global alignment
    - Local alignment
  - ‣ BLAST heuristic approach

## Slide 4

ALIGNMENT FOUNDATIONS

- **Why…**
  - ‣ Why compare biological sequences?

- **What…**
  - ‣ Alignment view of sequence changes during evolution (matches, mismatches and gaps)

- **How…**
  - ‣ Dot matrices
  - ‣ Dynamic programing
    - Global alignment
    - Local alignment
  - ‣ BLAST heuristic approach

**Basic Idea:** Display one sequence above another with spaces (termed gaps) inserted in both to reveal similarity of nucleotides or amino acids.

Seq1: C A T T C A C

Seq2: C T C G C A G C

[Screencast Material]

---

**Basic Idea:** Display one sequence above another with spaces (termed gaps) inserted in both to reveal similarity of nucleotides or amino acids.

Seq1: C A T T C A C

Seq2: C T C G C A G C

mismatch
match

Two types of character correspondence

---

**Basic Idea:** Display one sequence above another with spaces (termed gaps) inserted in both to reveal similarity of nucleotides or amino acids.

Seq1: C A T – T C A – C

Seq2: C – T C G C A G C

match
mismatch

Add gaps to increase number of matches

gaps

---

**Basic Idea:** Display one sequence above another with spaces (termed gaps) inserted in both to reveal similarity of nucleotides or amino acids.

Seq1: C A T – T C A – C

Seq2: C – T C G C A G C

match
mismatch } mutation
insertion
deletion } indels

Gaps represent 'indels'
mismatch represent mutations

## Why compare biological sequences?

- To obtain functional or mechanistic insight about a sequence by inference from another potentially better characterized sequence
- To find whether two (or more) genes or proteins are evolutionarily related
- To find structurally or functionally similar regions within sequences (e.g. catalytic sites, binding sites for other molecules, etc.)
- Many practical bioinformatics applications…

## Practical applications include...

- Similarity searching of databases
  - Protein structure prediction, annotation, etc...
- Assembly of sequence reads into a longer construct such as a genomic sequence
- Mapping sequencing reads to a known genome
  - "Resequencing", looking for differences from reference genome - SNPs, indels (insertions or deletions)
  - Mapping transcription factor binding sites via ChIP-Seq (chromatin immuno-precipitation sequencing)
  - Pretty much all next-gen sequencing data analysis

## Practical applications include...

- Similarity searching of databases
  - Protein structure prediction
- Assembly of sequences construct such as a bact
- Mappin to a known genome
  - king for differences from reference
  - s, indels (insertions or deletions)
  - ing transcription factor binding sites via ChIP-Seq (chromatin immuno-precipitation sequencing)
  - Pretty much all next-gen sequencing data analysis

N.B. Pairwise sequence alignment is arguably the most fundamental operation of bioinformatics!

## ALIGNMENT FOUNDATIONS

- **Why…**
  - Why compare biological sequences?
- **What…**
  - Alignment view of sequence changes during evolution (matches, mismatches and gaps)
- **How…**
  - Dot matrices
  - Dynamic programing
    - Global alignment
    - Local alignment
  - BLAST heuristic approach

## Sequence changes during evolution

There are three major types of sequence change that can occur during evolution.
- Mutations/Substitutions
- Deletions
- Insertions

Common Ancestor **(C)** CTCGTTA

Time

Recent Species **(B)** CATGTTA **(A)** CACTGTA

## Mutations, deletions and insertions

There are three major types of sequence change that can occur during evolution.
- Mutations/Substitutions       CTCGTTA → C**A**CGTTA
- Deletions
- Insertions

CTCGTTA

C**A**CGTTA ← **Mutation**

CATGTTA          CACTGTA

Likely occurred prior to speciation

## Mutations, deletions and insertions

There are three major types of sequence change that can occur during evolution.
- Mutations/Substitutions       CTCGTTA → C**A**CGTTA
- Deletions
- Insertions

CTCGTTA

C**A**CGTTA ← **Mutation**

CACGTTA     CACGTTA       (**speciation**)

CATGTTA       CACTGTA

## Mutations, deletions and insertions

There are three major types of sequence change that can occur during evolution.
- Mutations/Substitutions       CTCGTTA → C**A**CGTTA
- Deletions                     CAC**G**TTA → CACTTA
- Insertions

CTCGTTA

C**A**CGTTA ← **Mutation**

CACGTTA     CAC✗TTA
                      ←**Deletion**
              CACTTA

CATGTTA       CACTGTA

## Mutations, deletions and insertions

There are three major types of sequence change that can occur during evolution.
- Mutations/Substitutions
- Deletions
- Insertions

CTCGTTA ⟶ CACGTTA
CACGTTA ⟶ CACTTA
CACTTA ⟶ CACTGTA

CTCGTTA
CACGTTA — **Mutation**
CACGTTA   CACGTTA
CACTTA
**Deletion**
CATGTTA   CACTGTA **Insertion**

## Mutations, deletions and insertions

There are three major types of sequence change that can occur during evolution.
- Mutations/Substitutions
- Deletions
- Insertions

CTCGTTA ⟶ CACGTTA
CACGTTA ⟶ CATGTTA

CTCGTTA
CACGTTA — **Mutation**
CACGTTA
**Mutation** CATGTTA   CACTGTA

## Alignment view

Alignments are great tools to visualize sequence similarity and evolutionary changes in homologous sequences.
- Mismatches represent mutations/substitutions
- Gaps represent insertions and deletions (indels)

CTCGTTA
(C)

(B)   (A)
CATGTTA   CACTGTA

Substitution        Indels

(A) CAC-TGTA
(B) CATGT-TA

| Match 5 | Mismatch 1 | Gap 2 |
|---|---|---|

## Alternative alignments

- Unfortunately, finding the correct alignment is difficult if we do not know the evolutionary history of the two sequences

Q. Which of these 3 possible alignments is best?

1.
CACTGTA
CATGTTA

2.
CACTGT-A
CA-TGTTA

3.
CAC-TGTA
CATGT-TA

## Alternative alignments

- One way to judge alignments is to compare their number of matches, insertions, deletions and mutations

🟢 4 matches    🟢 6 matches    🟢 5 matches
🟠 3 mismatches    🟠 0 mismatches    🟠 1 mismatches
⚪ 0 gaps    ⚪ 2 gaps    ⚪ 2 gaps

```
CACTGTA    CACTGT-A    CAC-TGTA
|| :::||    || |||  |    || :  |  ||
CATGTTA    CA-TGTTA    CATGT-TA
```

## Scoring alignments

- We can assign a score for each match (+3), mismatch (+1) and indel (-1) to identify the optimal alignment for this scoring scheme

🟢 4 (+3)    🟢 6 (+3)    🟢 5 (+3)
🟠 3 (+1)    🟠 0 (+1)    🟠 1 (+1)
⚪ 0 (-1) = 15    ⚪ 2 (-1) = 16    ⚪ 2 (-1) = 14

```
CACTGTA    CACTGT-A    CAC-TGTA
|| :::||    || |||  |    || :  |  ||
CATGTTA    CA-TGTTA    CATGT-TA
```

## Optimal alignments

- Biologists often prefer parsimonious alignments, where the number of postulated sequence changes is minimized.

🟢 4 matches    🟢 6 matches    🟢 5 matches
🟠 3 mismatches    🟠 0 mismatches    🟠 1 mismatches
⚪ 0 gaps    ⚪ 2 gaps    ⚪ 2 gaps

```
CACTGTA    CACTGT-A    CAC-TGTA
|| :::||    || |||  |    || :  |  ||
CATGTTA    CA-TGTTA    CATGT-TA
```

## Optimal alignments

- Biologists often prefer parsimonious alignments, where the number of postulated sequence changes is minimized.

🟢 4 matches    🟢 6 matches    🟢 5 matches
🟠 3 mismatches    🟠 0 mismatches    🟠 1 mismatches
⚪ 0 gaps    ⚪ 2 gaps    ⚪ 2 gaps

```
CACTGTA    CACTGT-A    CAC-TGTA
|| :::||    || |||  |    || :  |  ||
CATGTTA    CA-TGTTA    CATGT-TA
```

# Optimal alignments

- Biologists often prefer parsimonious alignments, where the number of postulated sequence changes is minimized.

🟢 4 matches
🟠 3 mismatches
⚪ 0 gaps

🟢 6 matches
🟠 0 mismatches
⚪ 2 gaps

🟢 5 matches
🟠 1 mismatches
⚪ 2 gaps

```
CACTGTA    CACTG-TA    CAC-TGTA
| |:::||    || || ||    || : | ||
CATGTTA    CA-TGTTA    CATGT-TA
```

---

# Optimal alignments

- Biologists often prefer parsimonious alignments, where the number sequence changes is

Warning: There may be more than one optimal alignment and these may not reflect the true evolutionary history of our sequences!

🟢 4 matches
🟠 3 mismatches

🟠 1 mismatches
⚪ 2 gaps

```
CACTGTA    CACTGT-A    CAC-TGTA
| |:::||    || ||| |    || : | ||
CATGTTA    CA-TGTTA    CATGT-TA
```

---

# ALIGNMENT FOUNDATIONS

- **Why…**
  - Why compare biological sequences?
- **What…**
  - Alignment view of sequence changes during evolution (matches, mismatches and gaps)
- **How…**
  - ‣ Dot matrices
  - ‣ Dynamic programing
    - Global alignment
    - Local alignment
  - ‣ BLAST heuristic approach

---

# ALIGNMENT FOUNDATIONS

- **Why…**
  - Why compare biological sequences?
- **What…**
  - Alignment view of sequence changes during evolution (matches, mismatches and gaps)
- **How…**
  - ‣ Dot matrices
  - ‣ D

    How do we compute the optimal alignment between two sequences?

  - ‣ BLAST heuristic approach

## Dot plots: simple graphical approach

- Place one sequence on the vertical axis of a 2D grid (or matrix) and the other on the horizontal



## Dot plots: simple graphical approach

- Now simply put dots where the horizontal and vertical sequence values match



## Dot plots: simple graphical approach

- Diagonal runs of dots indicate matched segments of sequence



## Dot plots: simple graphical approach

Q. What would the dot matrix of a two identical sequences look like?

## Dot plots: simple graphical approach

• Dot matrices for long sequences can be noisy



## Dot plots: window size and match stringency

Solution: use a window and a threshold
– compare character by character within a window
– require certain fraction of matches within window in order to display it with a dot.
  • You have to choose window size and stringency



## Dot plots: window size and match stringency

Solution: use a window and a threshold
– compare character by character within a window
– require certain fraction of matches within window in order to display it with a dot.
  • You have to choose window size and stringency



## Dot plots: window size and match stringency

Solution: use a window and a threshold
– compare character by character within a window
– require certain fraction of matches within window in order to display it with a dot.
  • You have to choose window size and stringency

## Dot plots: window size and match stringency

Solution: use a window and a threshold
– compare character by character within a window
– require certain fraction of matches within window in order to display it with a dot.
• You have to choose window size and stringency



Filter
Window = 3
Stringency = 3

## Dot plots: window size and match stringency

Solution: use a window and a threshold
– compare character by character within a window
– require certain fraction of matches within window in order to display it with a dot.
• You have to choose window size and stringency



Filter
Window = 3
Stringency = 3

## Dot plots: window size and match stringency

Solution: use a window and a threshold
– compare character by character within a window
– require certain fraction of matches within window in order to display it with a dot.
• You have to choose window size and stringency



Filter
Window = 3
Stringency = 3

## Dot plots: window size and match stringency

Solution: use a window and a threshold
– compare character by character within a window
– require certain fraction of matches within window in order to display it with a dot.
• You have to choose window size and stringency



Filter
Window = 3
Stringency = 3

# Dot plots: window size and match stringency

Solution: use a window and a threshold
– compare character by character within a window
– require certain fraction of matches within window
  in order to display it with a dot.
  • You have to choose window size and stringency



Filter
Window = 3
Stringency = 3

# Dot plots: window size and match stringency

Solution: use a window and a threshold
– compare character by character within a window
– require certain fraction of matches within window
  in order to display it with a dot.
  • You have to choose window size and stringency



Filter
Window = 3
Stringency = 3

# Dot plots: window size and match stringency

Solution: use a window and a threshold
– compare character by character within a window
– require certain fraction of matches within window
  in order to display it with a dot.
  • You have to choose window size and stringency



Filter
Window = 3
Stringency = 3

# Dot plots: window size and match stringency

Solution: use a window and a threshold
– compare character by character within a window
– require certain fraction of matches within window
  in order to display it with a dot.
  • You have to choose window size and stringency



Filter
Window = 3
Stringency = 2

# Dot plots: window size and match stringency

Solution: use a <u>window</u> and a <u>threshold</u>
– compare character by character within a window
– require certain fraction of matches within window in order to display it with a dot.
 • You have to choose window size and stringency



Filter

Window = 3
Stringency = 2

## Window size = 5 bases



A dot plot simply puts a dot where two sequences match. In this example, dots are placed in the plot if 5 bases in a row match perfectly. Requiring a 5 base perfect match is a **<u>heuristic</u>** – only look at regions that have a certain degree of identity.

Do you expect evolutionarily related sequences to have more word matches (matches in a row over a certain length) than random or unrelated sequences?

## Window size = 7 bases



This is a dot plot of the same sequence pair. Now 7 bases in a row must match for a dot to be place. Noise is reduced.

Using windows of a certain length is very similar to using words (kmers) of N characters in the heuristic alignment search tools

Bigger window (kmer) fewer matches to consider

Web site used: http://www.vivo.colostate.edu/molkit/dnadot/

## Ungapped alignments



indels

Only **diagonals** can be followed.

Downward or rightward paths represent **insertion** or **deletions** (gaps in one sequence or the other).

Web site used: http://www.vivo.colostate.edu/molkit/dnadot/

## Uses for dot matrices

- Visually assessing the similarity of two protein or two nucleic acid sequences
- Finding local repeat sequences within a larger sequence by comparing a sequence to itself
  - Repeats appear as a set of diagonal runs stacked vertically and/or horizontally

## Repeats



Human LDL receptor protein sequence (Genbank P01130)

W = 1
S = 1

(Figure from Mount, "Bioinformatics sequence and genome analysis")

## Repeats



Human LDL receptor protein sequence (Genbank P01130)

W = 23
S = 7

(Figure from Mount, "Bioinformatics sequence and genome analysis")

## Your Turn!

Exploration of dot plot parameters (hands-on worksheet **Section 1**)

http://bio3d.ucsd.edu/dotplot/        https://bioboot.shinyapps.io/dotplot/



https://bioboot.shinyapps.io/dotplot2/

## ALIGNMENT FOUNDATIONS

- **Why…**
  - Why compare biological sequences?
- **What…**
  - Alignment view of sequence changes during evolution (matches, mismatches and gaps)
- **How…**
  - ‣ Dot matrices
  - ‣ Dynamic programing
    - Global alignment
    - Local alignment
  - ‣ BLAST heuristic approach

---

## The Dynamic Programming Algorithm

- The dynamic programming algorithm can be thought of an extension to the dot plot approach
  - One sequence is placed down the side of a grid and another across the top
  - Instead of placing a dot in the grid, we compute a score for each position
  - Finding the optimal alignment corresponds to finding the path through the grid with the best possible score



**Needleman, S.B. & Wunsch, C.D.** (1970) "A general method applicable to the search for similarities in the amino acid sequences of two proteins." J. Mol. Biol. 48:443-453.

---

## Different paths represent different alignments



```
Seq1: D P L E      Seq1: D P M E      Seq1: D P - E
      | | : |            | | |              | | |
Seq2: D P M E      Seq2: D P - E      Seq2: D P L E
```

Matches are represented by <u>diagonal paths</u> & indels with <u>horizontal or vertical path</u> segments

---

## Algorithm of Needleman and Wunsch

- The Needleman–Wunsch approach to global sequence alignment has three basic steps:
  - (1) setting up a 2D-grid (or alignment matrix),
  - (2) scoring the matrix, and
  - (3) identifying the optimal path through the matrix



**Needleman, S.B. & Wunsch, C.D.** (1970) "A general method applicable to the search for similarities in the amino acid sequences of two proteins." J. Mol. Biol. 48:443-453.

## Scoring the alignment matrix

- Start by filling in the first row and column – these are all indels (gaps).
  - Each step you take you will add the gap penalty to the score $(S_{i,j})$ accumulated in the previous cell

| | j | Sequence 2 | | | |
|---|---|---|---|---|---|
| | - | D | P | L | E |
| - | 0 | -2 | -4 | -6 | -8 |
| D | -2 | | | | |
| P | -4 | | | | |
| M | -6 | | | | |
| E | -8 | | | | |

Sequence 1 (label on left axis)

Scores: match = +1, mismatch = -1, gap = -2

---

## Scoring the alignment matrix

- Start by filling in the first row and column – these are all indels (gaps).
  - Each step you take you will add the gap penalty to the score $(S_{i,j})$ accumulated in the previous cell

| | j | Sequence 2 | | | |
|---|---|---|---|---|---|
| | - | D | P | L | E |
| - | 0 | -2 | -4 | -6 | -8 |
| D | -2 | | | | |
| P | -4 | | | | |
| M | -6 | | | | |
| E | -8 | | | | |

Sequence 1 (label on left axis)

Scores: match = +1, mismatch = -1, gap = -2

$$S_{i+4} = (-2) + (-2) + (-2) + (-2)$$

```
Seq1: DPME
Seq2: ----
```

---

## Scoring the alignment matrix

- Then go to the empty corner cell (upper left). It has filled in values in up, left and diagonal directions
  - Now can ask which of the three directions gives the highest score?
  - keep track of this score and direction

| | j | | | | |
|---|---|---|---|---|---|
| | - | D | P | L | E |
| - | 0 | -2 | -4 | -6 | -8 |
| D | -2 | ? | | | |
| P | -4 | | | | |
| M | -6 | | | | |
| E | -8 | | | | |

Scores: match = +1, mismatch = -1, gap = -2

| | j-1 | j |
|---|---|---|
| i-1 | S(i-1, j-1) ①→ | S(i-1, j) ②↓ |
| i | S(i, j-1) ③→ | S(i, j) |

---

## Scoring the alignment matrix

- Then go to the empty corner cell (upper left). It has filled in values in up, left and diagonal directions
  - Now can ask which of the three directions gives the highest score?
  - keep track of this score and direction

| | j | | | | |
|---|---|---|---|---|---|
| | - | D | P | L | E |
| - | 0 | -2 | -4 | -6 | -8 |
| D | -2 | ? | | | |
| P | -4 | | | | |
| M | -6 | | | | |
| E | -8 | | | | |

Scores: match = +1, mismatch = -1, gap = -2

$$S(i, j) = Max \begin{cases} S(i-1, j-1) + (mis)match & ① \\ S(i-1, j) + gap\ penalty & ② \\ S(i, j-1) + gap\ penalty & ③ \end{cases}$$

## Scoring the alignment matrix

- Then go to the empty corner cell (upper left). It has filled in values in up, left and diagonal directions
  - Now can ask which direction gives the highest score
  - keep track of direction and score



Scores: match = +1, mismatch = -1, gap = -2

① $(0)+(+1) = +1$  <= (D-D) match!

② $(-2)+(-2) = -4$

③ $(-2)+(-2) = -4$

Alignment
```
D
D
```

## Scoring the alignment matrix

- At each step, the score in the current cell is determine by the scores in the neighboring cells
  - The maximal score and the direction that gave that score is stored (we will use these later to determine the optimal alignment)



Scores: match = +1, mismatch = -1, gap = -2

① $(-2)+(-1) = -3$  <= (D-P) mismatch!

② $(-4)+(-2) = -6$

③ $(1)+(-2) = -1$

Alignment
```
D-
DP
```

## Scoring the alignment matrix

- We will continue to store the alignment score ($S_{i,j}$) for all possible alignments in the alignment matrix.



Scores: match = +1, mismatch = -1, gap = -2

① $(-4)+(-1) = -5$  <= (D-L) mismatch

② $(-6)+(-2) = -8$

③ $(-1)+(-2) = -3$

Alignment
```
D--
DPL
```

## Scoring the alignment matrix

- For the highlighted cell, the corresponding score ($S_{i,j}$) refers to the score of the optimal alignment of the first i characters from sequence1, and the first j characters from sequence2.



Scores: match = +1, mismatch = -1, indel = -2

① $(-1)+(-1) = -2$

② $(-3)+(-2) = -5$

③ $(2)+(-2) = 0$

Alignment
```
DP-
DPL
```

## Scoring the alignment matrix

- At each step, the score in the current cell is determine by the scores in the neighboring cells
  - The maximal score and the direction that gave that score is stored

Scores: match = +1, mismatch = -1, indel = -2

|   | - | D | P | L | E |
|---|---|---|---|---|---|
| - | 0 | -2 | -4 | -6 | -8 |
| D | -2 | 1 | -1 | -3 | -5 |
| P | -4 | -1 | 2 | 0 | -2 |
| M | -6 | -3 | 0 | 1 | |
| E | -8 | | | | |

① (2)+(-1) = 0 <= mismatch

② (0)+(-2) = -2

③ (0)+(-2) = -2

Alignment

DPM
DPL

---

## Scoring the alignment matrix

- The score of the best alignment of the entire sequences corresponds to $S_{n,m}$
  - (where n and m are the length of the sequences)

Scores: match = +1, mismatch = -1, indel = -2

|   | - | D | P | L | E |
|---|---|---|---|---|---|
| - | 0 | -2 | -4 | -6 | -8 |
| D | -2 | 1 | -1 | -3 | -5 |
| P | -4 | -1 | 2 | 0 | -2 |
| M | -6 | -3 | 0 | 1 | -1 |
| E | -8 | -5 | -2 | -1 | 2 |

① (+1)+(+1) = +2

② (-1)+(-2) = -3

③ (-1)+(-2) = -3

Alignment

DPME
DPLE

---

## Scoring the alignment matrix

- To find the best alignment, we retrace the arrows starting from the bottom right cell
  - N.B. The optimal alignment score and alignment are dependent on the chosen scoring system

Scores: match = +1, mismatch = -1, indel = -2

|   | - | D | P | L | E |
|---|---|---|---|---|---|
| - | 0 | -2 | -4 | -6 | -8 |
| D | -2 | 1 | -1 | -3 | -5 |
| P | -4 | -1 | 2 | 0 | -2 |
| M | -6 | -3 | 0 | 1 | -1 |
| E | -8 | -5 | -2 | -1 | 2 |

Alignment

DPME
DPLE

---

## Questions:

- What is the optimal score for the alignment of these sequences and how do we find the optimal alignment?

|   | - | C | A | T | G | T | T | A |
|---|---|---|---|---|---|---|---|---|
| - | 0 | -2 | -4 | -6 | -8 | -10 | -12 | -14 |
| C | -2 | 1 | -1 | -3 | -5 | -7 | -9 | -11 |
| A | -4 | -1 | 2 | 0 | -2 | -4 | -6 | -8 |
| C | -6 | -3 | 0 | 1 | -1 | -3 | -5 | -7 |
| T | -8 | -5 | -2 | 1 | 0 | 0 | -2 | -4 |
| G | -10 | -7 | -4 | -1 | 2 | 0 | -1 | -3 |
| T | -12 | -9 | -6 | -3 | 0 | 3 | 1 | -1 |
| A | -14 | -11 | -8 | -5 | -2 | 1 | 2 | 2 |

# Questions:

- What is the optimal score for the alignment of these sequences and how do we find the optimal alignment?



# Questions:

- To find the best alignment we retrace the arrows starting from the bottom right cell



# More than one alignment possible

- Sometimes more than one alignment can result in the same optimal score



Alignment

CACTGT-A
CA-TGTTA

CACTG-TA
CA-TGTTA

# The alignment and score are dependent on the scoring system

- Here we increase the gap penalty from -2 to -3



Alignment

CACTGT-A
CA-TGTTA

CACTG-TA
CA-TGTTA

CACTGTA
CATGTTA

# The alignment and score are dependent on the scoring system

- Here we increase the gap penalty from -2 to -3



|   | - | C | A | T | G | T | T |
|---|---|---|---|---|---|---|---|
| - | 0 | -3 | -6 | -9 | | | |
| C | -3 | 1 | | | | | |
| A | | | | | | | |
| | | | | -8 | -11 | | |
| | | | 0 | -1 | -4 | -7 | |
| | -11 | -7 | -3 | 1 | -1 | -2 | -5 |
| T | -18 | -14 | -10 | -6 | -2 | 2 | 0 | -3 |
| A | -21 | -17 | -13 | -9 | -5 | -1 | 1 | 1 |

alignment

CACTGT-A
CA-TGTTA

CACTG-TA
CA-TGTTA

- - - - - - - - - - - - - -

CACTGTA
CATGTTA

**Key point: Optimal alignment solutions and their scores are not necessarily unique and depend on the scoring system!**

---

# NW DYNAMIC PROGRAMMING

Match: +2
Mismatch: -1
Gap: -2

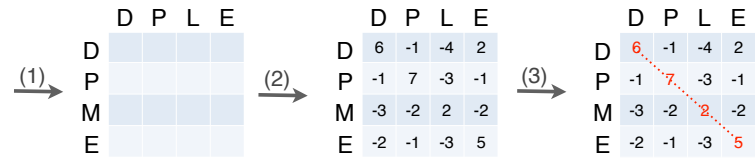|   |   | A | G | T | T | C |
|---|---|---|---|---|---|---|
|   | 0 |   |   |   |   |   |
| A |   |   |   |   |   |   |
| T |   |   |   |   |   |   |
| T |   |   |   |   |   |   |
| G |   |   |   |   |   |   |
| C |   |   |   |   |   |   |

---

# ALIGNMENT FOUNDATIONS

- **Why…**
  - Why compare biological sequences?
- **What…**
  - Alignment view of sequence changes during evolution (matches, mismatches and gaps)
- **How…**
  - ‣ Dot matrices
  - ‣ Dynamic programing
    - Global alignment
    - Local alignment
  - ‣ BLAST heuristic approach

---

# Global vs local alignments

- Needleman-Wunsch is a global alignment algorithm
  - Resulting alignment spans the complete sequences end to end
  - This is appropriate for closely related sequences that are similar in length
- For many practical applications we require local alignments
  - Local alignments highlight sub-regions (e.g. protein domains) in the two sequences that align well

Global

Local



112

# Local alignment: Definition

- Smith & Waterman proposed simply that a local alignment of two sequences allow arbitrary-length segments of each sequence to be aligned, with no penalty for the unaligned portions of the sequences. Otherwise, the score for a local alignment is calculated the same way as that for a global alignment

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

113

# The Smith-Waterman algorithm

- Three main modifications to Needleman-Wunsch:
  - Allow a node to start at 0
  - The score for a particular cell cannot be negative
    - if all other score options produce a negative value, then a zero must be inserted in the cell
  - Record the highest- scoring node, and trace back from there

$$S(i, j) = Max \begin{cases} S(i-1, j-1) + (mis)match & \text{①} \\ S(i-1, j) - gap\ penalty & \text{②} \\ S(i, j-1) - gap\ penalty & \text{③} \\ 0 & \text{④} \end{cases}$$



114

# Local alignments can be used for database searching

- Goal: Given a query sequence (Q) and a sequence database (D), find a list of sequences from D that are most similar to Q
  - Input: Q, D and scoring scheme
  - Output: Ranked list of hits



| Score | Ranked hit list | Annotation |
|---|---|---|
| 100 | **GTATGGTCA** | Ras |
| 90 | **TGATGGTCA** | Ras |
| 40 | CGATCTGCA | HSP90 |
| 38 | TCGTTGCTA | P450 |

116



Sequence 1

Local alignment
```
GCC-UCG
GCCAUUG
```

## The database search problem

- Due to the rapid growth of sequence databases, search algorithms have to be both efficient and sensitive
  - Time to search with SW is proportional to m x n (m is length of query, n is length of database), too slow for large databases!



To reduce search time heuristic algorithms, such as BLAST, first remove database sequences without a strong local similarity to the query sequence in a quick initial scan.

117

---

## The database search problem

- Due to the rapid growth of sequence databases, search algorithms have to be both efficient and sensitive
  - Time to search with SW is proportional to m x n (m is length of query, n is length of database), too slow for large databases!

Query **RGGVKRIKLMR**

GAQRGLA → Smaller database
**RGGVKRI**
**FKLLGRI**
**MGLGVKA**          **MGLGVKA**
MPQRGLA          **RGGVKRI**

To reduce search time heuristic algorithms, such as BLAST, first remove database sequences without a strong local similarity to the query sequence in a quick initial scan.

118

---

# ALIGNMENT FOUNDATIONS

- **Why…**
  - Why compare biological sequences?
- **What…**
  - Alignment view of sequence changes during evolution (matches, mismatches and gaps)
- **How…**
  - Dot matrices
  - Dynamic programing
    - Global alignment
    - Local alignment
  - BLAST heuristic approach

---

# Rapid, heuristic versions of Smith–Waterman: BLAST

- BLAST (Basic Local Alignment Search Tool) is a simplified form of Smith-Waterman (SW) alignment that is popular because it is fast and easily accessible
  - BLAST is a heuristic approximation to SW - It examines only part of the search space
  - BLAST saves time by restricting the search by scanning database sequences for likely matches before performing more rigorous alignments
  - Sacrifices some sensitivity in exchange for speed
  - In contrast to SW, BLAST is not guaranteed to find optimal alignments

120

## Rapid, heuristic versions of Smith–Waterman: BLAST

- BLAST (Basic Local Alignment Search To~~~
  simplified form of Smith-Waterman~~~
  that is popular because it is f~~~
  - BLAST finds regions~~~
    sequences
  - BLAST~~~ne search by scanning
    ~~~ly matches before performing
    ~~~nts
  - ~~~me sensitivity in exchange for speed
  - ~~~ntrast to SW, BLAST is not guaranteed to find optimal alignments

*"The central idea of the BLAST algorithm is to confine attention to sequence pairs that contain an initial word pair match"*
*Altschul et al. (1990)*

121

---

- BLAST uses this pre-screening heuristic approximation resulting in an an approach that is about 50 times faster than the Smith-Waterman

Query **RGGVKRIKLMR**

word match

GAQRGLA
**RGGVKRI**
FKLLGRI
**MGLGVKA**
MPQRGLA

**MGLGVKA**
**RGGVKRI**

Initial Database → Smaller Database

Speed

Exact SW    BLAST

122

---

## How BLAST works

- Four basic phases
  - Phase 1: compile a list of query word pairs (w=3)

**RGGVKRI**    Query sequence
**RGG**
  **GGV**
   **GVK**
    **VKR**
     **KRI**

generate list of w=3 words for query

123

---

## Blast

- Phase 2: expand word pairs to include those similar to query (defined as those above a similarity threshold to original word, i.e. match scores in substitution matrix)

**RGGVKRI**    Query sequence
**RGG RAG RIG RLG ...**
  **GGV GAV GTV GCV ...**
   **GVK GAK GIK GGK ...**
    **VKR VRR VHR VER ...**
     **KRI KKI KHI KDI ...**

extend list of words similar to query

124

# Blast

– Phase 3: a database is scanned to find sequence entries that match the compiled word list

**GNYGLKVISLDVE**   Database sequence

**RGGVKRI**   Query sequence

RGG RAG RIG RLG ...

GGV GAV GTV GCV ...

GVK GLK GIK GGK ...

VKR VRR VHR VER ...

KRI KKI KHI KDI ...

search for perfect matches in the database sequence

125

# Blast

– Phase 4: the initial database hits are extended in both directions using dynamic programing

**GNYGLKVISLDVE**   Database sequence

**RGGVKRI**   Query sequence

matched word is used as a local alignment seed



126



Alignment seed

**GRGGVKRISGL**   Query sequence

**GNYGLKVISLDV**   Database sequence

127



dynamic programming

Search for high scoring gapped alignment

Alignment seed

**GRGGVKRISGL**   Query sequence

**GNYGLKVISLDV**   Database sequence

128

## Slide 129



BLAST returns the highest scoring database hits in a ranked list

Alignment seed

GRGGVKRISGL — Query sequence
GNYGLKVIS-L — Database sequence

## BLAST output

- BLAST returns the highest scoring database hits in a ranked list along with details about the target sequence and alignment statistics

| Description | Max score | Query cover | E value | Max ident | Accession |
|---|---|---|---|---|---|
| kinesin-1 heavy chain [Homo sapiens] | 677 | 100% | 0 | 100% | NP_004512.1 |
| Kif5b protein [Mus musculus] | 676 | 100% | 0 | 98% | AAA20133.1 |
| Kinesin-14 heavy chain [Danio rerio] | 595 | 88% | 0 | 78% | XP_00320703 |
| hypothetical protein EGK_18589 | 48.2 | 40% | 0.03 | 32% | ELK35081.1 |
| mKIAA4102 protein [Mus musculus] | 42.7 | 38% | 3.02 | 24% | EHH28205.1 |

## Statistical significance of results

- An important feature of BLAST is the computation of statistical significance for each hit. This is described by the E value (expect value)

| Description | Max score | Query cover | E value | Max ident | Accession |
|---|---|---|---|---|---|
| kinesin-1 heavy chain [Homo sapiens] | 677 | 100% | 0 | 100% | NP_004512.1 |
| Kif5b protein [Mus musculus] | 676 | 100% | 0 | 98% | AAA20133.1 |
| Kinesin-14 heavy chain [Danio rerio] | 595 | 88% | 0 | 78% | XP_00320703 |
| hypothetical protein EGK_18589 | 48.2 | 40% | 0.03 | 32% | ELK35081.1 |
| mKIAA4102 protein [Mus musculus] | 42.7 | 38% | 3.02 | 24% | EHH28205.1 |

## BLAST scores and E-values

- The E value is the expected number of hits that are as good or better than the observed local alignment score (with this score or better) if the query and database are random with respect to each other
  - i.e. the number of alignments expected to occur by chance with equivalent or better scores
- Typically, only hits with E value below a significance threshold are reported
  - This is equivalent to selecting alignments with score above a certain score threshold

**Slide 133**

- Ideally, a threshold separates all query related sequences (yellow) from all unrelated sequences (gray)



Alignment scores of unrelated sequences

Threshold

Alignment scores of related sequences

Number of sequences

Local alignment score

133

**Slide 134**

- Unfortunately, often both score distributions overlap
  - The E value describes the expected number of hits with a score above the threshold if the query and database are unrelated



Alignment scores of unrelated sequences

Threshold

Alignment scores of related sequences

Number of sequences

Local alignment score

134

**Slide 135**

- Unfortunately, often both score distributions overlap
  - The E value describes the expected number of hits with a score above the threshold if the query and database are unrelated



Alignment scores of unrelated sequences

Threshold

Number of sequences

The E-value provides an estimate of the number of false positive hits!

Local alignment score

135

**Slide 136**

| Description | Max score | Query cover | E value | Max ident | Accession |
|---|---|---|---|---|---|
| kinesin-1 heavy chain [Homo sapiens] | 677 | 100% | 0 | 100% | NP_004512.1 |
| Kif5b protein [Mus musculus] | 676 | 100% | 0 | 98% | AAA20133.1 |
| Kinesin-14 heavy chain [Danio rerio] | 595 | 88% | 0 | 78% | XP_00320703 |
| hypothetical protein EGK_18589 | 42.7 | 40% | 0.03 | 32% | ELK35081.1 |



Alignment scores of unrelated sequences

Number of sequences

A score of 42.7 or better is expected to occur by chance 3 in 100 times (E-value = 0.03)

Local alignment score

42.7

136

**Slide 1 (top-left):**

| Description | Max score | Total score | Query cover | E value | Max ident | Accession |
|---|---|---|---|---|---|---|
| kinesin-1 heavy chain [Homo | 677 | 677 | 100% | 0 | 100% | NP_004512.1 |
| Kif5b protein [Mus musculus] | 676 | 676 | 100% | 0 | 98% | AAA20133.1 |

In general *E* values < 0.005 are usually significant.

To find out more about *E* values see: "*The Statistics of Sequence Similarity Scores*" available in the help section of the NCBI BLAST site:

http://www.ncbi.nlm.nih.gov/blast/tutorial/Altschul-1.html

Num

0

Local alignment score

42.7

137

---

**Slide 2 (top-right):**

# Your Turn!

## Hands-on worksheet **Sections 6**

‣ Please do answer the last lab review question (**Q19**).
‣ We encourage discussion and exploration!

---

**Slide 3 (bottom-left):**

# Summary

• Sequence alignment is a fundamental operation underlying much of bioinformatics.

• Even when optimal solutions can be obtained they are not necessarily unique or reflective of the biologically correct alignment.

• Dynamic programming is a classic approach for solving the pairwise alignment problem.

• Reviewed global and local alignment approaches and their major application areas.

• Heuristic approaches are often necessary for large database searches and many genomic applications.

---

**Slide 4 (bottom-right):**

# FOR NEXT CLASS…

Check out the online:

☑ **Reading**: Sean Eddy's "What is dynamic programming?"
☑ **Homework**: (1) **Quiz**, (2) **Alignment Exercise**.

## Homework Grading
Both (1) quiz questions and (2) alignment exercise carry equal weights (*i.e.* 50% each).

| (Homework 2)  Assessment Criteria | Points | |
|---|---|---|
| Setup labeled alignment matrix | 1 | |
| Include initial column and row for GAPs | 1 | |
| All alignment matrix elements scored (i.e. filled in) | 1 | |
| Evidence for correct use of scoring scheme | 1 | |
| Direction arrows drawn between all cells | 1 | |
| Evidence of multiple arrows to a given cell if appropriate | 1 | D |
| Correct optimal score position in matrix used | 1 | C |
| Correct optimal score obtained for given scoring scheme | 1 | B |
| Traceback path(s) clearly highlighted | 1 | A |
| Correct *alignment(s)* yielding optimal score listed | 1 | A+ |

---

# REFERENCE SLIDES...

**Additional reference slides for the motivated student**

---

## Practical database searching with BLAST

- There are four basic components to a traditional BLAST search
  - (1) Choose the sequence (query)
  - (2) Select the BLAST program
  - (3) Choose the database to search
  - (4) Choose optional parameters
- Then click "BLAST"

---

## Step 1: Choose your sequence

- Sequence can be input in FASTA format or as accession number

## Step 2: Choose the BLAST program

| Query | | Database |
|-------|---|----------|
| **blastn** | DNA → (1) | DNA |
| **blastp** | protein → (1) | protein |
| **blastx** | DNA ← (6) → | protein |
| **tblastn** | protein → (6) | DNA |
| **tblastx** | DNA ← (36) → | DNA |

## DNA potentially encodes six proteins

```
        5' CAT CAA
         5' ATC AAC
          5' TCA ACT

5' CATCAACTACAACTCCAAAGACACCCTTACACATCAACAAACCTACCCAC 3'
3' GTAGTTGATGTTGAGGTTTCTGTGGGAATGTGTAGTTGTTTGGATGGGTG 5'

                        5' GTG GGT
                         5' TGG GTA
                          5' GGG TAG
```
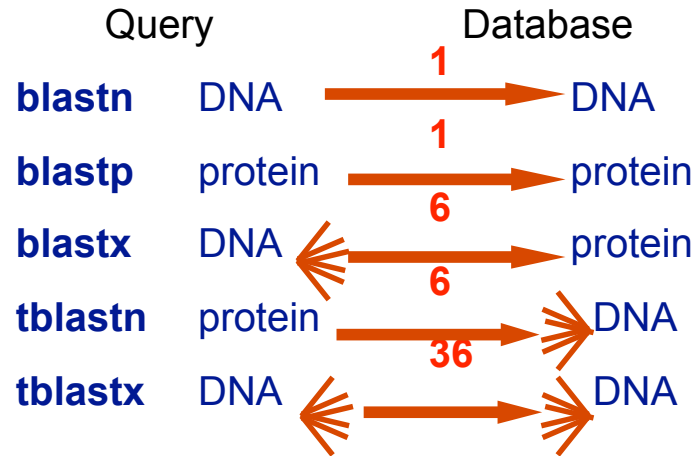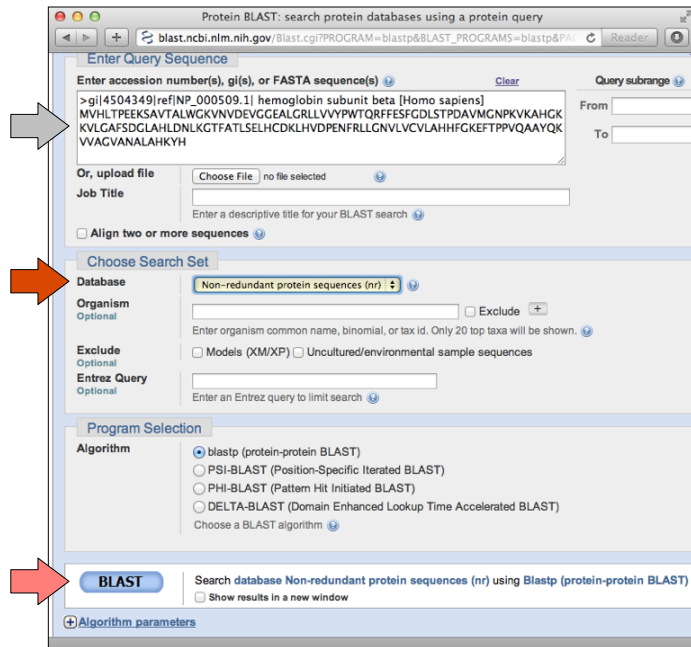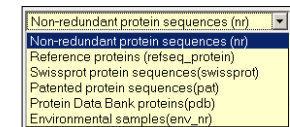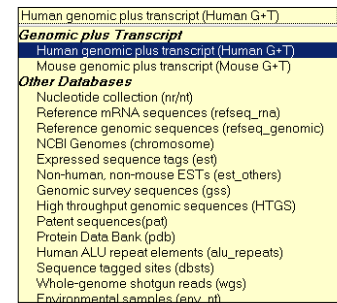
Protein BLAST: search protein databases using a protein query

blast.ncbi.nlm.nih.gov/Blast.cgi?PROGRAM=blastp&BLAST_PROGRAMS=blastp&PA...

Enter Query Sequence

Enter accession number(s), gi(s), or FASTA sequence(s)   Clear        Query subrange

>gi|4504349|ref|NP_000509.1| hemoglobin subunit beta [Homo sapiens]   From
MVHLTPEEKSAVTALWGKVNVDEVGGEALGRLLVVYPWTQRFFESFGDLSTPDAVMGNPKVKAHGK
KVLGAFSDGLAHLDNLKGTFATLSELHCDKLHVDPENFRLLGNVLVCVLAHHFGKEFTPPVQAAYQK   To
VVAGVANALAHKYH

Or, upload file   Choose File   no file selected

Job Title        Enter a descriptive title for your BLAST search

☐ Align two or more sequences

Choose Search Set

Database   Non-redundant protein sequences (nr) ▼
Organism   Optional                            ☐ Exclude +
           Enter organism common name, binomial, or tax id. Only 20 top taxa will be shown.
Exclude    Optional   ☐ Models (XM/XP)  ☐ Uncultured/environmental sample sequences
Entrez Query   Optional
           Enter an Entrez query to limit search

Program Selection

Algorithm   ● 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

BLAST   Search database Non-redundant protein sequences (nr) using Blastp (protein-protein BLAST)
        ☐ Show results in a new window
⊞ Algorithm parameters
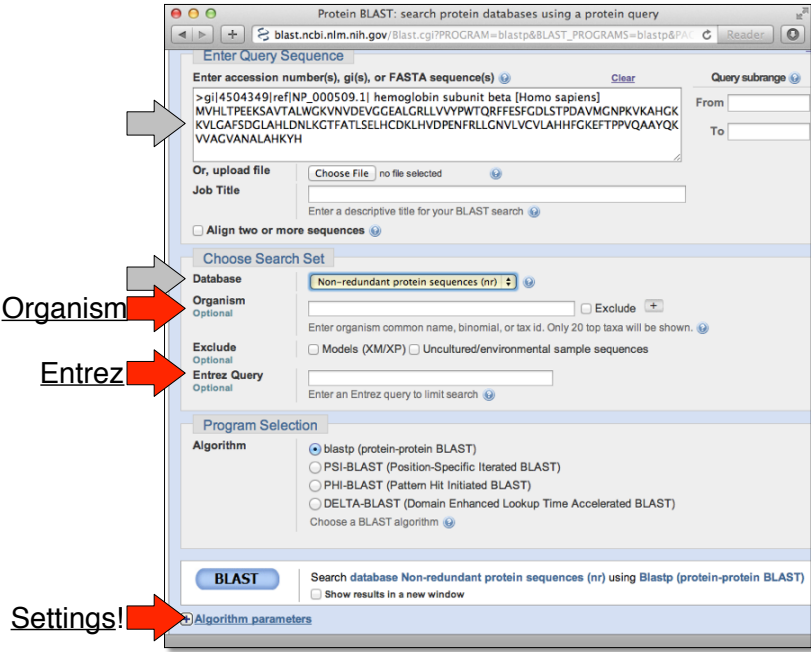
## Step 3: Choose the database

nr = non-redundant (most general database)

dbest = database of expressed sequence tags

dbsts = database of sequence tag sites

gss = genomic survey sequences



Human genomic plus transcript (Human G+T) ▼
**Genomic plus Transcript**
  Human genomic plus transcript (Human G+T)
  Mouse genomic plus transcript (Mouse G+T)
**Other Databases**
  Nucleotide collection (nr/nt)
  Reference mRNA sequences (refseq_rna)
  Reference genomic sequences (refseq_genomic)
  NCBI Genomes (chromosome)
  Expressed sequence tags (est)
  Non-human, non-mouse ESTs (est_others)
  Genomic survey sequences (gss)
  High throughput genomic sequences (HTGS)
  Patent sequences(pat)
  Protein Data Bank (pdb)
  Human ALU repeat elements (alu_repeats)
  Sequence tagged sites (dbsts)
  Whole-genome shotgun reads (wgs)
  Environmental samples (env_nt)

nucleotide databases

Non-redundant protein sequences (nr) ▼
Non-redundant protein sequences (nr)
Reference proteins (refseq_protein)
Swissprot protein sequences(swissprot)
Patented protein sequences(pat)
Protein Data Bank proteins(pdb)
Environmental samples(env_nr)

protein databases

# Slide 1 (top-left)

Protein BLAST: search protein databases using a protein query

blast.ncbi.nlm.nih.gov/Blast.cgi?PROGRAM=blastp&BLAST_PROGRAMS=blastp&PA...

**Enter Query Sequence**

Enter accession number(s), gi(s), or FASTA sequence(s)    Clear    Query subrange

>gi|4504349|ref|NP_000509.1| hemoglobin subunit beta [Homo sapiens]
MVHLTPEEKSAVTALWGKVNVDEVGGEALGRLLVVYPWTQRFFESFGDLSTPDAVMGNPKVKAHGK
KVLGAFSDGLAHLDNLKGTFATLSELHCDKLHVDPENFRLLGNVLVCVLAHHFGKEFTPPVQAAYQK
VVAGVANALAHKYH

From
To

Or, upload file    Choose File  no file selected
Job Title                Enter a descriptive title for your BLAST search
☐ Align two or more sequences

**Choose Search Set**

Database    Non-redundant protein sequences (nr)
Organism    ☐ Exclude  +
Optional    Enter organism common name, binomial, or tax id. Only 20 top taxa will be shown.
Exclude     ☐ Models (XM/XP) ☐ Uncultured/environmental sample sequences
Optional
Entrez Query
Optional    Enter an Entrez query to limit search

**Program Selection**

Algorithm   ◉ 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

**BLAST**   Search database Non-redundant protein sequences (nr) using Blastp (protein-protein BLAST)
            ☐ Show results in a new window

⊕ Algorithm parameters

Organism
Entrez
Settings!

149

# Slide 2 (top-right): Step 4a: Select optional search parameters

▼ 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    **Expect**
Word size    3    **Word size**
Max matches in a query range    0

**Scoring Parameters**

Matrix    BLOSUM62    **Scoring matrix**
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

**BLAST**   Search database Non-redundant protein sequences (nr) using Blastp
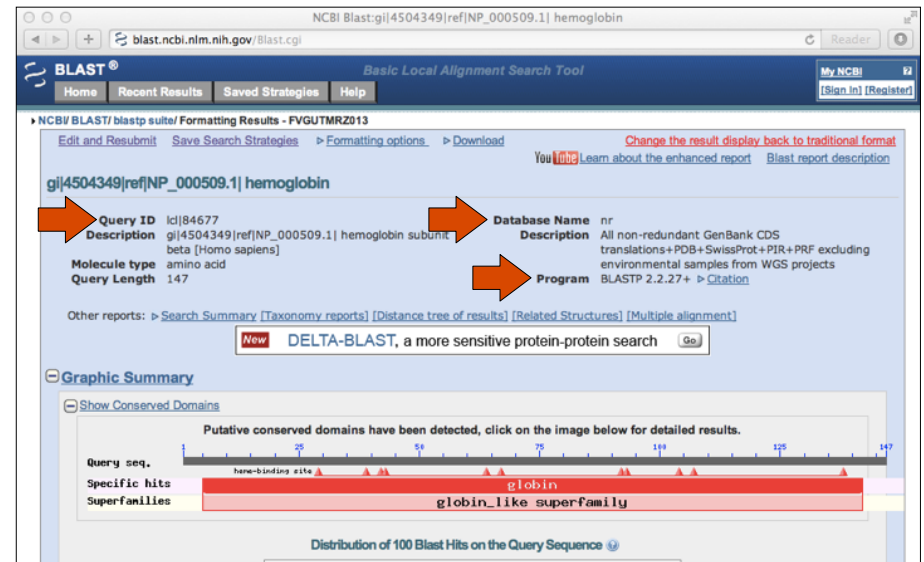            ☐ Show results in a new window

150

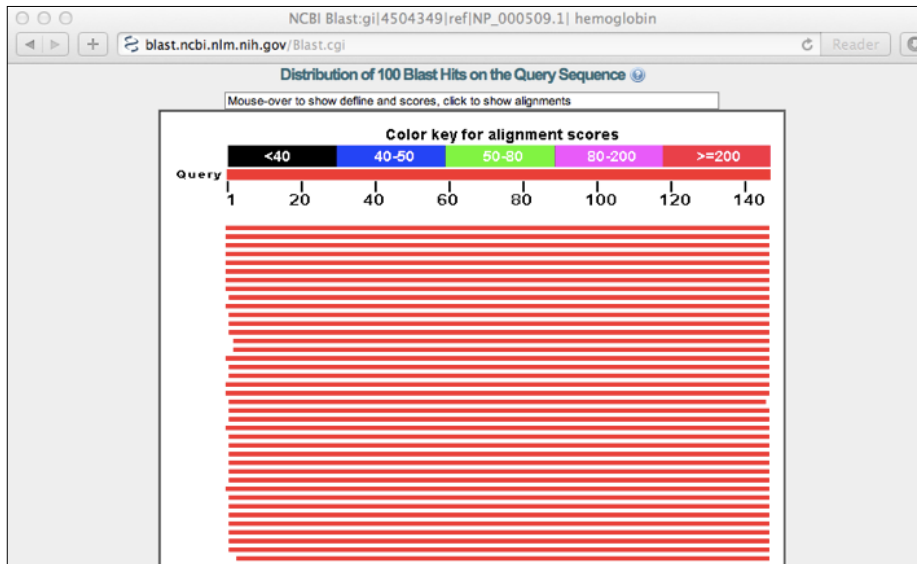# Slide 3 (bottom-left): Step 4: Optional parameters

- You can...
  - choose the organism to search
  - change the substitution matrix
  - change the expect (E) value
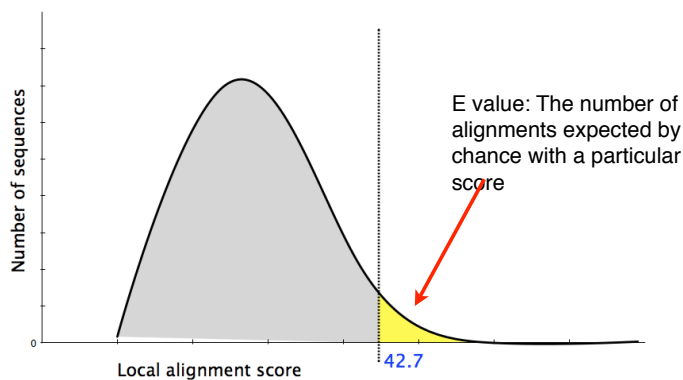  - change the word size
  - change the output format

151

# Slide 4 (bottom-right): Results page

NCBI Blast:gi|4504349|ref|NP_000509.1| hemoglobin

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

**BLAST®**    Basic Local Alignment Search Tool    My NCBI
Home  Recent Results  Saved Strategies  Help    [Sign In] [Register]

NCBI/ BLAST/ blastp suite/ Formatting Results - FVGUTMRZ013

Edit and Resubmit   Save Search Strategies   ▷ Formatting options   ▷ Download    Change the result display back to traditional format
You Tube Learn about the enhanced report    Blast report description

gi|4504349|ref|NP_000509.1| hemoglobin

Query ID  lcl|84677                              Database Name  nr
Description  gi|4504349|ref|NP_000509.1| hemoglobin subunit    Description  All non-redundant GenBank CDS
             beta [Homo sapiens]                                translations+PDB+SwissProt+PIR+PRF excluding
Molecule type  amino acid                                      environmental samples from WGS projects
Query Length  147                               Program  BLASTP 2.2.27+ ▷ Citation

Other reports: ▷ Search Summary [Taxonomy reports] [Distance tree of results] [Related Structures] [Multiple alignment]

New  DELTA-BLAST, a more sensitive protein-protein search  Go

⊟ Graphic Summary

⊟ Show Conserved Domains

Putative conserved domains have been detected, click on the image below for detailed results.

Query seq.
Specific hits    globin
Superfamilies    globin_like superfamily

Distribution of 100 Blast Hits on the Query Sequence

## Further down the results page...



## Further down the results page...



| Description | Max score | Total score | Query cover | E value | Max ident | Accession |
|---|---|---|---|---|---|---|
| kinesin-1 heavy chain [Homo | 677 | 677 | 100% | 0 | 100% | NP_004512.1 |
| Kif5b protein [Mus musculus] | 676 | 676 | 100% | 0 | 98% | AAA20133.1 |
| Kinesin-14 heavy chain [Danio rerio] | 595 | 595 | 88% | 0 | 78% | XP_00320703 |
| hypothetical protein EGK_18589 | 48.2 | 52 | 40% | 0.03 | 32% | ELK35081.1 |
| mKIAA4102 protein [Mus musculus] | 42.7 | 42.7 | 38% | 3.02 | 24% | EHH28205.1 |

E value: The number of alignments expected by chance with a particular score



155

## E values in BLAST

- Each alignment gets a score determined from the alignment and doesn't take into account the full length of the query, target or database
- The E value is what you want to look at
- E value = Expect
  - How often do I expect an alignment with this score give the length of my query and the size of the database
  - $E = Kmne^{-\lambda s}$
    - K and $\lambda$ are scaling factors
    - S is the score
    - m – length of query, n – length of database
  - E corrects for multiple comparisons, i.e., query compared to many sequences – proportional to length of database and query for a given S (score)

# Further down the results page...



# Different output formats are available



# E.g. Query anchored alignments



# ... and alignments with dots for identities

## Common problems

- Selecting the wrong version of BLAST
- Selecting the wrong database
- Too many hits returned
- Too few hits returned
- Unclear about the significance of a particular result - are these sequences homologous?

161

## How to handle too many results

- Focus on the question you are trying to answer
  - select "refseq" database to eliminate redundant matches from "nr"
  - Limit hits by organism
  - Use just a portion of the query sequence, when appropriate
  - Adjust the expect value; lowering E will reduce the number of matches returned

162

## How to handle too few results

- Many genes and proteins have no significant database matches
  - remove Entrez limits
  - raise E-value threshold
  - search different databases
  - try scoring matrices with lower BLOSUM values (or higher PAM values)
  - use a search algorithm that is more sensitive than BLAST (e.g. PSI-BLAST or HMMer)

163