*Current Topics in Genome Analysis*
*Fall 2003*

*Week 4*
*Biological Sequence Analysis I*

*Andy Baxevanis, Ph.D.*

# Overview

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

## 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 infer
  - 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)
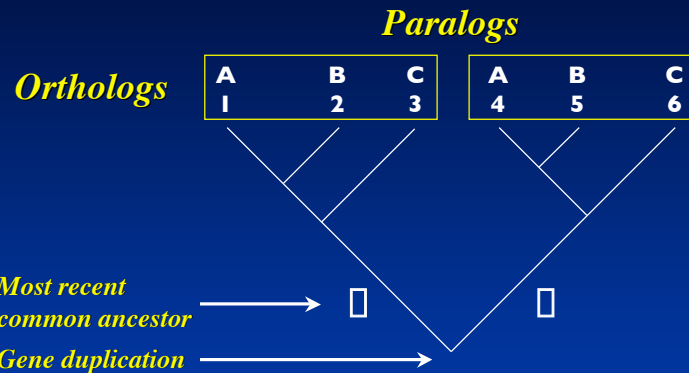
## Defining the Terms

- 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



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

## Determining Sequence Similarity

- Global sequence alignments
  - Sequence comparison along the entire length of the two sequences being aligned
  - Best for highly-similar sequences of similar length

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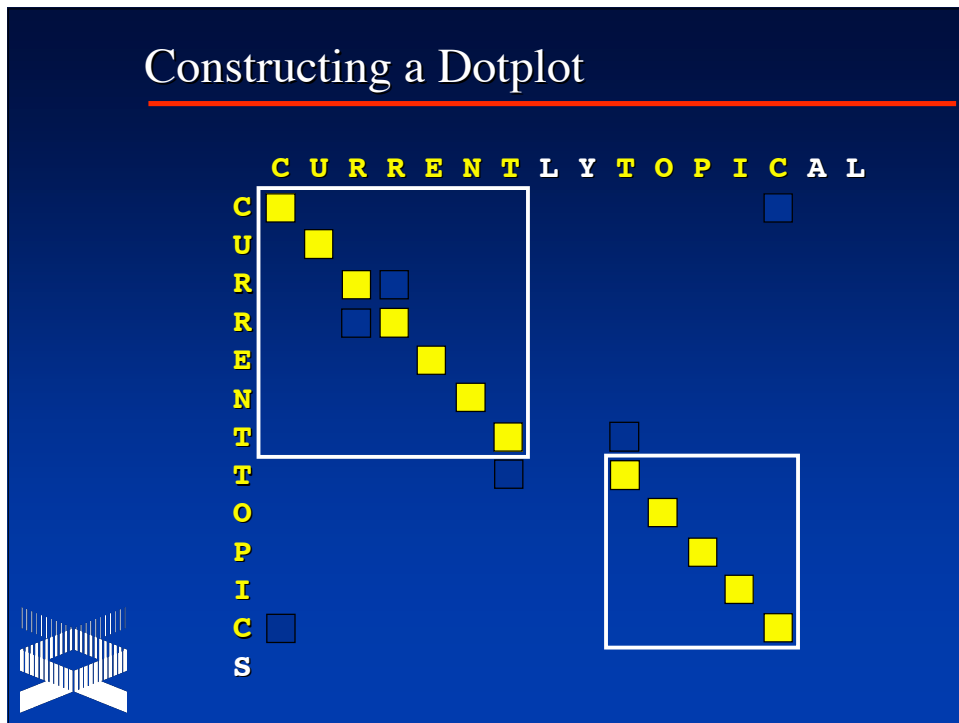
## Dotplots

- Visual method for comparing two sequences

- Allows for quick identification of
  - Regions of local alignment
  - Direct or inverted repeat regions
  - Insertions
  - Deletions
  - Low-complexity regions

- No statistical measure of the overall quality of the alignment

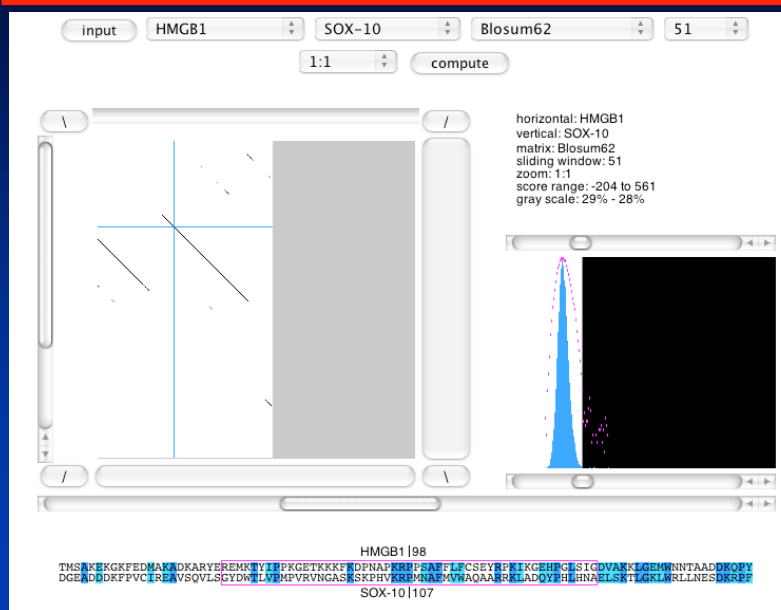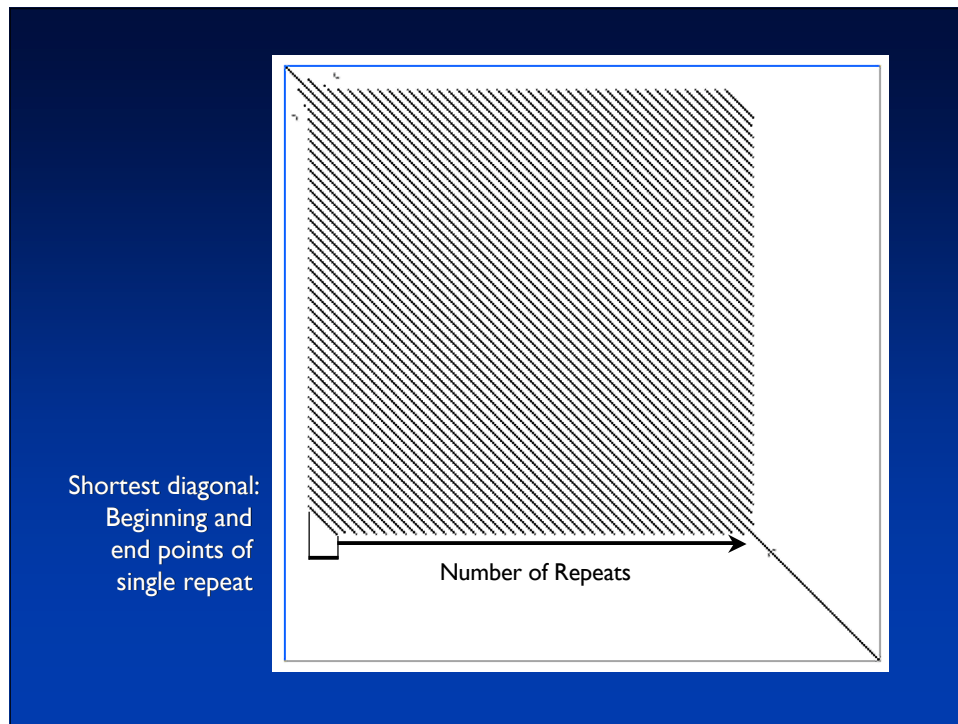## Constructing a Dotplot



## Constructing a Dotplot

## Tools for Constructing Dotplots

- Dotlet  (Java applet)
  *http://www.isrec.isb-sib.ch/java/dotlet/Dotlet.html*

- Dotter
  *http://www.cgr.ki.se/cgr/groups/sonhammer/Dotter.html*

- Dottup  (for complete genomes)
  *http://www.emboss.org*

- Dotplot subroutines also available through
  several software suites (GCG, DNA Strider)

## Finding Regions of Local Alignment

# Identifying Repeats

```
>gi|189599|gb|AAA60019.1| mucin
MTPGTQSPFFLLLLLTVLTVVTGSGHASSTPGGEKETSATQRSSVPSSTEKNAVSMTSSVLSSHSPGSGSSTTQGQDVTL
APATEPASGSAATWGQDVTSVPVTRPALGSTTPPAHDVTSAPDNKPAPGSTAPPAHGVTSAPDTRPAPGSTAPPAHGVTS
APDTRPAPGSTAPPAHGVTSAPDTRPAPGSTAPPAHGVTSAPDTRPAPGSTAPPAHGVTSAPDTRPAPGSTAPPAHGVTS
APDTRPAPGSTAPPAHGVTSAPDTRPAPGSTAPPAHGVTSAPDTR PAPGSTAPPAHGVTSAPDTR PAPGSTAPPAHGVTS
APDTRPAPGSTAPPAHGVTSAPDTRPAPGSTAPPAHGVTSAPDTRPAPGSTAPPAHGVTSAPDTRPAPGSTAPPAHGVTS
APDTRPAPGSTAPPAHGVTSAPDTRPAPGSTAPPAHGVTSAPDTRPAPGSTAPPAHGVTSAPDTRPAPGSTAPPAHGVTS
APDTRPAPGSTAPPAHGVTSAPDTRPAPGSTAPPAHGVTSAPDTRPAPGSTAPPAHGVTSAPDTRPAPGSTAPPAHGVTS
APDTRPAPGSTAPPAHGVTSAPDTRPAPGSTAPPAHGVTSAPDTRPAPGSTAPPAHGVTSAPDTRPAPGSTAPPAHGVTS
APDTRPAPGSTAPPAHGVTSAPDTRPAPGSTAPPAHGVTSAPDTRPAPGSTAPPAHGVTSAPDTRPAPGSTAPPAHGVTS
APDTRPAPGSTAPPAHGVTSAPDTRPAPGSTAPPAHGVTSAPDTRPAPGSTAPPAHGVTSAPDTRPAPGSTAPPAHGVTS
APDTRPAPGSTAPPAHGVTSAPDTRPAPGSTAPPAHGVTSAPDTRPAPGSTAPPAHGVTSAPDNRPALGSTAPPVHNVTS
ASGSASGSASTLVHNGTSARATTTPASKSTPFSJPSHHSDTPTTLASHSTKTDASSTHHSSVPPLTSSNHSTSPQLSTGV
SFFFLSFHISNLQFNSSLEDPSTDYYQELQRDJSEMFLQIYKQGGFLGLSNIKFRPGSVVVQJTLAFREGTINVHDVETQ
FNQYKTEAASRYNLTISDVSVSDVPFPFSAQGGAGVPGWGIALLVLVCVLVALAIVYLIALAVCQCRRKNYGQLDIFPAR
DTYHPMSEYPTYHTHGRYVPPSSTDRSPYEKVSAGNGGSSLSYTNPAVAAASANL
```

**PAPGSTAPPAHGVTSAPDTR**

40 tandem repeats of 20 amino acids



# Identifying Repeats

Shortest diagonal:
Beginning and
end points of
single repeat

Number of Repeats

# Identifying Low-Complexity Regions

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

- 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)

# Identifying Low-Complexity Regions

Example: *Drosophila* achaete-scute

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

*Homopolymeric
alanine-glutamine tract*

# Identifying Low-Complexity Regions

## 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 a particular theory of evolution
  - 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 |

## 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, 1978
- Point Accepted Mutation (PAM)
  - Look at patterns of substitutions in highly related proteins (> 85% similar), based on multiple sequence alignments
  - 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 are equally mutable
- Replacement is independent of surrounding residues
- Replacement is independent of previous mutations at the same position (Markov model)
- Sequences being compared are of average composition
- Forces responsible for sequence evolution over shorter time spans are the same as those for longer evolutionary time spans

## PAM Matrices: Sources of Error

- Small, globular proteins used to derive matrices (departure from average composition)
- Errors in PAM 1 are magnified up to PAM 250
- Does not account for conserved blocks or motifs

# BLOSUM Matrices

- Henikoff and Henikoff, 1992
- <u>Bloc</u>ks <u>Su</u>bstitution <u>M</u>atrix
  - Look only for differences in conserved, ungapped regions of a protein family ("blocks")
  - Directly calculated, using no extrapolations
  - More sensitive to 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
KKLEKEEEEGISQESSEEE
KKLEKEEEEGISQESSEEE
KKLEKEEEEGISQESSEEE
KPAQEETEETSSQESAEED
KKPAQETEETSSQESAEED
```

80% →

```
TGNQEEYGNTSSDSSDEDY

KKLEKEEEEGISQESSEEE
KKLEKEEEEGISQESSEEE
KKLEKEEEEGISQESSEEE

KPAQEETEETSSQESAEED
KKPAQETEETSSQESAEED
```

*A+T Hook Domain (Block IPB000637B)*

## 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% |
| PAM 120 | | > 50% |
| PAM 250 | Longer, weaker local alignments | > 30% |

BLOSUM *(Henikoff, 1993)*

| | | |
|---|---|---|
| BLOSUM 90 | Short alignments, highly similar | > 60% |
| BLOSUM 80 | | > 50% |
| BLOSUM 62 | Most effective in detecting known members of a protein family | > 35% |
| BLOSUM 30 | Longer, weaker local alignments | |

## 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

## 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: | GSQSLAALLNKCKT**PQG**QRLVNQWIKQPLMDKNRIEERLNLVEAFVED |

Neighborhood Words

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

$= 7 + 5 + 6$

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*

Cumulative Score

X

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* |
| $\lambda S$ | *normalized score* |
| $k$ | *minor constant* |

Cumulative Score

X

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

## 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

```
○○○        RID=1063758631-3905-654746.BLASTQ3, sp|P29617|PROS_DROME Protein prospero – Netscape        ○

   ◀  ▶  ⌂  ✕      🜋 http://www.ncbi.nlm.nih.gov/blast/Blast.cgi                      ▽  🔍 Search     

  🜋 RID=1063758631-3905-65474...                                                                   ✕

 ☐ >gi|6179901|gb|AAF05703.1|AF190403_1  🇱 homeodomain transcription factor Prospero [Drosophi
           Length = 1403

  Score =  880 bits (2273), Expect = 0.0                        ┌──────────────────────┐
  Identities = 493/627 (78%), Positives = 493/627 (78%)  ◀────  │ ≥ 25% for proteins   │
                                                                │ ≥ 70% for nucleotides│
 Query: 777   HVATAAPRPQMHHPAPARLPTRMGGAAGHTALKSELSEKFQMLRANNNSSMMRMSGTDLE 83
              HVATAAPRPQMHHPAPARLPTRMGGAAGHTALKSELSEKFQMLRANNNSSMMRMSGTDLE      ─ Gap
 Sbjct: 777   HVATAAPRPQMHHPAPARLPTRMGGAAGHTALKSELSEKFQMLRANNNSSMMRMSGTDLE 83  x Low-
                                                                        Complexity
 Query: 837   GLADVLKSEITTSLSALVDTIVTRFVHQRRLFSKQADSVTAAAEQLNKDLLLASQILDRK 89
              GLADVLKSEITTSLSALVDTIVTRFVHQRRLFSKQADSVTAAAEQLNKDLLLASQILDRK
 Sbjct: 837   GLADVLKSEITTSLSALVDTIVTRFVHQRRLFSKQADSVTAAAEQLNKDLLLASQILDRK 89

 Query: 897   SPRTKVADRPQNGPTPATQSAAAMFQAPKTPQGMNPVAAAALYNSMTGPFCLPPDXXXXX 956
              SPRTKVADRPQNGPTPATQSAAAMFQAPKTPQGMNPVAAAALYNSMTGPFCLPPD
 Sbjct: 897   SPRTKVADRPQNGPTPATQSAAAMFQAPKTPQGMNPVAAAALYNSMTGPFCLPPDQQQQQ 956

 Query: 957   XXXXXXXXXXXXXXXXXXXXXXXXXLEQNEALSLVVTPKKKRHKVTDTRITPRTVSRILAQDX 1016
                             LEQNEALSLVVTPKKKRVTDTRITPRTVSRILAQD
 Sbjct: 957   QTAQQQQSAQQQQQSSQQTQQQLEQNEALSLVVTPKKRHKVTDTRITPRTVSRILAQDG 1016

 Query: 1017  XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXASNGGNSNATPAQSPTRSSGGAAYHXXX 1076
                             ASNGGNSNATPAQSPTRSSGGAAYH
 Sbjct: 1017  VVPPTGGPPSTPQQQQQQQQQQQQQQQQQQQQQASNGGNSNATPAQSPTRSSGGAAYHPQP 1076

 Query: 1077  XXXXXXXXXVSLPTSVAIPNPSLHESKVFSPYSPFFNPXXXXXXXXXXXXXXXXXXXXXXXX 1136
                             VSLPTSVAIPNPSLHESKVFSPYSPFFNP
 Sbjct: 1077  PPPPPPMMPVSLPTSVAIPNPSLHESKVFSPYSPFFNPHAAAGQATAAQLHQHHQQHHPH 1136

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

 Query: 1197  DRQSECNSADMQFDGMAPTISFYKQMQLKTEHQESLMAKHCESLTPLHSSTLTPMHLRKA 1256    ◀▶
 ○ ✉ 🔒 ☎ ⊡ 🗋                                                                        ▭ 🖻 🔒
```

```
○○○        RID=1063758631-3905-654746.BLASTQ3, sp|P29617|PROS_DROME Protein prospero – Netscape        ○

   ◀  ▶  ⌂  ✕      🜋 http://www.ncbi.nlm.nih.gov/blast/Blast.cgi#6179901            ▽  🔍 Search     

  🜋 RID=1063758631-3905-65474...                                                                   ✕

  Score =  758 bits (1957), Expect = 0.0
  Identities = 454/704 (64%), Positives = 461/704 (65%)

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

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

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

 Query: 181  CSDRSLEAAAADVAGGSPPRAASVSSLNGGASSGEQHQSQLQHDLVAHHMLRNILQGKKE 240
             CSDRSLEAAAADVAGGSPPRAASVSSLNGGASSGEQHQSQLQHDLVAHHMLRNILQGKKE
 Sbjct: 181  CSDRSLEAAAADVAGGSPPRAASVSSLNGGASSGEQHQSQLQHDLVAHHMLRNILQGKKE 240

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

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

 Query: 361  XXXXXXXXXPTGQRSESRAPEEPQLPTKKESVDDMLDEVELLGLHSRGSDMDSLASPSHSX 420
                     PTGQRSESRAPEEPQLPTKKESVDDMLDEVELLGLHSRGSDMDSLASPSHS
 Sbjct: 361  AQDEEDAAPTGQRSESRAPEEPQLPTKKESVDDMLDEVELLGLHSRGSDMDSLASPSHSD 420

 Query: 421  XXXXXXXXXXXXXXXXXXCVEQKTSGSGCLKKPGMDLKRARVENIVSGMRCSPSSGLAQAG 480
                          CVEQKTSGSGCLKKPGMDLKRARVENIVSGMRCSPSSGLAQAG
 Sbjct: 421  MMLLDKDDVLDEDDDDDCVEQKTSGSGCLKKPGMDLKRARVENIVSGMRCSPSSGLAQAG 480   ◀▶
 ○ ✉ 🔒 ☎ ⊡ 🗋                                                                        ▭ 🖻 🔒
```

>gi|6179901|gb|AAF05703.1|AF190403_1  L homeodomain transcription factor Prospero
          Length = 1403

 Score =  880 bits (2273), Expect = 0.0 ✔
 Identities = 493/627 (78%) ✔ Positives = 493/627 (78%)

 Score =  758 bits (1957), Expect = 0.0 ✔
 Identities = 454/704 (64%) ✔ Positives = 461/704 (65%)

**HSP 2**          **HSP 1**

Color Key for Alignment Scores
<40  40-50  50-80  80-200  >=200

1_3905
0    250    500    750    1000    1250

# Suggested BLAST Cutoffs

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

**PICK THE RIGHT MATRIX AND
ALWAYS LOOK AT THE ALIGNMENTS!!!**

## Database Searching Artifacts

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

- Repetitive elements
  - LINE, SINE, Alu
  - Automatic masking "still under development"
  - RepeatMasker
    *http://repeatmasker.genome.washington.edu*

## Database Searching Artifacts

- "Hypothetical protein" hits
  - Some entries result from gene prediction or translation of transcripts
  - An ORF does not imply translation into a real protein

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

# BLAST2SEQUENCES

- 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/*

# MegaBLAST

- Optimized for aligning long and/or highly-similar sequences ("greedy algorithm")

- Good for batch nucleotide searches

- Search targets
  - Entire eukaryotic genomes
  - Trace Archives (125 million sequence traces)

- 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**

## Discontiguous MegaBLAST

- Designed specifically for the comparison of diverged sequences, particularly from different organisms

- Since these types of comparison may yield low degrees of identity, this variant performs better than the original MegaBLAST, which is optimized for sequences that are highly similar

# FASTA

- SSEARCH
  Smith-Waterman algorithm
  Rigorous and quite sensitive, but slow

- FASTA
  Regions of local alignment
  Approximation of Smith-Waterman algorithm
  Faster, but sacrifices sensitivity

- Bill Pearson, University of Virginia
  *http://fasta.bioch.virginia.edu*