## Overview

- Week 2
  - Similarity *vs*. Homology
  - Global *vs*. Local Alignments
  - Scoring Matrices
  - BLAST
  - BLAT

- Week 3
  - Profiles, Patterns, Motifs, and Domains
  - Structures: VAST, Cn3D, and *de novo* Prediction
  - Multiple Sequence Alignment

## 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
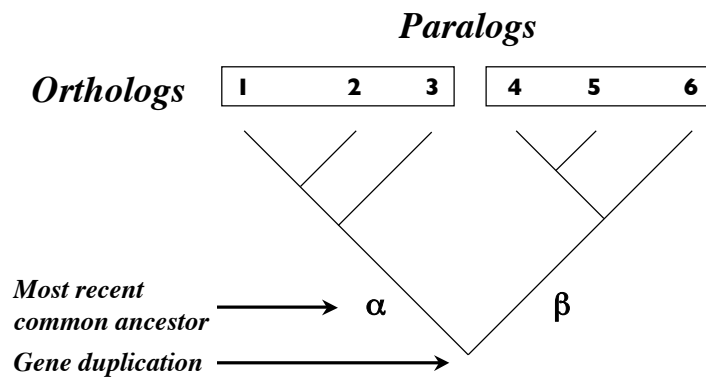    duplication (*paralogy*)

# 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

*Paralogs*

*Orthologs*   | 1     2     3 | 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)

## 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 alignment could be generated for any two sequences being compared
- Best for sequences that share some similarity, or for sequences of different lengths

## Scoring Matrices

- Empirical weighting scheme representing physicochemical and biological characteristics of nucleotides and amino acids
  - Side chain structure and chemistry
  - Side chain function

- Amino acid-based examples:
  - 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

- Why is understanding scoring matrices important?

  - Appear in all analyses involving sequence comparison

  - Implicitly represent particular evolutionary patterns

  - Choice of matrix can strongly influence outcomes of analyses

## 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:*

  Match          + 5
  Mismatch       – 4

- *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 -1 -1  0 -2 -1 -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  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

# BLOSUM Matrices

- Henikoff and Henikoff, 1992

- <u>Bloc</u>ks <u>S</u>ubstitution <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 detecting structural or functional substitutions

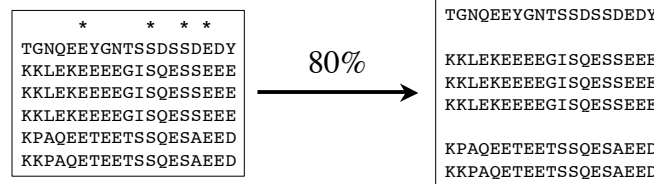  - Generally perform better than PAM matrices for local similarity searches *(Henikoff and Henikoff, 1993)*

# BLOSUM *n*

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

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

```
        *      *  * *
TGNQEEYGNTSSDSSDEDY
KKLEKEEEEGISQESSEEE
KKLEKEEEEGISQESSEEE          80%
KKLEKEEEEGISQESSEEE
KPAQEETEETSSQESAEED
KKPAQETEETSSQESAEED
```

```
TGNQEEYGNTSSDSSDEDY

KKLEKEEEEGISQESSEEE
KKLEKEEEEGISQESSEEE
KKLEKEEEEGISQESSEEE

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...

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

*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 |
|  | $n$ = | length of the gap |  |  |
| and | $G > L$ |  |  |  |

*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 with one another
  - 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          Neighborhood Score
PHG    13          Threshold
PMG    13          ($T = 13$)
PSG    13
PQA    12
PQN    12
etc.
```

# 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
```

**13**

## 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* |
| *λS* | *normalized score* |
| *k* | *minor constant* |

Cumulative Score

*X*

*S*

*T*

Extension

**14**

# 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

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
```

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

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

Cumulative Score

X

S

T

Extension

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



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

Available protein databases include:

| | |
|---|---|
| nr | Non-redundant |
| refseq | Reference Sequences |
| swissprot | SWISS-PROT |
| pat | Patents |
| pdb | Protein Data Bank |
| env_nr | Environmental samples |



Limit by organism or taxonomic group

# 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)
MESSAKMESGGAGQQPQPQPQQPFLPPAACFFATAAAAAAAAAAAAAAQSAQQQQQQQQQQQQAPQLRPAA
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 (but *not* enabled by default)

```
gb|AAB30541.1|  Prox 1=homeobox gene prospero homolog [mice, e...   219   1e-54
ref|XP_001375160.1|  PREDICTED: similar to FLJ36749 protein [Mono   216   1e-53
ref|XP_608175.3|  PREDICTED: similar to FLJ36749 protein [Bos tau   216   1e-53
ref|XP_234418.4|  PREDICTED: similar to RIKEN cDNA 1700058C01 ...   214   4e-53
ref|XP_001491332.1|  PREDICTED: similar to prospero homeobox 2 [E   213   1e-52
dbj|BAE06658.1|  transcription factor protein [Ciona intestinalis   210   6e-52
ref|XP_001517520.1|  PREDICTED: hypothetical protein [Ornithorhyn   210   7e-52
ref|XP_522907.2|  PREDICTED: hypothetical protein [Pan troglodyte   209   9e-52
ref|XP_001088672.1|  PREDICTED: similar to prospero-related ho...   209   1e-51
ref|NP_001073877.1|  prospero homeobox 2 [Homo sapiens] >ref|X...   207   4e-51
gb|AAI05928.1|  PROX2 protein [Homo sapiens] >gb|AAI05721.1|  P...   207   4e-51
ref|XP_692862.2|  PREDICTED: hypothetical protein [Danio rerio]    204   4e-50
emb|CAF92934.1|  unnamed protein product [Tetraodon nigroviridis]   201   2e-49
ref|NP_001071961.1|  transcription factor protein [Ciona intes...   199   1e-48
emb|CAG04605.1|  unnamed protein product [Tetraodon nigroviridis]   198   3e-48
emb|CAF95276.1|  unnamed protein product [Tetraodon nigroviridis]   196   1e-47
emb|CAG10630.1|  unnamed protein product [Tetraodon nigroviridis]   195   2e-47
gb|AAC28353.1|  Prox1 [Xenopus laevis]                             187   4e-45
emb|CAG09138.1|  unnamed protein product [Tetraodon nigroviridis]   175   2e-41
ref|XP_547908.2|  PREDICTED: similar to RIKEN cDNA 1700058C01 [Ca   168   2e-39
dbj|BAB17311.1|  Prox 1 [Cynops pyrrhogaster]                      162   3e-37
gb|EAW81198.1|  hCG22353 [Homo sapiens]                            158   2e-36
dbj|BAC04278.1|  unnamed protein product [Homo sapiens]            157   5e-36
gb|AAC59781.1|  prospero_like protein                              156   9e-36
gb|EDL02840.1|  RIKEN cDNA 1700058C01, isoform CRA_a [Mus musculu   154   4e-35
ref|XP_849216.1|  PREDICTED: similar to prospero-related homeo...   154   6e-35
ref|XP_001344006.1|  PREDICTED: similar to homeodomain protein [D   153   1e-34
emb|CAG09167.1|  unnamed protein product [Tetraodon nigroviridis]   150   6e-34
emb|CAG13403.1|  unnamed protein product [Tetraodon nigroviridis]   100   8e-19
gb|AAD30180.1|AC006530_2  homeobox prospero-like protein [Homo sa  97.4   6e-18
ref|XP_547411.2|  PREDICTED: similar to prospero-related homeo...  80.1   1e-12
pir||JC5496  Prox 1 protein 671 - chicken                         80.1   1e-12
ref|XP_001339994.1|  PREDICTED: similar to Prox 1 protein 671 ...  79.7   2e-12
ref|XP_970352.1|  PREDICTED: similar to Protein prospero, part...  57.4   9e-06
ref|NP_001100671.1|  prospero-related homeobox 1 [Rattus norve...  44.7   0.055
emb|CAF94749.1|  unnamed protein product [Tetraodon nigroviridis]  43.5   0.10
emb|CAP58279.1|  Prox1 protein [Xenopus tropicalis]               42.0   0.38
gb|AAF13029.1|AF070733_1  transcription factor Prox1 [Notophthalm  40.4   1.1
gb|ABG29070.1|  transcription factor Prox1 [Pleurodeles waltl]    38.9   3.3
ref|XP_663294.1|  hypothetical protein AN5690.2 [Aspergillus n...  37.7   7.2
```

Accept
(for now)

Reject



```
>ref|NP_731565.2|  prospero CG17228-PA, isoform A [Drosophila melanogaster]
 gb|AAN13501.2|  CG17228-PA, isoform A [Drosophila melanogaster]
Length=1535

 GENE ID: 41363 pros | prospero [Drosophila melanogaster]
(Over 100 PubMed links)

 Score =  938 bits (2425),  Expect = 0.0, Method: Compositional matrix ad
 Identities = 687/688 (99%), Positives = 687/688 (99%), Gaps = 0/688 (0%)

Query  17   LFQPQSVSTAnssssnnnnssTPAALATHsptsnspvsgassassllltaaFGNLFGGSSA   76
            LFQPQSVSTANSSSSNNNNSSTPAALATHSPTSNSPVSGASSASSLLTAAFGNLFGGSSA
Sbjct  17   LFQPQSVSTANSSSSNNNNSSTPAALATHSPTSNSPVSGASSASSLLTAAFGNLFGGSSA   76

Query  77   KMLNELFGRQMKQAQDATSGLPQSLDNAMLAAAMETATSAELLIGSLNSTSKLLQQQHNN   136
            KMLNELFGRQMKQAQDATSGLPQSLDNAMLAAAMETATSAELLIGSLNSTSKLLQQQHNN
Sbjct  77   KMLNELFGRQMKQAQDATSGLPQSLDNAMLAAAMETATSAELLIGSLNSTSKLLQQQHNN   136

Query  137  NSIAPANSTPMSNGTNasispgsahssshshqgvspKGSRRVSACSDRSLEAAAADVAGG   196
            NSIAPANSTPMSNGTNASISPGSAHSSSHSHQGVSPKGSRRVSACSDRSLEAAAADVAGG
Sbjct  137  NSIAPANSTPMSNGTNASISPGSAHSSSHSHQGVSPKGSRRVSACSDRSLEAAAADVAGG   196

Query  197  SPPRAASVSSLNGGASSGEQHQSQLQHDLVAHHMLRNILQGKKELMQLDQELRTAMqqqq   256
            SPPRAASVSSLNGGASSGEQHQSQLQHDLVAHHMLRNILQGKKELMQLDQELRTAMQQQQ
Sbjct  197  SPPRAASVSSLNGGASSGEQHQSQLQHDLVAHHMLRNILQGKKELMQLDQELRTAMQQQQ   256

Query  257  qqlqekeqlHSKLnnnnnnniaatannnnnttMESINLIDDSEMADIKIKSEPQTAPQPQ   316
            QQLQEKEQLHSKLNNNNNNNIAATANNNNNTTMESINLIDDSEMADIKIKSEPQTAPQPQ
Sbjct  257  QQLQEKEQLHSKLNNNNNNNIAATANNNNNTTMESINLIDDSEMADIKIKSEPQTAPQPQ   316

Query  317  QsphgsshssrsgsgsgshssmasdgslrrkssdsldsHGaqddaqdeedaaPTGQRSES   376
            QSPHGSSHSSRSGSGSGSHSSMASDGSLRRKSSDSLDSHGAQDDAQDEEDAAPTGQRSES
Sbjct  317  QSPHGSSHSSRSGSGSGSHSSMASDGSLRRKSSDSLDSHGAQDDAQDEEDAAPTGQRSES   376

Query  377  RAPEEPQLPTKKESVDDMLDEVELLGLHSRGSDMDSLASPSQSdmmlldkddvldedddd   436
            RAPEEPQLPTKKESVDDMLDEVELLGLHSRGSDMDSLASPS SDMMLLDKDDVLDEDDDD
Sbjct  377  RAPEEPQLPTKKESVDDMLDEVELLGLHSRGSDMDSLASPSHSDMMLLDKDDVLDEDDDD   436

Query  437  dCVEQKTSGSGCLKKPGMDLKRARVENIVSGMRCSPSSGLAQAGQLQVNGCKKRKLYQPQ   496
            DCVEQKTSGSGCLKKPGMDLKRARVENIVSGMRCSPSSGLAQAGQLQVNGCKKRKLYQPQ
Sbjct  437  DCVEQKTSGSGCLKKPGMDLKRARVENIVSGMRCSPSSGLAQAGQLQVNGCKKRKLYQPQ   496
```

≥ 25% for proteins
≥ 70% for nucleotides

— Gap
a Low-
  Complexity

```
                                                    NCBI Blast:Query sequence
            http://www.ncbi.nlm.nih.gov/blast/Blast.cgi#16768018                              Google
   Protein BLAST: search protein ...        NCBI Blast:Query sequence

Query   617    NHKEETGQERpgssspspsplkpktslgESSDSGANMLSQMMSKMMSGKLHNPLVGVGHP    676
               NHKEETGQERPGSSSPSPSPLKPKTSLGESSDSGANMLSQMMSKMMSGKLHNPLVGVGHP
Sbjct   617    NHKEETGQERPGSSSPSPSPLKPKTSLGESSDSGANMLSQMMSKMMSGKLHNPLVGVGHP    676

Query   677    ALPQGFPPLLQHMGDMSHAAAMYQQFFF    704
               ALPQGFPPLLQHMGDMSHAAAMYQQFFF
Sbjct   677    ALPQGFPPLLQHMGDMSHAAAMYQQFFF    704

  Score =  635 bits (1639),  Expect = 6e-180, Method: Compositional matrix adjust.
  Identities = 461/498 (92%), Positives = 463/498 (92%), Gaps = 32/498 (6%)

Query   906    PQNGPTPATQSAAAMFQAPKTPQGMNPVAAAALYNSMTGPFCLPPDqqqqqqtaqqqqsa    965
               P   P+P    +AAAMFQAPKTPQGMNPVAAAALYNSMTGPFCLPPDQQQQQQTAQQQQSA
Sbjct   1070   PHIRPSP---TAAAMFQAPKTPQGMNPVAAAALYNSMTGPFCLPPDQQQQQQTAQQQQSA    1126

Query   966    qqqqqssqqtqqqLEQNEALSLVVTPKKKRHKVTDTRITPRTVSRILAQDgvvpptggpp    1025
               QQQQQSSQQTQQQLEQNEALSLVVTPKKKRHKVTDTRITPRTVSRILAQDGVVPPTGGPP
Sbjct   1127   QQQQQSSQQTQQQLEQNEALSLVVTPKKKRHKVTDTRITPRTVSRILAQDGVVPPTGGPP    11

Query   1026   stpqqqqqqqqqqqqqqqqqqqqqASNGGNSNATPAQSPTRSSGGAAYHpqppppppppmmp    10
               STPQQQQQQQQQQQQQQQQQQQQQASNGGNSNATPAQSPTRSSGGAAYHPQPPPPPPPMMP
Sbjct   1187   STPQQQQQQQQQQQQQQQQQQQQQASNGGNSNATPAQSPTRSSGGAAYHPQPPPPPPPMMP    12

Query   1086   VSLPTSVAIPNPSLHESKVFSPYSPFFNPhaaagqataaqlhqhhqqhhphhqsmqlsss    1
               VSLPTSVAIPNPSLHESKVFSPYSPFFNPHAAAGQATAAQLHQHHQQHHPHHQSMQLSSS
Sbjct   1247   VSLPTSVAIPNPSLHESKVFSPYSPFFNPHAAAGQATAAQLHQHHQQHHPHHQSMQLSSS    1306

Query   1146   ppgslgALMDSRDspplphppsmlhpallaaahhggspDYKTCLRAVMDAQRQSECNSA    1205
               PPGSLGALMDSRDSPPLPHPPSMLHPALLAAAHHGGSPDYKTCLRAVMDAQDRQSECNSA
Sbjct   1307   PPGSLGALMDSRDSPPLPHPPSMLHPALLAAAHHGGSPDYKTCLRAVMDAQDRQSECNSA    1366

Query   1206   DMQFDGMAPTISFYKQMQLKTEHQESLMAKHCESLTPLLGSTLTPMHLRKAKLMFFWVRY    1265
               DMQFDGMAPT                     SSTLTPMHLRKAKLMFFWVRY
Sbjct   1367   DMQFDGMAPT----------------------------SSTLTPMHLRKAKLMFFWVRY    1397

Query   1266   PSSAVLKMYFPDIKFNKNNTAQLVKWFSNFREFYYIQMEKYARQAVTEGIKTPDDLLIAG    1325
               PSSAVLKMYFPDIKFNKNNTAQLVKWFSNFREFYYIQMEKYARQAVTEGIKTPDDLLIAG
Sbjct   1398   PSSAVLKMYFPDIKFNKNNTAQLVKWFSNFREFYYIQMEKYARQAVTEGIKTPDDLLIAG    1457

Query   1326   DSELYRVLNLHYNRNNHIEVPQNFRFVVESTLREFFRAIQGGKDTEQSWKKSIYKIISRM    1385
Done
```

No definition line ∴ second HSP identified

— Gap
a Low-Complexity

```
>┌ref|NP_731565.2| UG prospero CG17228-PA, isoform A [Drosophila melanogaster]
 gb|AAN13501.2| G CG17228-PA, isoform A [Drosophila melanogaster]
Length=1535

 Score =  938 bits (2425),  Expect = 0.0 Method: Compositional matrix adjust.
 Identities = 687/688 (99%), Positives = 687/688 (99%), Gaps = 0/688 (0%)


 Score =  635 bits (1639),  Expect = 6e-180 Method: Compositional matrix adjust.
 Identities = 461/498 (92%), Positives = 463/498 (92%), Gaps = 32/498 (6%)
```

HSP 1
Q: 17–704
S: 17–704

HSP 2
Q:  906–1403
S: 1070–1535

Color key for alignment scores
<40   40-50   50-60   80-200   >=200
Query
0     250    500    750   1000   1250

# Suggested BLAST Cutoffs

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

- *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

- Repetitive elements
    - LINEs, SINEs, retroviral repeats
    - Choose "Filter: Species-Specific Repeats" when using BLASTN
    - RepeatMasker
      *http://www.repeatmasker.org*

- 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



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

# 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**



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

# 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



http://genome.ucsc.edu

# 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*