# Lecture 7: Sequencing a Cloned Gene – Analysis and Annotation

# Universal Genetic Code

| First Letter | Second Letter | | | | | | | | Third Letter |
|---|---|---|---|---|---|---|---|---|---|
| | **U** | | **C** | | **A** | | **G** | | |
| U | UUU<br>UUC | Phenylalanine | UCU<br>UCC | Serine | UAU<br>UAC | Tyrosine | UGU<br>UGC | Cysteine | U<br>C |
| | UUA<br>UUG | Leucine | UCA<br>UCG | | UAA<br>UAG | Stop<br>Stop | UGA<br>UGG | Stop<br>Tryptophan | A<br>G |
| C | CUU<br>CUC | Leucine | CCU<br>CCC | Proline | CAU<br>CAC | Histidine | CGU<br>CGC | Arginine | U<br>C |
| | CUA<br>CUG | | CCA<br>CCG | | CAA<br>CAG | Glutamine | CGA<br>CGG | | A<br>G |
| A | AUU<br>AUC<br>AUA | Isoleucine | ACU<br>ACC | Threonine | AAU<br>AAC | Asparagine | AGU<br>AGC | Serine | U<br>C<br>A |
| | AUG | Methionine; Start | ACA<br>ACG | | AAA<br>AAG | Lysine | AGA<br>AGG | Arginine | A<br>G |
| G | GUU<br>GUC | Valine | GCU<br>GCC | Alanine | GAU<br>GAC | Aspartate | GGU<br>GGC | Glycine | U<br>C |
| | GUA<br>GUG | | GCA<br>GCG | | GAA<br>GAG | Glutamate | GGA<br>GGG | | A<br>G |

# Test Sequences

**_btr-3_  Nucleotide Sequence**

**ATG**AAAGTGGAGAGTTGGTTGCACTTGGGTTGGTTGCTGGGGTTGCTGCTGGTCCTGTTGCCGTTGGTC
CGATGCCAAGGATGGGGCGAACCACGGTTCGAGACGGGAAATGTGGAAAATATATCACTCGCCGCATAC
AACGAGGCGCAGCTGCAGCAAGATGTCTGGATGGTGGAGGAGATGGATGCACCGTTCGTGCTGCTCTAC
ATCAATTACCAAGGACCGTCCGAGCCTACGATACGCGAGTCACCGGCCGATCTTGACGCAAGGCTACAGC
TGTCCGAGGCTGGCCGCTGGTCGATCGTAATCAATCGCCGTCAGGACTACGAGGTGCATCAGCGTAGCA
GTCTCATTCTGCTGGCCGTCGAATCCACGGCTATCCCGTACGCGATCGTGGTCAACTTGGTGAACGTGCTG
GACAATGCGCCCGTCATGACGGCCCAAGGTAGCTGTGAGATTGAGGAGTTGCGCGGGGACTTTGTGAC
GGACTGTCTGTTTAACGTGTACCATGCGGACGGGTTCGAGGAGAATGGCATTGGCAATTCGAGCACGAA
CGAGCTGTCGTTCGAGATCGGTGATGTGGCCGGTGCGCGGGACCACTTTACGTACGTGCCCTCCACGGT
GACCCCTTCCCAGCCGATCTACAACAAGCTGTTCAATTTGAAAGTTTTAAAGCAGCTGGACTACACCGAG
AACGCTATATTTAACTTCATCACCACCGTGTACGACCTAGACCGGACGCACTCCTTCAAGATGAGTACGAT
CGTTCAGGTNCGCAACGTCGATAGCCGGCCTCCGATCTTTAGCCGACCGTTCNCCAGCGAACGNATCATG
NAAANGgAANCATTTTACGCgANCGTGATCGCANtCGAcCGTGACACTgGaCTAAACAAACCGATCTGTT
ACGAGCTGACGGCTCTAGTACCGGAATATCAGAAATATTTCGATATTGGACAAACTGATGGAAAGCTGAC
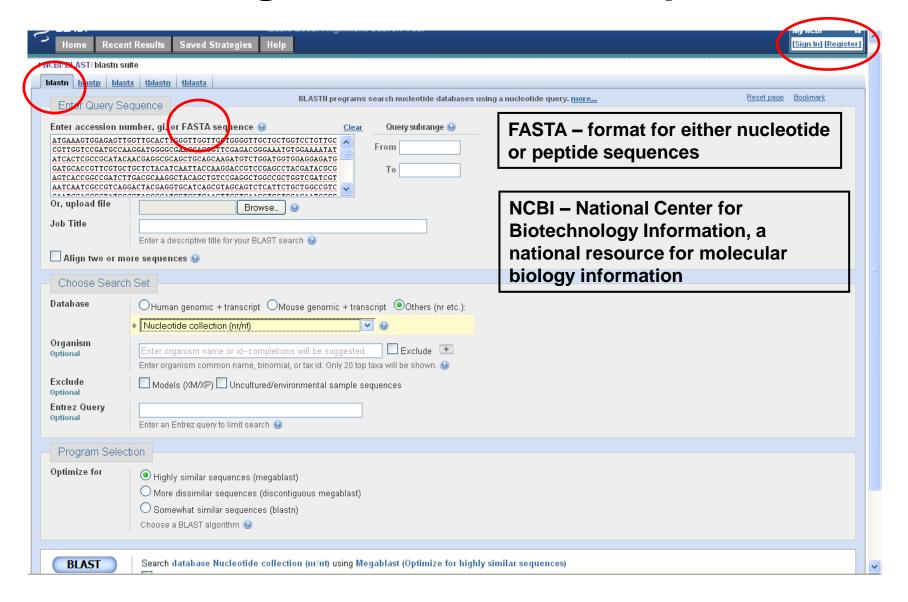CGTGCACCCGATTGATCGAGATGCGG

**BTR-3 Predicted  Protein  Sequence**

MKVESWLHLGWLLGLLLVLLPLVRCQGWGEPRFETGNVENISLAAYNEAQLQQDVWMVEEMDAPFVLLYI
NYQGPSEPTIRESPADLDARLQLSEAGRWSIVINRRQDYEVHQRSSLILLAVESTAIPYAIVVNLVNVLDNAPV
MTAQGSCEIEELRGDFVTDCLFNVYHADGFEENGIGNSSTNELSFEIGDVAGARDHFTYVPSTVTPSQPIYN
KLFNLKVLKQLDYTENAIFNFITTVYDLDRTHSFKMSTIVQVRNVDSRPPIFSRPFXSERIMXXEXFYAXVIAX
DRDTGLNKPICYELTALVPEYQKYFDIGQTDGKLTVHPIDRDA

# National Center for Biotechnology Information (NCBI) BLAST Family of Programs

- **Blastp – compares an amino acid query sequence against a protein sequence database; recognizes evolutionary conservation**

- **Blastn – compares a nucleotide query sequence against a nucleotide sequence database**

- **Blastx – compares a nucleotide query sequence translated in all reading frames against a protein sequence database**

- **Tblastn – compares a protein query sequence against a nucleotide sequence database dynamically translated in all reading frames**

- **Tblastx - compares the six-frame translations of a nucleotide query sequence against the six-frame translations of a nucleotide sequence database**

- **BLAST Line (BLink) –  A link option on protein records that displays the results of a pre-computed BLAST search of that protein against all other protein sequences at NCBI**

# Entering the Nucleotide Sequence



FASTA – format for either nucleotide or peptide sequences

NCBI – National Center for Biotechnology Information, a national resource for molecular biology information

# Important Terms Used in BLAST

- Accession number: unique identification number for the sequence

- Score: score calculated for this particular match, using a scoring matrix

- Query coverage: percentage of the query sequence that matched the sequence in the database. (In this case, nucleotide 280 to nucleotide 1000 = 72% of the query sequence)

- E value (Expect value): the number of hits one can "expect" to see by chance when searching a database of a particular size. (In the following example, 0% implies that it is a significant match)

- Maximum identity: percentage of nucleotides matched being identical (In the following example, 706 of the 727 nucleotides matched were identical to each other, which is 97%)

# Results Obtained

# Sequence Comparison

# Entering the Protein Sequence



Paste protein sequence

# Sequence Identification Numbers

- **There are two types of sequence identification numbers, GI and VERSION, each of which have different formats and were implemented at different times.**

- **GI number (sometimes written in lower case, "gi") is simply a series of digits that are assigned consecutively to each sequence record processed by NCBI.**

- **The GI number bears no resemblance to the Accession number of the sequence record.**

- **Nucleotide sequence GI number is shown in the VERSION field of the database record.**

- **Protein sequence GI number is shown in the CDS/db_xref field of a nucleotide database record, and the VERSION field of a protein database record.**

# Sequence Identification Numbers

- **VERSION is made of the accession number of the database record followed by a dot and a version number (and is, therefore, sometimes referred to as the "accession.version").**

- **Nucleotide sequence version contains two letters followed by six digits, a dot, and a version number (or for older nucleotide sequence records, the format is one letter followed by five digits, a dot, and a version number)**

- **Protein sequence version contains three letters followed by five digits, a dot, and a version number.**

- **The GI number has been used for many years by NCBI to track sequence histories in GenBank and the other sequence databases it maintains.**

- **The VERSION system of identifiers was adopted in February, 1999, by the International Nucleotide Sequence Database Collaboration (GenBank, the NIH genetic sequence database, an annotated collection of all publicly available DNA sequences; the European Molecular Biology Laboratory, EMBL; the DNA Databank of Japan, DDBJ).**

- **The two systems of identifiers run in parallel to each other, i.e., when any change is made to a sequence, it receives a new GI number AND an increase to its version number.**

# Results Obtained



- **Red bars denote alignment score ≥200**

- **Pink bars denote alignment score of 80-200**

**Protein domain found in the sequence**

# Results Obtained – Proteins in the Database that Matched the Query Sequence

```
                                                                  Score      E
Sequences producing significant alignments:                      (Bits)   Value

ref|XP_312086.4|  AGAP002828-PA [Anopheles gambiae str. PEST] ...   471    4e-131   U G
ref|XP_001864657.1|  conserved hypothetical protein [Culex qui...   211    7e-53      G
ref|XP_001652804.1|  hypothetical protein AaeL_AAEL007478 [Aed...   203    3e-50      G
gb|AAM21151.1|  cadherin [Manduca sexta]                            106    3e-21
dbj|BAA99406.1|  cadherin-like membrane protein [Bombyx mori]       105    5e-21      G
dbj|BAA99405.1|  cadherin-like membrane protein [Bombyx mori]       105    5e-21      G
ref|NP_001037682.1|  cadherin-like membrane protein [Bombyx mo...   105    5e-21    U G
dbj|BAA99404.1|  cadherin-like membrane protein [Bombyx mori]       105    5e-21      G
gb|ABI55354.1|  cadherin-like protein [Helicoverpa armigera] >...   105    7e-21
gb|ACY69034.1|  mutant cadherin [Helicoverpa armigera]              105    9e-21
gb|ACY69035.1|  mutant cadherin [Helicoverpa armigera]              104    1e-20
gb|ABF69363.1|  truncated cadherin-like protein [Helicoverpa a...   104    2e-20
gb|ABI55346.1|  cadherin-like protein [Helicoverpa armigera]        103    2e-20
gb|ABI55350.1|  cadherin-like protein [Helicoverpa armigera]        103    2e-20
gb|ACF94775.1|  cadherin protein [Helicoverpa armigera]             103    3e-20
gb|ABS90362.1|  truncated cadherin [Helicoverpa armigera] >gb|...   103    3e-20
gb|ABF69362.1|  cadherin-like protein [Helicoverpa armigera]        103    3e-20
gb|ACY69033.1|  mutant cadherin [Helicoverpa armigera]              103    4e-20
gb|AAG37912.1|AF319973_2  cadherin-related protein receptor BT...   103    4e-20
gb|ACY69032.1|  mutant cadherin [Helicoverpa armigera]              102    4e-20
gb|ACK37450.1|  cadherin Il [Ostrinia nubilalis]                    102    4e-20
gb|AAT67416.1|  cadherin-like protein [Helicoverpa armigera]        102    5e-20
gb|ABS59299.1|  cadherin-like protein [Ostrinia furnacalis]         102    5e-20
gb|AAU50667.1|  E-cadherin [Helicoverpa armigera]                   102    5e-20
gb|ABI55359.1|  cadherin-like protein [Helicoverpa armigera]        102    5e-20
gb|AAT67417.1|  cadherin-like protein [Helicoverpa armigera]        102    5e-20
gb|ABU41413.1|  cadherin-like protein [Plutella xylostella]         102    5e-20
gb|AAU50666.1|  E-cadherin [Helicoverpa armigera] >gb|ABI55349...   102    5e-20
gb|ACZ06063.1|  cadherin [Helicoverpa armigera]                     102    5e-20
gb|ABI55355.1|  cadherin-like protein [Helicoverpa armigera]        102    5e-20
gb|AAT37678.1|  cadherin Al [Ostrinia nubilalis]                    102    5e-20
gb|ABI55348.1|  cadherin-like protein [Helicoverpa armigera]        102    6e-20
gb|ABI55358.1|  cadherin-like protein [Helicoverpa armigera]        102    6e-20
```

**Unigene (U) and Entrez Gene (G) links available.**

# *UniGene*

- *UniGene* provides an organized view of the transcriptome.
- Each *UniGene* entry is a set of transcript sequences that appear to come from the same transcription locus (gene or expressed pseudogene), together with information on protein similarities, gene expression, cDNA clone reagents, and genomic location.
- In addition to sequences of well-characterized genes, hundreds of thousands of novel expressed sequence tag (EST) sequences are available.
- *UniGene* is a resource for gene discovery.
- *UniGene* has also been used by experimentalists to select reagents for gene mapping projects and large-scale expression analysis.

# Entrez Gene

- **Entrez Gene (http://www.ncbi.nim.nih.gov/gene) is the NCBI's database for gene-specific information.**

- **Entrez Gene maintains records from genomes which have been completely sequenced, which have an active research community to submit gene-specific information, or which are scheduled for intense sequence analysis.**

- **The content represents the integration of curation and automated processing from NCBI's Reference Sequence project (RefSeq), collaborating model organism databases, consortia such as Gene Ontology and other databases within NCBI.**

- **Records in Entrez Gene are assigned unique, stable and tracked integers as identifiers.**

- **The content (nomenclature, genomic location, gene products and their attributes, markers, phenotypes and links to citations, sequences, variation details, maps, expression, homologs, protein domains and external databases) is available via interactive browsing through NCBI's Entrez system, via NCBI's Entrez programming utilities (E-Utilities) and for bulk transfer by File Transfer Protocol (FTP).**

# Sequence Comparisons

# CLUSTAL

- **CLUSTAL is a widely used  multiple sequence alignment computer program.**

- **There are two main variations:**
  - ✓**ClustalW: command line interface**
  - ✓**ClustalX: has a graphical user interface**

# CLUSTALW

- The basic information that multiple alignments of protein sequences provide is **identification** of **conserved** sequence regions.

- This is very useful in designing experiments to **test and modify** the **function** of specific proteins, in **predicting** the **function and structure** of proteins and in identifying **new** members of protein families.

- In ClustalW, a **pairwise score** is calculated for every pair of sequences that are to be aligned.

- Pairwise scores are calculated as the number of **identities** in the best alignment divided by the **number of residues** compared (gap positions are excluded).

- As the pairwise score is calculated independently of the matrix and gaps chosen, it will always be the **same value** for a particular pair of sequences.

- Alignment score is calculated in **two** ways – **fast** and **slow** (more accurate mode).

- The scores are calculated from **separate** pairwise alignments.

- The scores can be calculated using **two** methods: **dynamic programming** (slow but accurate) or by the method of **Wilbur and Lipman** (extremely fast but approximate).

- The Wilbur-Lipman Method constructs tables of **prime K-tuples** to find regions of **similarity** between **two or more** DNA sequence pairs.

- Prime *k*-tuple is a **finite** collection of values representing a **repeatable pattern of differences** between **prime numbers**.

# CLUSTALW



**Paste all sequences in FASTA format. Three sequences were used:**

1. **Anopheles (derived)**
2. **Anopheles (from database)**
3. **Culex (from database).**

**FASTA format is a text-based format for representing either nucleotide sequences or peptide sequences, in which nucleotides or amino acids are represented using single-letter codes.**

# Results Obtained

## CLUSTALW Result

[clustalw.aln][clustalw.dnd][readme]

```
CLUSTAL W (1.81) Multiple Sequence Alignments


Sequence type explicitly set to Protein
Sequence format is Pearson
Sequence 1: 1                  333 aa
Sequence 2: 2                 1561 aa
Sequence 3: 3                 1129 aa
Start of Pairwise alignments
Aligning...


Sequences (1:2) Aligned. Score: 68.1682
Sequences (1:3) Aligned. Score: 35.7357
Sequences (2:2) Aligned. Score: 100
Sequences (2:3) Aligned. Score: 32.5952
Sequences (3:2) Aligned. Score: 32.5952
Sequences (3:3) Aligned. Score: 100
Guide tree      file created:   [clustalw.dnd]
Start of Multiple Alignment
There are 2 groups
Aligning...
Group 1: Sequences:    2      Score:4201
Group 2: Sequences:    3      Score:4689
Alignment Score 4539
CLUSTAL-Alignment file created  [clustalw.aln]
```

**Pairwise scores calculated for each pair of sequences**

**Alignment score calculated from pairwise scores**

# Sequence Comparisons

**Three sequences compared**

```
1   ------------------------------------------------------------
2   ------------------------------------------------------------
3   MSPPLMLLLITTSTTLTGAHLSRIQYNVCPKWLMMCANVEWKNVAVASITNCQHHRARCA

1   ---------------------------------MKVESWLHLGWLLGLLLVLLPLVR
2   ------------------------------------------------------------
3   VPKFPSRATKNNCQLLPVGSHRIHPDQITLVGGNGHADKIIRKQTGWLLLLLLPSLVFAQ

1   CQG--WGEPRFETGNVENISLAAYNEAQLQQDVWMVEEMDAPFVLLYINYQGPSEPTIRE
2   ------------------------------------------------------------
3   DPPGTWQQPYAIPVDAEKVSFLGYDSLSSELRVSMWEEMVVPFKLVELNYHGP-EADIKI

1   SPADLDARLQLSEAGRWSIVINRRQDYEVHQRSSLILLAVESTAIPYAIVVNLVNVLDNA
2   -----------SEGGRWSIVINRRQDYEVHQRSSLILLAVESTAIPYAIVVNLVNVLDNA
3   TNSGQTGAVLHLEGGKHFIVINNKMDYEVAAHRTSMVYLSVGNSQ-IFLAIDLINILDNV
             *.*:   ****.: ****  : : ::      ..:     :.::*:*:***.

1   PVMTAQGSCEIEELRGDFVTDCLFNVYHADGFEENGIGNSSTNELSFEIGDVAGARDHFT
2   PVMTAQGSCEIEELRGDFVTDCLFNVYHADGFEENGIGNSSTNELSFEIGDVAGARDHFT
3   PVMSSAGPCSVDEGLENYLSNCEYTVFHADGFVTNGILGNDTNAVGFDLPETNAELFKFE
    ***:: *.*.::*   :::::* :.*:*****  *** ...** :.*:: :. .   :*

1   YVPSTVTPSQPIYNKLFNLKVLKQLDYTENAIFNFITTVYDLDRTHSFKMSTIVQVRNVD
2   YVPSTVTPSQPIYNKLFNLKVLKQLDYTENAIFNFITTVYDLDRTHSFKMSTIVQVRNVD
3   EVVSGGDN----YNKKFKLKVLKKLDYTQNAVYSFLVTVYDLNRTHTATQNIVVQVINVE
     * *       *** *:*****:****:**::.*:.*****:***: . .  :*** **:

1   SRPPIFSRPFXSERIMXXEXFYAXVIAXDRDTGLNKPICYELTALVPEY--QKYFDIGQT
2   SRPPIFSRPFTSERIMEKEPFYATVIAIDRDTGLNKPICYELTALVPECKQAKYFEIGQT
3   SRDPVFTRPFTTQRIDEKSPYSTIVQAIDGDTGLGRPICYEIVTEQEKY--AEYFSIGRE
    ** *:*:*** ::;**   . : : * * * ****.:*****:.:    :    :**.**:

1   DGKLTVHPIDRDA-----------------------------------------------
2   DGKLTVHPIDRDAEQNELYTFTIVAYKCHNRLLNTSSEGAIILLDKNDNIPEIYMKPLEL
3   TGELNVKPINRDHEQNEFYQFTIWAYKCHNREFNESNVGAIILNDLNDSPPVFSVEPTQL
     *:*.*:*:**:**
```

**Symbols Used in CLUSTALW:**

- ✓ **\* Indicates identical amino acid residues in all sequences (or identical bases if DNA sequences are aligned)**

- ✓ **: indicates different but highly conserved (very similar) amino acids**

- ✓ **. Indicates different amino acids that are somewhat similar**

- ✓ **Blank indicates dissimilar amino acids or gaps (or different bases if DNA sequences are aligned)**

# Dayhoff Point Accepted Mutation (PAM)

- Point Accepted Mutation (PAM) is a **set of matrices** used to score sequence alignments.

- The PAM matrices are based on **1572** observed mutations in **71** families of closely related proteins.

- Each matrix is **twenty-by-twenty** (for the twenty standard amino acids).

- The value in a given cell represents the **probability** of a **substitution** of one amino acid for another.

- This type of matrix is commonly known as a **substitution matrix**.

- The PAM matrices imply a **Markov** chain model of protein mutation.

- A Markov model is a stochastic model used to model randomly changing systems where it is assumed that future states depend only on the current state, not on the events that occurred before it

- The PAM matrices are **normalized** so that the **PAM1 matrix** has **one point mutation per hundred amino acids**, and is appropriate for scoring sequences which are very similar.

- PAM matrices for comparing sequences of **lower** similarity are calculated from **repeated multiplication** of the PAM1 matrix by itself.

- **PAM250** is equivalent to **250 substitutions per hundred amino acids**.

# BLOck SUbstitution Matrix (BLOSUM)

- PAM compares **closely related** species and **does not** work very well for aligning **evolutionarily divergent** sequences.

- Sequence changes **over long** evolutionary time-scales are **not** approximated very efficiently by compounding **small** changes that occur over **short** time-scales.

- The BLOSUM series of matrices uses **multiple alignments** of evolutionarily **divergent** proteins and  the probabilities used in the matrix calculation are computed by looking at "**blocks**" of conserved sequences found in **multiple** protein alignments.

- These conserved sequences are assumed to be of functional importance within **related** proteins.

- To **reduce bias** from closely related sequences, segments in a block with a sequence identity **above** a certain threshold are **clustered** giving weight value of **1** to each such cluster.

- For the **BLOSUM45** matrix, this threshold is set at **45**%.

- Pairs frequencies are counted between clusters and, therefore, pairs are counted **only** between segments **less** than 45% identical.

- A  higher numbered BLOSUM matrix is used for aligning two closely related sequences and a lower number for more divergent sequences.

- **BLOSUM62** matrix does an excellent job detecting similarities in **distant** sequences, and is the **default**  matrix used most for **alignment applications** such as BLAST.

# Differences Between PAM and BLOSUM

- PAM matrices are based on an **explicit evolutionary model** (i.e., replacements are counted on the branches of a phylogenetic tree) whereas the BLOSUM matrices are based on an **implicit** model of evolution.

- The PAM matrices are based on **mutations** observed throughout a **global** alignment and includes both **highly conserved** and **highly mutable** regions.

- The BLOSUM matrices are based **only** on highly conserved regions in series of alignments forbidden to contain gaps.

- The method used to count the replacements is different: **unlike** the PAM matrix, the BLOSUM procedure uses **groups** of sequences within which **not all** mutations are counted the **same**.

- **Higher** numbers in the **PAM** matrix-naming scheme denote **larger evolutionary distance** whereas **higher** numbers in the **BLOSUM** matrix-naming scheme denote **higher sequence similarity** and, therefore, **smaller evolutionary distance**.

- Example: PAM150 is used for more distant sequences than PAM100; BLOSUM62 is used for closer sequences than Blosum50.