***Current Topics in Genome Analysis***
***Spring 2005***

***Week 4***
***Biological Sequence Analysis I***

*Andy Baxevanis, Ph.D.*

## Overview

- Week 4: Comparative methods and concepts
  - Similarity *vs*. Homology
  - Global *vs*. Local Alignments
  - Scoring Matrices
  - BLAST
  - BLAT

- Week 5: Predictive methods and concepts
  - Profiles, patterns, motifs, and domains
  - Secondary structure prediction
  - Structures: VAST, Cn3D, and *de novo* prediction

# Why do sequence alignments?

- Provide a measure of relatedness between nucleotide or amino acid sequences

- Determining relatedness allows one to draw biological inferences regarding
  - structural relationships
  - functional relationships
  - evolutionary relationships

    → *importance of using correct terminology*

# Defining the Terms

- The quantitative measure: ***Similarity***
  - Always based on an observable
  - Usually expressed as percent identity
  - Quantify changes that occur as two sequences diverge
    - substitutions
    - insertions
    - deletions
  - Identify residues crucial for maintaining a protein's structure or function

- High degrees of sequence similarity *might* imply
  - a common evolutionary history
  - possible commonality in biological function

## Defining the Terms

- The conclusion: ***Homology***
  - Genes *are* or *are not* homologous
    (not measured in degrees)
  - Homology implies an evolutionary relationship

- The term "homolog" may apply to the relationship
  - between genes separated by the event of speciation
    (*orthology*)
  - between genes separated by the event of genetic
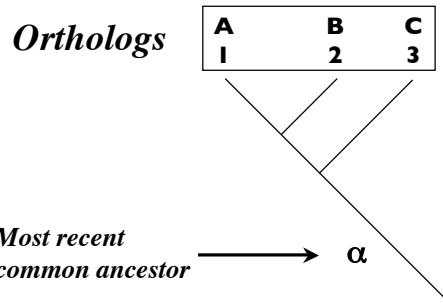    duplication (*paralogy*)
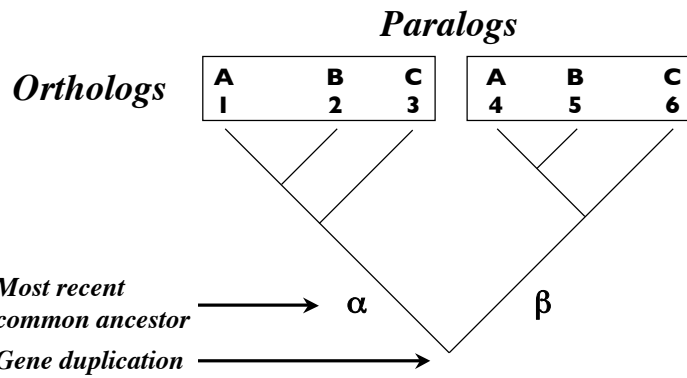
## Defining the Terms

- Orthologs
  - Sequences are direct descendants of a sequence in a
    common ancestor
  - Most likely have similar domain structure, three-
    dimensional structure, and biological function

- Paralogs
  - Related through a gene duplication event
  - Provides insight into "evolutionary innovation"
    (adapting a pre-existing gene product for a new
    function)

# Defining the Terms

*Orthologs*

| A | B | C |
|---|---|---|
| 1 | 2 | 3 |

*Most recent common ancestor* → α

# Defining the Terms

*Paralogs*

*Orthologs*

| A | B | C |
|---|---|---|
| 1 | 2 | 3 |

| A | B | C |
|---|---|---|
| 4 | 5 | 6 |

*Most recent common ancestor* → α    β

*Gene duplication* →

- Genes 1-3 are orthologous
- Genes 4-6 are orthologous
- Any pair of α and β genes are paralogous
  (genes related through a gene duplication event)

## Overview

- Week 4: Comparative methods and concepts
  - Similarity *vs.* Homology
  - Global *vs.* Local Alignments
  - Scoring Matrices
  - BLAST
  - BLAT

- Week 5: Predictive methods and concepts
  - Profiles, patterns, motifs, and domains
  - Secondary structure prediction
  - Structures: VAST, Cn3D, and *de novo* prediction

## Global Sequence Alignments

- Sequence comparison along the entire length of the two sequences being aligned
- Best for highly-similar sequences of similar length
- As the degree of sequence similarity declines, global alignment methods tend to miss important biological relationships

## Local Sequence Alignments

- Sequence comparison intended to find the most similar regions in the two sequences being aligned ("paired subsequences")

- Regions outside the area of local alignment are excluded

- More than one local alignments could be generated for any two sequences being compared

- Best for sequences that share some similarity, or for sequences of different lengths

## Overview

- Week 4: Comparative methods and concepts
  - Similarity *vs.* Homology
  - Global *vs.* Local Alignments
  - Scoring Matrices
  - BLAST
  - BLAT

- Week 5: Predictive methods and concepts
  - Profiles, patterns, motifs, and domains
  - Secondary structure prediction
  - Structures: VAST, Cn3D, and *de novo* prediction

## Scoring Matrices

- Empirical weighting scheme to represent biology (side chain chemistry, structure, and function)
  - Cys/Pro important for structure and function
  - Trp has bulky side chain
  - Lys/Arg have positively-charged side chains

## Scoring Matrices

- *Conservation:* What residues can substitute for another residue and not adversely affect the function of the protein?
  - Ile/Val - both small and hydrophobic
  - Ser/Thr - both polar
  - *Conserve charge, size, hydrophobicity, other physicochemical factors*

- *Frequency:* How often does a particular residue occur amongst the entire constellation of proteins?

# Scoring Matrices

- Importance of understanding scoring matrices
  - Appear in all analyses involving sequence comparison
  - Implicitly represent particular evolutionary patterns
  - Choice of matrix can strongly influence outcomes

# Matrix Structure: Nucleotides

|   | A | T | G | C | S | W | R | Y | K | M | B | V | H | D | N |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| A | 5 | -4 | -4 | -4 | -4 | 1 | 1 | -4 | -4 | 1 | -4 | -1 | -1 | -1 | -2 |
| T | -4 | 5 | -4 | -4 | -4 | 1 | -4 | 1 | 1 | -4 | -1 | -4 | -1 | -1 | -2 |
| G | -4 | -4 | 5 | -4 | 1 | -4 | 1 | -4 | 1 | -4 | -1 | -1 | -4 | -1 | -2 |
| C | -4 | -4 | -4 | 5 | 1 | -4 | -4 | 1 | -4 | 1 | -1 | -1 | -1 | -4 | -2 |
| S | -4 | -4 | 1 | 1 | -1 | -4 | -2 | -2 | -2 | -2 | -1 | -1 | -3 | -3 | -1 |
| W | 1 | 1 | -4 | -4 | -4 | -1 | -2 | -2 | -2 | -2 | -3 | -3 | -1 | -1 | -1 |
| R | 1 | -4 | 1 | -4 | -2 | -2 | -1 | -4 | -2 | -2 | -3 | -1 | -3 | -1 | -1 |
| Y | -4 | 1 | -4 | 1 | -2 | -2 | -4 | -1 | -2 | -2 | -1 | -3 | -1 | -3 | -1 |
| K | -4 | 1 | 1 | -4 | -2 | -2 | -2 | -2 | -1 | -4 | -1 | -3 | -3 | -1 | -1 |
| M | 1 | -4 | -4 | 1 | -2 | -2 | -2 | -2 | -4 | -1 | -3 | -1 | -1 | -3 | -1 |
| B | -4 | -1 | -1 | -1 | -1 | -3 | -3 | -1 | -1 | -3 | -1 | -2 | -2 | -2 | -1 |
| V | -1 | -4 | -1 | -1 | -1 | -3 | -1 | -3 | -3 | -1 | -2 | -1 | -2 | -2 | -1 |
| H | -1 | -1 | -4 | -1 | -3 | -1 | -3 | -1 | -3 | -1 | -2 | -2 | -1 | -2 | -1 |
| D | -1 | -1 | -1 | -4 | -3 | -1 | -1 | -3 | -1 | -3 | -2 | -2 | -2 | -1 | -1 |
| N | -2 | -2 | -2 | -2 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |

- *Simple match/mismatch scoring scheme*
- *Assumes each nucleotide occurs 25% of the time*

## Matrix Structure: Proteins

```
      A  R  N  D  C  Q  E  G  H  I  L  K  M  F  P  S  T  W  Y  V  B  Z  X  *
A     4 -1 -2 -2  0 -1 -1  0 -2 -1 -1 -1 -1 -2 -1  1  0 -3 -2  0 -2 -1  0 -4
R    -1  5  0 -2 -3  1  0 -2  0 -3 -2  2 -1 -3 -2 -1 -1 -3 -2 -3 -1  0 -1 -4
N    -2  0  6  1 -3  0  0  0  1 -3 -3  0 -2 -3 -2  1  0 -4 -2 -3  3  0 -1 -4
D    -2 -2  1  6 -3  0  2 -1 -1 -3 -4 -1 -3 -3 -1  0 -1 -4 -3 -3  4  1 -1 -4
C     0 -3 -3 -3  9 -3 -4 -3 -3 -1 -1 -3 -1 -2 -3 -1 -1 -2 -2 -1 -3 -3 -2 -4
Q    -1  1  0  0 -3  5  2 -2  0 -3 -2  1  0 -3 -1  0 -1 -2 -1 -2  0  3 -1 -4
E    -1  0  0  2 -4  2  5 -2  0 -3 -3  1 -2 -3 -1  0 -1 -3 -2 -2  1  4 -1 -4
G     0 -2  0 -1 -3 -2 -2  6 -2 -4 -4 -2 -3 -3 -2  0 -2 -2 -3 -3 -1 -2 -1 -4
H    -2  0  1 -1 -3  0  0 -2  8 -3 -3 -1 -2 -1 -2 -1 -2 -2  2 -3  0  0 -1 -4
I    -1 -3 -3 -3 -1 -3 -3 -4 -3  4  2 -3  1  0 -3 -2 -1 -3 -1  3 -3 -3 -1 -4
L    -1 -2 -3 -4 -1 -2 -3 -4 -3  2  4 -2  2  0 -3 -2 -1 -2 -1  1 -4 -3 -1 -4
K    -1  2  0 -1 -3  1  1 -2 -1 -3 -2  5 -1 -3 -1  0 -1 -3 -2 -2  0  1 -1 -4
M    -1 -1 -2 -3 -1  0 -2 -3 -2  1  2 -1  5  0 -2 -1 -1 -1 -1  1 -3 -1 -1 -4
F    -2 -3 -3 -3 -2 -3 -3 -3 -1  0  0 -3  0  6 -4 -2 -2  1  3 -1 -3 -3 -1 -4
P    -1 -2 -2 -1 -3 -1 -1 -2 -2 -3 -3 -1 -2 -4  7 -1 -1 -4 -3 -2 -2 -1 -2 -4
S     1 -1  1  0 -1  0  0  0 -1 -2 -2  0 -1 -2 -1  4  1 -3 -2 -2  0  0  0 -4
T     0 -1  0 -1 -1 -1 -1 -2 -2 -1 -1 -1 -1 -2 -1  1  5 -2 -2  0 -1 -1  0 -4
W    -3 -3 -4 -4 -2 -2 -3 -2 -2 -3 -2 -3 -1  1 -4 -3 -2 11  2 -3 -4 -3 -2 -4
Y    -2 -2 -2 -3 -2 -1 -2 -3  2 -1 -1 -2 -1  3 -3 -2 -2  2  7 -1 -3 -2 -1 -4
V     0 -3 -3 -3 -1 -2 -2 -3 -3  3  1 -2  1 -1 -2 -2  0 -3 -1  4 -3 -2 -1 -4
B    -2 -1  3  4 -3  0  1 -1  0 -3 -4  0 -3 -3 -2  0 -1 -4 -3 -3  4  1 -1 -4
Z    -1  0  0  1 -3  3  4 -2  0 -3 -3  1 -1 -3 -1  0 -1 -3 -2 -2  1  4 -1 -4
X     0 -1 -1 -1 -2 -1 -1 -1 -1 -1 -1 -1 -1 -1 -2  0  0 -2 -1 -1 -1 -1 -1 -4
*    -4 -4 -4 -4 -4 -4 -4 -4 -4 -4 -4 -4 -4 -4 -4 -4 -4 -4 -4 -4 -4 -4 -4  1
```

BLOSUM62

## PAM Matrices

- Margaret Dayhoff and colleagues, 1978
  - Look at patterns of substitutions in highly related proteins (> 85% similar) within multiple sequence alignments
  - Analysis documented 1572 changes in 71 groups of proteins examined
  - Substitution tables constructed based on results
  - Given high degree of similarity within original sequence set, results represent substitution pattern that would be expected over short evolutionary distances

## PAM Matrices

- Short evolutionary distance
  ∴ change in function unlikely

- Point Accepted Mutation (PAM)
  - The new side chain must function the same way as the old one ("acceptance")
  - On average, 1 PAM corresponds to 1 amino acid change per 100 residues
  - 1 PAM ~ 1% divergence
  - Extrapolate to predict patterns at longer evolutionary distances

## PAM Matrices: Assumptions

- All sites assumed to be equally mutable
- Replacement of amino acids is independent of previous mutations at the same position
- Replacement is independent of surrounding residues
- Forces responsible for sequence evolution over shorter time spans are the same as those over longer time spans

# PAM Matrices: Sources of Error

- Small, globular proteins of average composition used to derive matrices

- Errors in PAM 1 are magnified up to PAM 250 (only PAM 1 is based on direct observation)

- Does not account for conserved blocks or motifs

# BLOSUM Matrices

- Henikoff and Henikoff, 1992

- Blocks Substitution Matrix

  - Look only for differences in conserved, ungapped regions of a protein family ("blocks")

  - Directly calculated, using no extrapolations

  - More sensitive to detecting structural or functional substitutions

  - Generally perform better than PAM matrices for local similarity searches *(Henikoff and Henikoff, 1993)*

# BLOSUM *n*

- Calculated from sequences sharing no more than *n%* identity

- Contribution of sequences > *n%* identical clustered and weighted to 1

```
          *      *  *  *                           TGNQEEYGNTSSDSSDEDY
TGNQEEYGNTSSDSSDEDY
KKLEKEEEGISQESSEEE                                 KKLEKEEEGISQESSEEE
KKLEKEEEGISQESSEEE          80%                     KKLEKEEEGISQESSEEE
KKLEKEEEGISQESSEEE        ───────►                 KKLEKEEEGISQESSEEE
KPAQEETEETSSQESAEED
KKPAQETEETSSQESAEED                                KPAQEETEETSSQESAEED
                                                   KKPAQETEETSSQESAEED
```

*A+T Hook Domain (Block IPB000637B)*

*2,000 blocks representing > 500 groups of related proteins*

---

# BLOSUM *n*

- Clustering reduces contribution of closely-related sequences (less bias towards substitutions that occur in the most closely related members of a family)

- Substitution frequencies are more heavily-influenced by sequences that are more divergent than this cutoff

- Reducing *n* yields more distantly-related sequences

# So many matrices...

### Triple-PAM Strategy *(Altschul, 1991)*

| | | |
|---|---|---|
| PAM 40 | Short alignments, highly similar | 70-90% |
| PAM 160 | Detecting known members of a protein family | 50-60% |
| PAM 250 | Longer, weaker local alignments | ~ 30% |

### BLOSUM *(Henikoff, 1993)*

| | | |
|---|---|---|
| BLOSUM 90 | Short alignments, highly similar | 70-90% |
| BLOSUM 80 | Detecting known members of a protein family | 50-60% |
| BLOSUM 62 | Most effective in finding all potential similarities | 30-40% |
| BLOSUM 30 | Longer, weaker local alignments | < 30% |

---

# So many matrices...

- Matrix Equivalencies

  PAM 250 ~ BLOSUM 45

  PAM 160 ~ BLOSUM 62

  PAM 120 ~ BLOSUM 80

- Specialized matrices
  - Transmembrane proteins
  - Species-specific matrices

*Wheeler, 2003*

## So many matrices...

*No single matrix is*
*the complete answer for*
*all sequence comparisons*

## Gaps

- Compensate for insertions and deletions

- Used to improve alignments between two sequences

- Must be kept to a reasonable number, to not reflect a biological implausible scenario (~1 gap per 20 residues good rule-of-thumb)

- Cannot be scored simply as a "match" or a "mismatch"

## Affine Gap Penalty

Fixed deduction for introducing a gap *plus*
an additional deduction proportional to the length of the gap

$$\text{Deduction for a gap} = G + Ln$$

|  |  |  |  | nuc | pro |
|---|---|---|---|---|---|
| where | $G$ | = | gap-opening penalty | 5 | 11 |
|  | $L$ | = | gap-extension penalty | 2 | 1 |
| and | $n$ | = | length of the gap |  |  |

Can adjust scores to make gap insertion more or less permissive, but most programs will use values of $G$ and $L$ most appropriate for the scoring matrix selected

## Overview

- Week 4: Comparative methods and concepts
  - Similarity *vs.* Homology
  - Global *vs.* Local Alignments
  - Scoring Matrices
  - BLAST
  - BLAT

- Week 5: Predictive methods and concepts
  - Profiles, patterns, motifs, and domains
  - Secondary structure prediction
  - Structures: VAST, Cn3D, and *de novo* prediction

## BLAST

- <u>B</u>asic <u>L</u>ocal <u>A</u>lignment <u>S</u>earch <u>T</u>ool

- Seeks high-scoring segment pairs (HSP)
  - pair of sequences that can be aligned without gaps
  - when aligned, have maximal aggregate score (score cannot be improved by extension or trimming)
  - score must be above score threshhold $S$
  - gapped or ungapped

- Results not limited to the "best HSP" for any given sequence pair

## BLAST Algorithms

| Program | Query Sequence | Target Sequence |
|---|---|---|
| **BLASTN** | **Nucleotide** | **Nucleotide** |
| **BLASTP** | **Protein** | **Protein** |
| **BLASTX** | **Nucleotide, six-frame translation** | **Protein** |
| TBLASTN | Protein | Nucleotide, six-frame translation |
| TBLASTX | Nucleotide, six-frame translation | Nucleotide, six-frame translation |

## Neighborhood Words

Query Word ($W$ = 3)

Query:    GSQSLAALLNKCKTPQGQRLVNQWIKQPLMDKNRIEERLNLVEAFVED

Neighborhood
Words

| | | |
|---|---|---|
| PQG | 18 | = 7 + 5 + 6 |
| PEG | 15 | |
| PRG | 14 | |
| PKG | 14 | |
| PNG | 13 | |
| PDG | 13 | |
| PHG | 13 | |
| PMG | 13 | |
| PSG | 13 | |
| PQA | 12 | |
| PQN | 12 | |
| *etc.* | | |

Neighborhood
Score Threshold
($T$ = 13)

## High-Scoring Segment Pairs

| | |
|---|---|
| PQG | 18 |
| PEG | 15 |
| PRG | 14 |
| PKG | 14 |
| PNG | 13 |
| PDG | 13 |
| PHG | 13 |
| PMG | 13 |
| PSG | 13 |
| PQA | 12 |
| PQN | 12 |
| *etc.* | |

```
Query:   325   SLAALLNKCKTPQGQRLVNQWIKQPLMDKNRIEERLNLVEA   365
               +LA++L    TP+G R++ +W+ +P+ D   + ER   + A
Sbjct:   290   TLASVLDCTVTPMGSRMLKRWLHMPVRDTRVLLERQQTIGA   330
```

## Extension

```
Query:   325   SLAALLNKCKTPQGQRLVNQWIKQPLMDKNRIEERLNLVEA   365
               +LA++L    TP+G R++ +W+ +P+ D    + ER    + A
Sbjct:   290   TLASVLDCTVTPMGSRMLKRWLHMPVRDTRVLLERQQTIGA   330
```

*Significance decay*
- *mismatches*
- *gap penalties*

X

Cumulative Score

S

T

Extension

## Scores and Probabilities

```
Query:   325   SLAALLNKCKTPQGQRLVNQWIKQPLMDKNRIEERLNLVEA   365
               +LA++L    TP+G R++ +W+ +P+ D    + ER    + A
Sbjct:   290   TLASVLDCTVTPMGSRMLKRWLHMPVRDTRVLLERQQTIGA   330
```

*Karlin-Altschul Equation*

$$E = kmNe^{-\lambda S}$$

| | |
|---|---|
| *m* | *# letters in query* |
| *N* | *# letters in database* |
| *mN* | *size of search space* |
| *λS* | *normalized score* |
| *k* | *minor constant* |

X

Cumulative Score

S

T

Extension

## Scores and Probabilities

```
Query:   325  SLAALLNKCKTPQGQRLVNQWIKQPLMDKNRIEERLNLVEA  365
              +LA++L    TP+G R++ +W+ +P+ D    + ER    + A
Sbjct:   290  TLASVLDCTVTPMGSRMLKRWLHMPVRDTRVLLERQQTIGA  330
```

$$E = kmNe^{-\lambda S}$$

*Number of HSPs found purely by chance*

*Lower values signify higher similarity*

Cumulative Score — Extension

*X*

*S*

*T*

## Scores and Probabilities

```
Query:   325  SLAALLNKCKTPQGQRLVNQWIKQPLMDKNRIEERLNLVEA  365
              +LA++L    TP+G R++ +W+ +P+ D    + ER    + A
Sbjct:   290  TLASVLDCTVTPMGSRMLKRWLHMPVRDTRVLLERQQTIGA  330
```

$$E \leq 10^{-6}$$
*for nucleotides*

$$E \leq 10^{-3}$$
*for proteins*

Cumulative Score — Extension

*X*

*S*

*T*

# Low-Complexity Regions

Defined as regions of biased composition

- Homopolymeric runs
- Short-period repeats
- Subtle over-representation of several residues

```
>gi|20455478|sp|P50553|ASC1_HUMAN Achaete-scute homolog 1 (HASH1)
MESSAKMESGGAGQQPQPQPQQPFLPPAACFFATAAAAAAAAAAAAAAQSAQQQQQQQQQQQAPQLRPAA
DGQPSGGGHKSAPKQVKRQRSSSPELMRCKRRLNFSGFGYSLPQQQAAAVARRNERERNRVKLVNLGFAT
LREHVPNGAANKKMSKVETLRSAVEYIRALQQLLDEHDAVSAAFQAGVLSPTISPNYSNDLNSMAGSPVS
SYSSDEGSYDPLSPEEQELLDFTNWF
```

*Homopolymeric*
*alanine-glutamine tract*

# Identifying Low-Complexity Regions

- Biological origins and role not well-understood
  - DNA replication errors (polymerase slippage)?
  - Unequal crossing-over?

- May confound sequence analysis
  - BLAST relies on uniformly-distributed amino acid frequencies
  - Often lead to false positives
  - Filtering is advised (and usually enabled by default)

First screenshot:

RID=1106230395-25537-60155114044.BLASTQ4, Protein query – Netscape

http://www.ncbi.nlm.nih.gov/BLAST/Blast.cgi#6179901

RID=1106230395-25537-6015...

```
>gi|6179901|gb|AAF05703.1|   homeodomain transcription factor Prospero [Drosophila mela
              Length = 1403

 Score =  938 bits (2424), Expect = 0.0
 Identities = 493/627 (78%), Positives = 493/627 (78%)
```

> ≥ 25% for proteins
> ≥ 70% for nucleotides

```
Query: 777   HVATAAPRPQMHHPAPARLPTRMGGAAGHTALKSELSEKFQMLRANNNSSMMRMSGTDLE 836
             HVATAAPRPQMHHPAPARLPTRMGGAAGHTALKSELSEKFQMLRANNNSSMMRMSGTDLE
Sbjct: 777   HVATAAPRPQMHHPAPARLPTRMGGAAGHTALKSELSEKFQMLRANNNSSMMRMSGTDLE 836

Query: 837   GLADVLKSEITTSLSALVDTIVTRFVHQRRLFSKQADSVTAAAEQLNKDLLLASQILDRK 896
             GLADVLKSEITTSLSALVDTIVTRFVHQRRLFSKQADSVTAAAEQLNKDLLLASQILDRK
Sbjct: 837   GLADVLKSEITTSLSALVDTIVTRFVHQRRLFSKQADSVTAAAEQLNKDLLLASQILDRK 896

Query: 897   SPRTKVADRPQNGPTPATQSAAAMFQAPKTPQGMNPVAAAALYNSMTGPFCLPPDXXXXX 956
             SPRTKVADRPQNGPTPATQSAAAMFQAPKTPQGMNPVAAAALYNSMTGPFCLPPD
Sbjct: 897   SPRTKVADRPQNGPTPATQSAAAMFQAPKTPQGMNPVAAAALYNSMTGPFCLPPDQQQQQ 956

Query: 957   XXXXXXXXXXXXXXXXXXXXXXXXLEQNEALSLVVTPKKKRHKVTDTRITPRTVSRILAQDX 1016
                                     LEQNEALSLVVTPKKKRHKVTDTRITPRTVSRILAQD
Sbjct: 957   QTAQQQQSAQQQQQSSQQTQQQLEQNEALSLVVTPKKKRHKVTDTRITPRTVSRILAQDG 1016

Query: 1017  XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXASNGGNSNATPAQSPTRSSGGAAYHXXX 1076
                                                 ASNGGNSNATPAQSPTRSSGGAAYH
Sbjct: 1017  VVPPTGGPPSTPQQQQQQQQQQQQQQQQQQQQQASNGGNSNATPAQSPTRSSGGAAYHPQP 1076

Query: 1077  XXXXXXXXXXVSLPTSVAIPNPSLHESKVFSPYSPFFNPXXXXXXXXXXXXXXXXXXXXXXXX 1136
                       VSLPTSVAIPNPSLHESKVFSPYSPFFNP
Sbjct: 1077  PPPPPPMMPVSLPTSVAIPNPSLHESKVFSPYSPFFNPHAAAGQATAACHQHHQQHHPH 1136

Query: 1137  XXXXXXXXXXXXXXXXXALMDSRDXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXDYKTCLRAVMDAQ 1196
                              ALMDSRD                            DYKTCLRAVMDAQ
Sbjct: 1137  HQSMQLSSSPPGSLGALMDSRDSPPLPHPPSMLHPALLAAAHHGGSPDYKTCLRAVMDAQ 1196

Query: 1197  DRQSECNSADMQFDGMAPTISFYKQMQLKTEHQESLMAKHCESLTPLHSSTLTPMHLRKA 1256
             DRQSECNSADMQFDGMAPTISFYKQMQLKTEHQESLMAKHCESLTPLHSSTLTPMHLRKA
Sbjct: 1197  DRQSECNSADMQFDGMAPTISFYKQMQLKTEHQESLMAKHCESLTPLHSSTLTPMHLRKA 1256

Query: 1257  KLMFFWVRYPSSAVLKMYFPDIKFNKNNTAQLVKWFSNFREFYYIQMEKYARQAVTEGIK 1316
             KLMFFWVRYPSSAVLKMYFPDIKFNKNNTAQLVKWFSNFREFYYIQMEKYARQAVTEGIK
Sbjct: 1257  KLMFFWVRYPSSAVLKMYFPDIKFNKNNTAQLVKWFSNFREFYYIQMEKYARQAVTEGIK 1316

Query: 1317  TPDDLLIAGDSELYRVLNLHYNRNNNHIEVPQNFRFVVESTLREFFRAIQGGKDTEQSWKK 1376
             TPDDLLIAGDSELYRVLNLHYNRNNNHIEVPQNFRFVVESTLREFFRAIQGGKDTEQSWKK
Sbjct: 1317  TPDDLLIAGDSELYRVLNLHYNRNNNHIEVPQNFRFVVESTLREFFRAIQGGKDTEQSWKK 1376

Query: 1377  SIYKIISRMDDPVPEYFKSPNFLEQLE 1403
             SIYKIISRMDDPVPEYFKSPNFLEQLE
Sbjct: 1377  SIYKIISRMDDPVPEYFKSPNFLEQLE 1403
```

> — Gap
> x Low-
>    Complexity

Second screenshot:

RID=1106230395-25537-60155114044.BLASTQ4, Protein query – Netscape

http://www.ncbi.nlm.nih.gov/BLAST/Blast.cgi#6179901

RID=1106230395-25537-6015...

> No definition line ∴
> second HSP identified

```
 Score =  826 bits (2134), Expect = 0.0
 Identities = 454/704 (64%), Positives = 461/704 (65%)

Query: 1     MSSXXXXXXXXXXXXLFQPQSVSTAXXXXXXXXXXXXTPAALATHXXXXXXXXXXXXXXXX 60
             MSS          LFQPQSVSTA            TPAALATH
Sbjct: 1     MSSAAAAAAGAAGGGALFQPQSVSTANSSSSNNNNSSTPAALATHSPTSNSPVSGASSAS 60

Query: 61    XXXXXXXFGNLFGGSSAKMLNELFGRQMKQAQDATSGLPQSLDNAMLAAAMETATSAELLI 120
                    FGNLFGGSS +       +        QSLDNAMLAAAMETATSAELL
Sbjct: 61    SLLTAAFGNLFGGSSGQDAERAVWPPDEAGPGRNEWPAQSLDNAMLAAAMETATSAELLN 120

Query: 121   GSLNSTSKLLQQQHNNNSIAPANSTPMSNGTNXXXXXXXXXXXXXXXXXXXXXXXKGSRRVSA 180
             +L   ++       ++ P    TPMSNGTN                         KGSRRVSA
Sbjct: 121   LALQFHVQVAAAAAITTALLPPIGTPMSNGTNASISPGSAHSSSHSHQGVSPKGSRRVSA 180

Query: 181   CSDRSLEAAAADVAGGSPPRAASVSSLNGGASSGEQHQSQLQHDLVAHHMLRNILQGKKE 240
             CSDRSLEAAAADVAGGSPPRAASVSSLNGGASSGEQHQSQLQHDLVAHHMLRNILQGKKE
Sbjct: 181   CSDRSLEAAAADVAGGSPPRAASVSSLNGGASSGEQHQSQLQHDLVAHHMLRNILQGKKE 240

Query: 241   LMQLDQELRTAMXXXXXXXXXXXXXXHSKLXXXXXXXXXXXXXXXXXXXXMESINLIDDSEM 300
             LMQLDQELRTAM              HSKL                    MESINLIDDSEM
Sbjct: 241   LMQLDQELRTAMQQQQQQLQEKEQLHSKLNNNNNNNIAATANNNNNTTMESINLIDDSEM 300

Query: 301   ADIKIKSEPQTAPQPQQXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXHGXXXX 360
             ADIKIKSEPQTAPQPQQ                                    HG
Sbjct: 301   ADIKIKSEPQTAPQPQQSPHGSSHSSRSGSGSGSHSSMASDGSLRRKSSDSLDSHGAQDD 360

Query: 361   XXXXXXXXPTGQRSESRAPEEPQLPTKKESVDDMLDEVELLGLHSRGSDMDSLASPSHSX 420
                     PTGQRSESRAPEEPQLPTKKESVDDMLDEVELLGLHSRGSDMDSLASPSHS
Sbjct: 361   AQDEEDAAPTGQRSESRAPEEPQLPTKKESVDDMLDEVELLGLHSRGSDMDSLASPSHSD 420

Query: 421   XXXXXXXXXXXXXXXXXCVEQKTSGSGCLKKPGMDLKRARVENIVSGMRCSPSSGLAQAG 480
                              CVEQKTSGSGCLKKPGMDLKRARVENIVSGMRCSPSSGLAQAG
Sbjct: 421   MMLLDKDDVLDEDDDDDCVEQKTSGSGCLKKPGMDLKRARVENIVSGMRCSPSSGLAQAG 480

Query: 481   QLQVNGCKKRKLYQPQQHAMERYVXXXXGLNFGLNLQSMMLDQEDSESNELESPQIQQKR 540
             QLQVNGCKKRKLYQPQQHAMERYV    GLNFGLNLQSMMLDQEDSESNELESPQIQQKR
Sbjct: 481   QLQVNGCKKRKLYQPQQHAMERYVAAAAGLNFGLNLQSMMLDQEDSESNELESPQIQQKR 540

Query: 541   VEKNALKSQLRSMQEQLAEMQQKYVQLCSRMEQESXXXXXXXXXXXXXXXXXXXXNGSSDHI 600
             VEKNALKSQLRSMQEQLAEMQQKYVQLCSRMEQES                    NGSSDHI
Sbjct: 541   VEKNALKSQLRSMQEQLAEMQQKYVQLCSRMEQESECQELDQDQDVEQEQEPDNGSSDHI 600

Query: 601   ELSPSPTLTGDGDVSPNHKEETGQERXXXXXXXXXXXXXXXXXXXXXESSDSGANMLSQMMSK 660
             ELSPSPTLTGDGDVSPNHKEETGQER                     ESSDSGANMLSQMMSK
Sbjct: 601   ELSPSPTLTGDGDVSPNHKEETGQERPGSSSPSPSPLKPKTSLGESSDSGANMLSQMMSK 660
```

```
>gi|6179901|gb|AAF05703.1|   homeodomain transcription factor Prospero
           Length = 1403

 Score =  938 bits (2424), Expect = 0.0  ✓
 Identities = 493/627 (78%) ✓ Positives = 493/627 (78%)



 Score =  826 bits (2134), Expect = 0.0  ✓
 Identities = 454/704 (64%) ✓ Positives = 461/704 (65%)
```

HSP 2            HSP 1

## Suggested BLAST Cutoffs

|  | *E* value | Sequence Identity |
|---|---|---|
| Nucleotide | $\leq 10^{-6}$ | $\geq 70\%$ |
| Protein | $\leq 10^{-3}$ | $\geq 25\%$ |

- *Do not use these cutoffs blindly!*
- *Pay attention to alignments on either side of the dividing line*
- *Do not ignore biology!*

## Database Searching Artifacts

- Low-complexity regions
  - Nucleotide searches: removed with DUST  (➔ N)
  - Protein searches: removed with SEG        (➔ X)

- Repetitive elements
  - LINE, SINE, Alu
  - Automatic masking "experimental and still under development"
  - RepeatMasker
    *http://www.repeatmasker.org*

## Database Searching Artifacts

- Low-quality sequence hits
  - Expressed sequence tags (ESTs)
  - Single-pass sequence reads from large-scale sequencing (possibly with vector contaminants)

# BLAST 2 Sequences

- Finds local alignments between two protein or nucleotide sequences of interest
  - All BLAST programs available
  - Select BLOSUM and PAM matrices available for protein comparisons
  - Same affine gap costs (adjustable)
  - Input sequences can be masked

- Implementations
  - NCBI Web interface
  - bl2seq downloadable executable
    *ftp://ncbi.nlm.nih.gov/blast/executables/*

Blast 2 Sequences – Netscape

http://www.ncbi.nlm.nih.gov/blast/bl2seq/wblast2.cgi

Search

Blast 2 Sequences

NCBI | Entrez | BLAST 2 sequences | BLAST | Example | Help

**BLAST 2 SEQUENCES**

This tool produces the alignment of two given sequences using BLAST engine for local alignment.
The stand-alone executable for blasting two sequences (bl2seq) can be retrieved from NCBI ftp site
**Reference:** Tatiana A. Tatusova, Thomas L. Madden (1999), "Blast 2 sequences - a new tool for comparing protein and nucleotide sequences",
FEMS Microbiol Lett. 174:247-250

Program [blastp ▾] Matrix [BLOSUM62 ▾] ◄

PAM30
PAM70
BLOSUM80
BLOSUM62
BLOSUM45

Parameters used in BLASTN program only:
**Reward for a match:** **Penalty for a mis**

□ Use Mega BLAST   Strand option [Not Applicable ▾]

Open gap [11] and extension gap [1] penalties
gap x_dropoff [50]   expect [10.0]   word size [3]   Fi

Sequence 1 Enter accession or GI [ ] or download from file [ ] Browse...
or sequence in FASTA format from:[0]   to:[0]
>NP_008872.1| SOX-10 [Homo sapiens]
MAEEQDLSEVELSPVGSEEPRCLSPGSAPSLGPDGGGGSGLRASPGPGELGKVKKEQQDG
CIREAVSQVLSGYDWTLVPMPVRVNGASKSKPHVKRPMNAFMVWAQAARRKLADQYPHLHN
LWRLLNESDKRPFIEEAERLRMQHKKDHPDYKYQPRRRKNGKAAQGEAECPGGEAEQGGTA
HLDHRHPGEGSPMSDGNPEHPSGQSHGPPTPPTTPKTELQSGKADPKRDGRSMGEGGKPHI
ISHEVMSNMETFDVAELDQYLPPNGHPGHVSSYSAAGYGLGSALAVASGHSAWISKPPGVA◄

Sequence 2 Enter accession or GI [ ] or download from file [ ] Browse...
or sequence in FASTA format from:[0]   to:[0]
>NP_003131.1| sex determining region Y [Homo sapiens]
MQSYASAMLSVFNSDDYSPAVQENIPALRRSSSFLCTESCNSKYQCETGENSKGNVQDRVKR
SRDQRRKMALENPRMRNSEISKQLGYQWKMLTEAEKWPFFQEAQKLQAMHREKYPNYKYRPR
NCSLLPADPASVLCSEVQLDNRLYRDDCTKATHSRMEHQLGHLPPINAASSPQQRDRYSHWT

[Align] [Clear Input]

*Comments and suggestions to blast-help@ncbi.nlm.nih.gov*

---

# MegaBLAST

- Optimized for aligning very long and/or highly-similar sequences

- Good for batch nucleotide searches

- Search targets include
  - Entire eukaryotic genomes
  - Complete chromosomes and contigs from RefSeq

- Run speeds approximately 10 times faster than BLASTN
  - Adjusted word size
  - Different gap scoring scheme

# BLASTN *vs.* MegaBLAST

- Word size
    - BLASTN default     = 11
    - MegaBLAST default   = 28

- *Non-affine* gap penalties

$$\text{Deduction for a gap} = r/2 - q$$

where      $r$ = match reward      (default 1)

               $q$ = mismatch penalty      (default -2)

and       **no penalty for opening the gap**

```
>ref|NT_094219.1|Mm5_93856_33 Mus musculus chromosome 5 genomic contig, strain C57BL/6J
         Length = 194175

 Score = 5388 bits (2802), Expect = 0.0
 Identities = 2802/2802 (100%)
 Strand = Plus / Plus

Query: 2       aaaggatcccaagcagtgcagggtgagcgcccctcgggtctggggacggtcagctgcagg 61
               ||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||
Sbjct: 107035  aaaggatcccaagcagtgcagggtgagcgcccctcgggtctggggacggtcagctgcagg 107094

Query: 62      gccgggagcaccctctggggtgcccgatggggccttccgcggggcgcaggggcgaagttg 121
               ||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||
Sbjct: 107095  gccgggagcaccctctggggtgcccgatggggccttccgcggggcgcaggggcgaagttg 107154

Query: 122     ggatcgcacgcagtgagcccgagctacccagcgcgacatgctgcctcctgggcgcaatgg 181
               ||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||
Sbjct: 107155  ggatcgcacgcagtgagcccgagctacccagcgcgacatgctgcctcctgggcgcaatgg 107214

Query: 182     caccgcacatcgggcgcggctgggggttgcagaggcaactggcgcaggtggacgcccccgg 241
               |||||||||||||||||||||||||||||||||||||||||||||||||||||||||||
Sbjct: 107215  caccgcacatcgggcgcggctgggggttgcagaggcaactggcgcaggtggacgcccccgg 107274

Query: 242     gggctcagcagccccactgggacccgcgcaggtggtcaccgctggcctcctgactctctt 301
               ||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||
Sbjct: 107275  gggctcagcagccccactgggacccgcgcaggtggtcaccgctggcctcctgactctctt 107334

Query: 302     aatcgtctggaccttgctcggcaacgtcctagtgtgtgctgctatcgtccgcagccgcca 361
               ||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||
Sbjct: 107335  aatcgtctggaccttgctcggcaacgtcctagtgtgtgctgctatcgtccgcagccgcca 107394

Query: 362     cctgcgcgccaaaatgaccaacatcttcatcgtatctctggccgtctcagacctcttcgt 421
               ||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||
Sbjct: 107395  cctgcgcgccaaaatgaccaacatcttcatcgtatctctggccgtctcagacctcttcgt 107454

Query: 422     ggcattgctggtcatgccttggaaggccgtggctgaggtggccgggtactggccctttgg 481
               ||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||
Sbjct: 107455  ggcattgctggtcatgccttggaaggccgtggctgaggtggccgggtactggccctttgg 107514

Query: 482     ggcattctgcgacatctgggtggcctttgacatcatgtgctccactgcttccatcctgaa 541
```
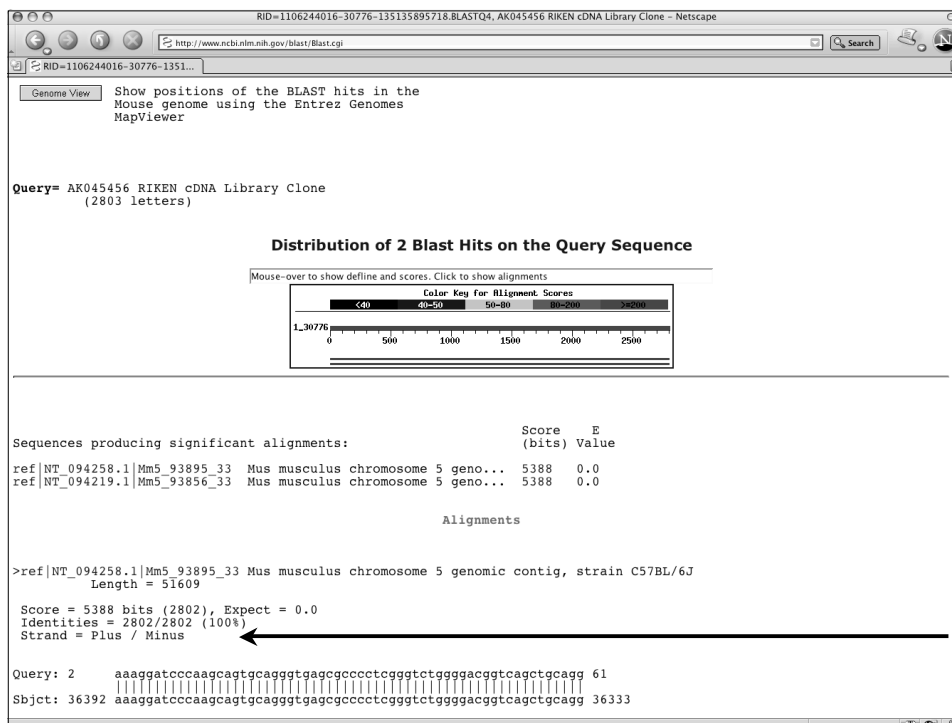


```
Genome View   Show positions of the BLAST hits in the
              Mouse genome using the Entrez Genomes
              MapViewer



Query= AK045456 RIKEN cDNA Library Clone
       (2803 letters)


              Distribution of 2 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
       1_30776
               0     500    1000    1500    2000    2500



                                             Score    E
Sequences producing significant alignments:  (bits) Value

ref|NT_094258.1|Mm5_93895_33  Mus musculus chromosome 5 geno...  5388   0.0
ref|NT_094219.1|Mm5_93856_33  Mus musculus chromosome 5 geno...  5388   0.0

                              Alignments


>ref|NT_094258.1|Mm5_93895_33 Mus musculus chromosome 5 genomic contig, strain C57BL/6J
         Length = 51609

 Score = 5388 bits (2802), Expect = 0.0
 Identities = 2802/2802 (100%)
 Strand = Plus / Minus

Query: 2       aaaggatcccaagcagtgcagggtgagcgcccctcgggtctggggacggtcagctgcagg 61
               ||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||
Sbjct: 36392  aaaggatcccaagcagtgcagggtgagcgcccctcgggtctggggacggtcagctgcagg 36333
```

## Overview

- Week 4: Comparative methods and concepts
  - Similarity *vs.* Homology
  - Global *vs.* Local Alignments
  - Scoring Matrices
  - BLAST
  - BLAT

- Week 5: Predictive methods and concepts
  - Profiles, patterns, motifs, and domains
  - Secondary structure prediction
  - Structures: VAST, Cn3D, and *de novo* prediction

# BLAT

- "BLAST-Like Alignment Tool"

- Designed to rapidly-align longer nucleotide sequences ($L \geq 40$) having > 95% sequence similarity

- Can find exact matches reliably down to $L = 33$

- Method of choice when looking for exact matches in nucleotide databases

- 500 times faster for mRNA/DNA searches

- May miss divergent or shorter sequence alignments

- Can be used on protein sequences

# When to Use BLAT

- To characterize an unknown gene or sequence fragment
  - Find its genomic coordinates
  - Determine gene structure (the presence and position of exons)
  - Identify markers of interest in the vicinity of a sequence

- To find highly-similar sequences
  - Identify gene family members
  - Identify putative homologs

- To display a specific sequence as a separate track

## Screenshot 1

**User Sequence vs Genomic – Netscape**

http://genome.ucsc.edu/cgi-bin/hgc?o=101460824&g=htcUserAli&i=../trash/hgSs_genome_1050_1106241806.pslx+..%2Ftrash%2FhgSs_genome_1050_1106241

User Sequence vs Genomic

**Alignment of CB312815**

CB312815
Rat.chr5
block1
together

### Alignment of CB312815 and chr5:101460825-101461549

Click on links in the frame to the left to navigate through the alignment. Matching bases in cDNA and genomic sequences are colored blue and capitalized. Light blue bases mark the boundaries of gaps in either sequence (often splice sites).

**cDNA CB312815**

```
GgGGCTCTCG CTGGCCTGTG TCTCAGAAGC TGCTTTCTCC ACCTCTTCCT  50
TGTGAATTTC CTAAACTCTC TACCTCTGGT TCATGTTCGC TCTTCTGGAT  100
AGTCTGTGTG CAATGAGCCC TTAAAGGAAT ATTGCAATGA GCTATAAGAG  150
TTGTGAGCCT GCGGTAGGCA AGGCCTGCAC TGGGACAGCA AAGGAAATTT  200
CATTGCATCT GCTCCTAAGT CACAGGTTAT CCAGAGCCCA CTTTACCCCA  250
AGAGACAGCC TCTCCCCCAT CCCTAGGAAA CAGTAGAGCT TAGGAAAATG  300
AATGACTCCA CCACATTCAA GAGGCTTCAA ATTGTATACT TGGCATTTCT  350
GATTTCAGTT CTGAAATTCT GTCCCTTAGT CGTGGGGAAA ATAAGAAATG  400
GAGTTACACC TTGTCATTTA AAAAACCATT GAATTAAGAG AAATGGAAAA  450
TCATGCCCAC ATAAAACATG TATGGAAGTG TTCATGTTTT GATCATGGCG  500
GGGGATATAG CTCAGTCATG GAGTGCTTGC ATAGCAATGT GCATAATCCG  550
AGGTTCAAGC CCCAGCACCG AAAAAGAGAA aCGGGAGGAG TGGAGGCATT  600
CACAGCAGCG TTTTCAGTAT AGGCGCAAAG GGGAAGGAGT TTAAACACCT  650
ACTGAGGgAA TGGATAAGCG GAGTGCCcTT GTCTATACTC GGGgatgGCT  700
AGTCATCAcg taAGAAAAGT TTGgaAAATG ATAaaatacc aatgggatgg  750
atcccccttta aaccatcc
```

**Genomic chr5 :**

```
cttggaagaa ggtaactata cattaatata gagccctctt tttctttgca  101460774
ggcccaggac acacaggacg gatgtttcca agtcactcca gggacagcat  101460824
GaGGCTCTCG CTGGCCTGTG TCTCAGAAGC TGCTTTCTCC ACCTCTTCCT  101460874
TGTGAATTTC CTAAACTCTC TACCTCTGGT TCATGTTCGC TCTTCTGGAT  101460924
AGTCTGTGTG CAATGAGCCC TTAAAGGAAT ATTGCAATGA GCTATAAGAG  101460974
TTGTGAGCCT GCGGTAGGCA AGGCCTGCAC TGGGACAGCA AAGGAAATTT  101461024
CATTGCATCT GCTCCTAAGT CACAGGTTAT CCAGAGCCCA CTTTACCCCA  101461074
AGAGACAGCC TCTCCCCCAT CCCTAGGAAA CAGTAGAGCT TAGGAAAATG  101461124
AATGACTCCA CCACATTCAA GAGGCTTCAA ATTGTATACT TGGCATTTCT  101461174
GATTTCAGTT CTGAAATTCT GTCCCTTAGT CGTGGGGAAA ATAAGAAATG  101461224
GAGTTACACC TTGTCATTTA AAAAACCATT GAATTAAGAG AAATGGAAAA  101461274
TCATGCCCAC ATAAAACATG TATGGAAGTG TTCATGTTTT GATCATGGCG  101461324
GGGGATATAG CTCAGTCATG GAGTGCTTGC ATAGCAATGT GCATAATCCG  101461374
AGGTTCAAGC CCCAGCACCG AAAAAGAGAA gCGGGAGGAG TGGAGGCATT  101461424
CACAGCAGCG TTTTCAGTAT AGGCGCAAAG GGGAAGGAGT TTAAACACCT  101461474
ACTGAGGAAT GGATAAGCGG AGTGCCTTGT CTATACTCGG CatGCTAGTC  101461524
```

Done

## Screenshot 2

**User Sequence vs Genomic – Netscape**

http://genome.ucsc.edu/cgi-bin/hgc?o=101460824&g=htcUserAli&i=../trash/hgSs_genome_1050_1106241806.pslx+..%2Ftrash%2FhgSs_genome_1050_1106241

User Sequence vs Genomic

**Alignment of CB312815**

CB312815
Rat.chr5
block1
together

**Side by Side Alignment**

```
000000001 ggggctctcgctggcctgtgtctcagaagctgctttctccacctcttcct 000000050
          >>>>>>>>> ||||||||||||||||||||||||||||||||||||||||||||||||| >>>>>>>>>
101460825 gaggctctcgctggcctgtgtctcagaagctgctttctccacctcttcct 101460874

000000051 tgtgaatttcctaaactctctacctctggttcatgttcgctcttctggat 000000100
          >>>>>>>>> |||||||||||||||||||||||||||||||||||||||||||||||||| >>>>>>>>>
101460875 tgtgaatttcctaaactctctacctctggttcatgttcgctcttctggat 101460924

000000101 agtctgtgtgcaatgagcccttaaaggaatattgcaatgagctataagag 000000150
          >>>>>>>>> |||||||||||||||||||||||||||||||||||||||||||||||||| >>>>>>>>>
101460925 agtctgtgtgcaatgagcccttaaaggaatattgcaatgagctataagag 101460974

000000151 ttgtgagcctgcggtaggcaaggcctgcactgggacagcaaaggaaattt 000000200
          >>>>>>>>> |||||||||||||||||||||||||||||||||||||||||||||||||| >>>>>>>>>
101460975 ttgtgagcctgcggtaggcaaggcctgcactgggacagcaaaggaaattt 101461024

000000201 cattgcatctgctcctaagtcacaggttatccagagcccactttacccca 000000250
          >>>>>>>>> |||||||||||||||||||||||||||||||||||||||||||||||||| >>>>>>>>>
101461025 cattgcatctgctcctaagtcacaggttatccagagcccactttacccca 101461074

000000251 agagacagcctctcccccatccctaggaaacagtagagcttaggaaaatg 000000300
          >>>>>>>>> |||||||||||||||||||||||||||||||||||||||||||||||||| >>>>>>>>>
101461075 agagacagcctctcccccatccctaggaaacagtagagcttaggaaaatg 101461124

000000301 aatgactccaccacattcaagaggcttcaaattgtatacttggcatttct 000000350
          >>>>>>>>> |||||||||||||||||||||||||||||||||||||||||||||||||| >>>>>>>>>
101461125 aatgactccaccacattcaagaggcttcaaattgtatacttggcatttct 101461174

000000351 gatttcagttctgaaattctgtcccttagtcgtggggaaataagaaatg 000000400
          >>>>>>>>> |||||||||||||||||||||||||||||||||||||||||||||||||| >>>>>>>>>
101461175 gatttcagttctgaaattctgtcccttagtcgtggggaaataagaaatg 101461224

000000401 gagttacaccttgtcatttaaaaaaccattgaattaagagaaatggaaaa 000000450
          >>>>>>>>> |||||||||||||||||||||||||||||||||||||||||||||||||| >>>>>>>>>
101461225 gagttacaccttgtcatttaaaaaaccattgaattaagagaaatggaaaa 101461274

000000451 tcatgcccacataaaacatgtatggaagtgttcatgttttgatcatggcg 000000500
          >>>>>>>>> |||||||||||||||||||||||||||||||||||||||||||||||||| >>>>>>>>>
101461275 tcatgcccacataaaacatgtatggaagtgttcatgttttgatcatggcg 101461324

000000501 ggggatatagctcagtcatggagtgcttgcatagcaatgtgcataatccg 000000550
          >>>>>>>>> |||||||||||||||||||||||||||||||||||||||||||||||||| >>>>>>>>>
101461325 ggggatatagctcagtcatggagtgcttgcatagcaatgtgcataatccg 101461374

000000551 aggttcaagccccagcaccgaaaaagagaaacgggaggagtggaggcatt 000000600
          >>>>>>>>> ||||||||||||||||||||||||||||| ||||||||||||||||||||| >>>>>>>>>
101461375 aggttcaagccccagcaccgaaaaagagaagcgggaggagtggaggcatt 101461424
```

Done

# FASTA

- Identifies regions of local alignment
- Employs an approximation of the Smith-Waterman algorithm to determine the best alignment between two sequences
- Method is significantly different from that used by BLAST
- Online implementations at
  *http://fasta.bioch.virginia.edu*
  *http://www.ebi.ac.uk/fasta33*

# Overview

- Week 4: Comparative methods and concepts
  - Similarity *vs*. Homology
  - Global *vs*. Local Alignments
  - Scoring Matrices
  - BLAST
  - BLAT

- Week 5: Predictive methods and concepts
  - Profiles, patterns, motifs, and domains
  - Secondary structure prediction
  - Structures: VAST, Cn3D, and *de novo* prediction

**Further Reading**

Altschul, S.F., Boguski, M.S., Gish, W., and Wootton, J.C.  1994.  Issues in searching molecular sequence databases. *Nat. Genet.* 6: 119-129. A review of the issues that are of importance in using sequence similarity search programs, including potential pitfalls.

Baxevanis, A.D. Assessing pairwise sequence similarity: BLAST and FASTA. In *Bioinformatics: A Practical Guide to the Analysis of Genes and Proteins*, third edition (Baxevanis, A.D. and Ouellette, B.F.F., eds.), John Wiley and Sons, 2005. An overview of the methods used to generate pairwise sequence alignments and assess the biological significance of results.

Henikoff, S. and Henikoff, J.G.  2000.  Amino acid substitution matrices. *Adv. Protein Chem.* 54: 73-97. A comprehensive review covering the factors critical to the construction of protein scoring matrices.

Korf, I., Yandell, M., and Bedell, J. BLAST. O'Reilly and Associates, 2003. An in-depth treatment of the BLAST algorithm, its applications, as well as installation, hardware, and software considerations. The book provides "documentation" that is not easily found elsewhere.

Pearson, W.R. Finding protein and nucleotide similarities with FASTA.  2003.  *Current Protocols in Bioinformatics* 3.9.1-3.9.23. An in-depth discussion of the FASTA algorithm, including worked examples and additional information regarding run options and use scenarios.

Wheeler, D.G. Selecting the right protein scoring matrix.  2003.  *Current Protocols in Bioinformatics* 3.5.1-3.5.6. A discussion of PAM, BLOSUM, and specialized scoring matrices, with guidance regarding the proper choice of matrices for particular types of protein-based analyses.

**References**

Altschul, S.F., Gish, W., Miller, W., Myers, E.W., and Lipman, D.J.  1991.  Basic local alignment search tool. *J. Mol. Biol.* 215: 403-410.

Altschul, S.F., Madden T.L., Schäffer, A.A., Zhang, J., Zhang, Z., Miller, W., Lipman, D.J.  1997. Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Res.* 25: 3389-3402.

Dayhoff, M.O., Schwartz, R.M., Orcutt, B.C.  1978.  A model of evolutionary change in proteins.  *In* Atlas of Protein Sequence and Structure, M.O. Dayhoff, ed., National Biomedical Research Foundation, Washington, 5: 345-352.

Henikoff, S. and Henikoff, J.G.  1991.  Automated assembly of protein blocks for database searching. *Nucleic Acids Res.* 19: 6565-6572.

Henikoff, S. and Henikoff, J.G.  1992.  Amino acid substitution matrices from protein blocks. *Proc. Natl. Acad. Sci. USA* 89: 10915-10919.

Henikoff, S. and Henikoff, J.G.  1993.  Performance evaluation of amino acid substitution matrices. *Proteins Struct. Funct. Genet.* 17: 49-61.

Henikoff, S. and Henikoff, J.G.  2000.  Amino acid substitution matrices. *Adv. Protein Chem.* 54: 73-97.

Karlin, S. and Altschul, S.F.  1990.  Methods for assessing the statistical significance of molecular sequence features by using general scoring schemes. *Proc. Natl. Acad. Sci. USA* 87: 2264-2268.

Kent, W.J.  2002.  BLAT: the BLAST-like alignment tool. *Genome Res.* 12: 656-664.

Pearson, W.R.  1995.  Comparison of methods for searching protein sequence databases. *Protein Sci.* 4: 1145-1160.

Pearson, W.R.  2000.  Flexible sequence similarity searching with the FASTA3 program package. *Methods Mol. Biol.* 132: 185-219.

Pearson, W.R. Finding protein and nucleotide similarities with FASTA.  2003.  *Current Protocols in Bioinformatics* 3.9.1-3.9.23.

Pearson, W.R. and Lipman, D.J.  1988.  Improved tools for biological sequence comparison. *Proc. Natl. Acad. Sci. USA* 85: 2444-2448.

Smith, T.F. and Waterman, M.S.  1981.  Identification of common molecular subsequences.  *J. Mol. Biol.* 147: 195-197.

Tatusova, T.A. and Madden, T.L.  1999.  BLAST2Sequences, a new tool for comparing protein and nucleotide sequences. *FEMS Microbio. Lett.* 174: 247-250.

Wootton, J.C. and Federhen, S.  1993.  Statistics of local complexity in amino acid sequences and sequence databases. *Comput. Chem.* 17: 149-163.