# Genome Annotation

**Bioinformatics 301**
**David Wishart**
*david.wishart@ualberta.ca*
*Notes at: http://wishartlab.com*

# Objectives*

- **To demonstrate the growing importance of gene and genome annotation in biology and the role bioinformatics plays**

- **To make students aware of new trends in gene and genome annotation (i.e. "deep" annotation)**

- **To make students aware of the methods, algorithms and tools used for gene and genome annotation**

# Genome Sequence

```
>P12345 Yeast chromosome1
GATTACAGATTACAGATTACAGATTACAGATTACAG
ATTACAGATTACAGATTACAGATTACAGATTACAGA
TTACAGATTACAGATTACAGATTACAGATTACAGAT
TACAGATTAGAGATTACAGATTACAGATTACAGATT
ACAGATTACAGATTACAGATTACAGATTACAGATTA
CAGATTACAGATTACAGATTACAGATTACAGATTAC
AGATTACAGATTACAGATTACAGATTACAGATTACA
GATTACAGATTACAGATTACAGATTACAGATTACAG
ATTACAGATTACAGATTACAGATTACAGATTACAGA
TTACAGATTACAGATTACAGATTACAGATTACAGAT
```

# Predict Genes

This server can accept sequences up to 1 million base pairs (1 Mbp) in length. If you have trouble with the web server or if you have a large number of sequences to process, request a local copy of the program (see instructions at the bottom of this page) or use the GENSCAN email server. If your browser (e.g., Lynx) does not support file upload or multipart forms, use the older version.

Organism: [Vertebrate ▼] Suboptimal exon cutoff (optional): [1.00 ▼]

Sequence name (optional): [                                          ]

Print options: [Predicted peptides only ▼]

Upload your DNA sequence file (one-letter code, upper or lower case, spaces/numbers ignored):
[                    ] [ Browse... ]

Or paste your DNA sequence here (one-letter code, upper or lower case, spaces/numbers ignored):

Document: Done

# The Result…

```
>P12346 Sequence 1
ATGTACAGATTACAGATTACAGATTACAGATTACAG
ATTACAGATTACAGATTACAGATTACAGATTACAGA
TTACAGATTACAGATTACAGATTACAGAT


>P12347 Sequence 2
ATGAGATTAGAGATTACAGATTACAGATTACAGATT
ACAGATTACAGATTACAGATTACAGATTACAGATTA
CAGATTACAGATTACAGATTACAGATTACAGATT


>P12348 Sequence 3
ATGTTACAGATTACAGATTACAGATTACAGATTACA
GATTACAGATTACAGATTACAGATTACA...
```

# Is This Annotated?

```
>P12346 Sequence 1
ATGTACAGATTACAGATTACAGATTACAGATTACAG
ATTACAGATTACAGATTACAGATTACAGATTACAGA
TTACAGATTACAGATTACAGATTACAGAT

>P12347 Sequence 2
ATGAGATTAGAGATTACAGATTACAGATTACAGATT
ACAGATTACAGATTACAGATTACAGATTACAGATTA
CAGATTACAGATTACAGATTACAGATTACAGATT

>P12348 Sequence 3
ATGTTACAGATTACAGATTACAGATTACAGATTACA
GATTACAGATTACAGATTACAGATTACA...
```

# How About This?

```
>P12346 Sequence 1
MEKGQASRTDHNMCLKPGAAERTPESTSPASDAAGG
IPQNLKGFYQALNNWLKDSQLKPPPSSGTREWAALK
LPNTHIALD

>P12347 Sequence 2
MKPQRTLNASELVISLIVESINTHISHOUSEPLEAS
EWILLITALLCEASE

>P12348 Sequence 3
MQWERTGHFDALKPQWERTYHEREISANTHERS...
```

# Gene Annotation*

- *<u>Annotation</u>* – **to identify and describe all the physico-chemical, functional and structural properties of a gene including its DNA sequence, protein sequence, sequence corrections, name(s), position, function(s), abundance, location, mass, pI, absorptivity, solubility, active sites, binding sites, reactions, substrates, homologues, $2^o$ structure, 3D structure, domains, pathways, interacting partners**

# Gene Annotation

$$\parallel$$

# Protein Annotation

# Protein/Gene vs. Proteome/ Genome Annotation

- **Gene/Protein annotation is concerned with one or a small number (<50) genes or proteins from one or several types of organisms**

- **Genome/Proteome annotation is concerned with entire proteomes (>2000 proteins) from a specific organism (or for all organisms)** *- need for speed*

# Different Levels of Annotation*

- **Sparse** – typical of archival databanks like GenBank, usually just includes name, depositor, accession number, dates, ID #

- **Moderate** – typical of many curated protein sequence databanks (UniProt or TrEMBL)

- **Detailed** – not typical (occasionally found in organism-specific databases)

# Different Levels of Database Annotation*

- **GenBank** (large # of sequences, minimal annotation)

- **TrEMBL** (large # of sequences, slightly better [computer] annotation)

- **UniProtKB** (small # of sequences, even better [hand] annotation)

- **Organsim-specific DB** (very small # of sequences, best annotation)

# GenBank Annotation (GST)

# UniProtKB Annotation (GST)

# The CCDB*

## The CyberCell Database (CCDB): a comprehensive, self-updating, relational database to coordinate and facilitate *in silico* modeling of *Escherichia coli*

Shan Sundararaj, Anchi Guo, Bahram Habibi-Nazhad, Melania Rouani[1], Paul Stothard, Michael Ellison[1] and David S. Wishart

Faculty of Pharmacy and Pharmaceutical Sciences and [1]Department of Biochemistry, University of Alberta, Edmonton, Alberta T6G 2N8, Canada

**http://ccdb.wishartlab.com/CCDB/**

# CCDB Annotation (GST)

# CCDB Annotation

| | |
|---|---|
| EC_Number | 2.5.1.18 |
| #_Amino_Acids_T | 201 (Translated Protein) |
| #_Amino_Acids_M | 200 (Mature Protein) |
| Calculated_Mw_(Daltons)_T | 22868.4 (Translated Protein) |
| Calculated_Mw_(Daltons)_M | 22737.2 (Mature Protein) |
| Theoretical_pI_T | 5.85 |
| Theoretical_pI_M | 5.86 |
| Observed_pI | Not av |
| Sequence_Verified | 1) Arc bac 2) Nish mut resi -325 |
| Protein_Sequence | >GT_E MKLFY TLLTE YKPTV FMQRM >GT_E KLFYK LLTEG PTVRA RMAEF |

| PROSITE_Motif | 1) PS00374 Methylated-DNA--protein-cysteine methyltransferase active site. [LIVMF]-P-C-H-R-[LIVMF]. |
|---|---|
| | 2) PS00462 Gamma-glutamyltranspeptidase signature. T-[STA]-H-x-[ST]-[LIVMA]-x(4)-G-[SN]-x-V-[STA]-x-T-x-T-[LIVM]-[NE]-PA x(1,2)-[FY]-G. |
| | 3) PS01311 Prolipoprotein diacylglyceryl transferase signature. G-R-x-[GA]-N-F-[LIVMF]-N-x-E-x(2)-G. |
| | 4) PS00197 2Fe-2S ferredoxins, iron-sulfur binding region signature. C-{C}-{C}-[GA]-{C}-C-[GAST]-{CPDEKRHFYW}-C. |
| | 5) PS01039 Bact ature. G-[FYIL]-[DE]- |
| Other_Sites | 1) Active Site (1( 2) Active Site (10 |
| #_Transmembrane_Regions | No |
| Cys/Met_Translated | 1.00 %Cys 2.49 %Met 3.49 %Cys+%Me |
| Cys/Met_Mature | 1.00 %Cys 2.00 %Met 3.00 %Cys+%Me |
| | MKLFYKPGACSLA CEEEEECCHHHHH QVPALLLDDGTLL CCEEEECCCCEEE |

| Sec_Structure_(PDB; 1A0F) | CCEEEECCCCEEECHHHHHHHHHHHCCCCCCCCCCCCCHHHHHHHHHHH IATELHKGFTPLFRPDTPEEYKPTVRAQLEKKLQYVNEALKDEHWICGQR HHHHHHHHHHHCCCCCCCCCHHHHHHHHHHHHHHHHHHCCCCCCCCCC FTIADAYLFTVLRWAYAVKLNLEGLEHIAAFMQRMAERPEVQDALSAEGL CCHHHHHHHHHHHHCCCCCCCCCHHHHHHHHHHHCHHHHHHHHHHCC K C |
|---|---|
| PDB_Accession | 1A0F |
| 3D_View | JAVA 3D View (PDB) |
| Resolution | 2.1 Angstroms |
| Structure_Class/Fold_Class | All Alpha |
| Quaternary_Structure | homodimeric $A_2$ Complex of gst |
| Interacting_Partners | 1) gst |
| Cofactor | None |
| Metal_Ion | None |
| Kcat_Value_[1/min] | Not available |
| Specific_Activity_uM/min/mg | Not available |

Document: Done

# CCDB Contents*

- **Functional info (predicted or known)**
- **Sequence information (sites, modifications, pI, MW, cleavage)**
- **Location information (in chromosome & cell)**
- **Interacting partners (known & predicted)**
- **Structure ($2^o$, $3^o$, $4^o$, predicted)**
- **Enzymatic rate and binding constants**
- **Abundance, copy number, concentration**
- **Links to other sites & viewing tools**
- **Integrated version of all major Db's**
- *70+ fields for each entry*

# GeneCards Content

- **Aliases**
- **Databases**
- **Disorders**
- **Domains**
- **Drugs/Cmpds**
- **Expression**
- **Function**
- **Location**

- **Orthologs/Paralogs**
- **Pathways and Interactions**
- **References**
- **Proteins/MAbs**
- **SNPs**
- **Transcripts**
- **Gene Maps**

**http://www.genecards.org/index.shtml**

# GeneCards Annotation

# GeneCards Annotation

# Ultimate Goal...

- **To achieve the same level of protein/ proteome annotation as found in CCDB or GeneCards for all genes/proteins -- _automatically_**

# How?

# Annotation Methods*

- **Annotation by homology (BLAST)**
  - **requires a large, well annotated database of protein sequences**
- **Annotation by sequence composition**
  - **simple statistical/mathematical methods**
- **Annotation by sequence features, profiles or motifs**
  - **requires sophisticated sequence analysis tools**

# Annotation by Homology*

- **Statistically significant sequence matches identified by BLAST searches against GenBank (nr), UniProt, DDBJ, PDB, InterPro, KEGG, Brenda, STRING**

- **Properties or annotation inferred by name, keywords, features, comments**

*Databases Are Key*

# Sequence Databases*

- **GenBank**
  - www.ncbi.nlm.nih.gov/

- **UniProt/trEMBL**
  - http://www.uniprot.org/

- **DDBJ**
  - http://www.ddbj.nig.ac.jp

# Structure Databases*

- **RCSB-PDB**
  - **http://www.rcsb.org/pdb/**
- **PDBe**
  - **http://www.ebi.ac.uk/pdbe/**
- **CATH**
  - **http://www.cathdb.info/**
- **SCOP**
  - **http://scop.mrc-lmb.cam.ac.uk/scop/**

# Interaction Databases*

- **STRING**
  - **http://string.embl.de/**
- **DIP**
  - **http://dip.doe-mbi.ucla.edu/**
- **PIM**
  - **http://www.ebi.ac.uk/intact/ main.xhtml**
- **MINT**
  - **http://mint.bio.uniroma2.it/ mint/Welcome.do**

# Bibliographic Databases

- **PubMed Medline**
  - **http://www.ncbi.nlm.nih.gov/PubMed/**
- **Google Scholar**
  - **http://scholar.google.ca/**
- **Your Local eLibrary**
  - **www.XXXX.ca**
- **Current Contents**
  - **http://science.thomsonreuters.com/**

# Annotation by Homology
## An Example

- **76 residue protein from *Methanobacter thermoautotrophicum* (newly sequenced)**
- **What does it do?**

- `MMKIQIYGTGCANCQMLEKNAREAVKELGIDAE`
  `FEKIKEMDQILEAGLTALPGLAVDGELKIMGRV`
  `ASKEEIKKILS`

# PSI BLAST



- **PSI-BLAST – position specific iterative BLAST**

- **Derives a position-specific scoring matrix (PSSM) from the multiple sequence alignment of sequences detected above a given score threshold using protein BLAST**

- **This PSSM is used to further search the database for new matches, and is updated for subsequent iterations with these newly detected sequences**

- **PSI-BLAST provides a means of detecting distant relationships between proteins**

# PSI-BLAST

# PSI-BLAST*

Run PSI-Blast iteration 2

end:

- means that the alignment score was below the threshold on the previous iteration
· means that the alignment was checked on the previous iteration

**Sequences with E-value BETTER than threshold**

| | | | Score | E |
|---|---|---|---|---|
| ences producing significant alignments: | | | (bits) | Value |
| pir\|\|F69219 | conserved hypothetical protein MTH895 – Methanobacte... | | 110 | 2e-24 |
| gb\|AAB52989.1\| | (U72238) ORFR5 [Anabaena PCC7120] | | 107 | 2e-23 |
| sp\|Q58001\|Y581 METJA | HYPOTHETICAL PROTEIN MJ0581 >gi\|2128389\|pir... | | 103 | 2e-22 |
| pir\|\|F72306 | conserved hypothetical protein – Thermotoga maritima... | | 99 | 4e-21 |
| pir\|\|H69530 | conserved hypothetical protein AF2248 – Archaeoglobu... | | 98 | 1e-20 |
| sp\|P42035\|THIO METTM | PROBABLE THIOREDOXIN (GLUTAREDOXIN-LIKE PRO... | | 42 | 9e-04 |
| sp\|O26898\|THIO METTH | PROBABLE THIOREDOXIN (GLUTAREDOXIN-LIKE PRO... | | 41 | 0.001 |

Run PSI-Blast iteration 2

# PSI-BLAST*

|  | | | Score | E |
|---|---|---|---|---|
| quences producing significant alignments: | | | (bits) | Value |
| 🟢 ☑ | pir\|\|S54843 | glutaredoxin-like protein – Pyrococcus furiosus >gi\|... | 99 | 3e-21 |
| 🟢 ☑ | pir\|\|H71239 | probable glutaredoxin-like protein – Pyrococcus hori... | 99 | 4e-21 |
| 🟢 ☑ | pir\|\|F69219 | conserved hypothetical protein MTH895 – Methanobacte... | 98 | 1e-20 |
| 🟢 ☑ | gb\|AAB52989.1\| | (U72238) ORFR5 [Anabaena PCC7120] | 96 | 5e-20 |
| 🟢 ☑ | pir\|\|F75204 | glutaredoxin-like protein PAB2245 – Pyrococcus abyss... | 96 | 5e-20 |
| 🟢 ☑ | pir\|\|G72322 | glutaredoxin-related protein – Thermotoga maritima (... | 89 | 3e-18 |
| 🟢 ☑ | sp\|Q58001\|Y581 METJA | HYPOTHETICAL PROTEIN MJ0581 >gi\|2128389\|pir... | 89 | 6e-18 |
| 🟢 ☑ | pir\|\|F72306 | conserved hypothetical protein – Thermotoga maritima... | 88 | 9e-18 |
| 🟢 ☑ | pir\|\|H69530 | conserved hypothetical protein AF2248 – Archaeoglobu... | 87 | 2e-17 |
| 🟢 ☑ | sp\|P42035\|THIO METTM | PROBABLE THIOREDOXIN (GLUTAREDOXIN-LIKE PRO... | 87 | 2e-17 |
| 🟢 ☑ | pir\|\|A72669 | probable glutaredoxin-like protein APE0775 – Aeropyr... | 86 | 4e-17 |
| 🟢 ☑ | sp\|O26898\|THIO METTH | PROBABLE THIOREDOXIN (GLUTAREDOXIN-LIKE PRO... | 86 | 5e-17 |
| 🟢 ☑ | sp\|O28137\|THIO ARCFU | PROBABLE THIOREDOXIN >gi\|7450264\|pir\|\|A6951... | 85 | 6e-17 |
| 🟢 ☑ | sp\|Q57755\|THIO METJA | THIOREDOXIN >gi\|2129305\|pir\|\|D64338 thiored... | 78 | 1e-14 |
| 🟢 ☑ | sp\|P22904\|YME3 THIFE | HYPOTHETICAL 9.0 KD PROTEIN IN MOBE 3'REGIO... | 73 | 3e-13 |
| ᴡ ☑ | pir\|\|E70340 | glutaredoxin-like protein – Aquifex aeolicus >gi\|298... | 45 | 1e-04 |

Run PSI-Blast iteration 6

# Conclusions

- **Protein is a thioredoxin or glutaredoxin (function, family)**

- **Protein has thioredoxin fold ($2^o$ and 3D structure)**

- **Active site is from residues 11-14 (active site location)**

- **Protein is soluble, cytoplasmic (cellular location)**

# Annotation Methods

- **Annotation by homology (BLAST)**
  - **requires a large, well annotated database of protein sequences**
- **Annotation by sequence composition**
  - **simple statistical/mathematical methods**
- **Annotation by sequence features, profiles or motifs**
  - **requires sophisticated sequence analysis tools**

# Annotation by Composition*

- **Molecular Weight**

- **Isoelectric Point**

- **UV Absorptivity**

- **Hydrophobicity**

# Where To Go



http://www.expasy.ch/tools/#proteome

# Molecular Weight*

- **Useful for SDS PAGE and 2D gel analysis**
- **Useful for deciding on SEC matrix**
- **Useful for deciding on MWC for dialysis**
- **<u>Essential</u> in synthetic peptide analysis**
- **<u>Essential</u> in peptide sequencing (classical or mass-spectrometry based)**
- **<u>Essential</u> in proteomics and high throughput protein characterization**

# Molecular Weight*

- **Crude MW calculation: MW = 110 X Numres**

- **Exact MW calculation: MW = $\Sigma nAA_i$ x $MW_i$**

- **Remember to add 1 water (18.01 amu) after adding all res.**

- **Corrections for CHO, PO4, Acetyl, CONH2**

| Amino Acid Residue Weights | | | |
|---|---|---|---|
| Residue | Weight | Residue | Weight |
| A | 71.08 | M | 131.21 |
| C | 103.14 | N | 114.11 |
| D | 115.09 | P | 97.12 |
| E | 129.12 | Q | 128.14 |
| F | 147.18 | R | 156.2 |
| G | 57.06 | S | 87.08 |
| H | 137.15 | T | 101.11 |
| I | 113.17 | V | 99.14 |
| K | 128.18 | W | 186.21 |
| L | 113.17 | Y | 163.18 |

# Molecular Weight & Proteomics



**2-D Gel**



**QTOF Mass Spectrometry**

# Isoelectric Point*

- **The pH at which a protein has a net charge=0**

- $$Q = \Sigma \, Ni/(1 + 10^{pH-pKi})$$

This is a transcendental equation



pH / Charge

| pKa Values for Ionizable Amino Acids | | | |
|---|---|---|---|
| Residue | pKa | Residue | pKa |
| C | 10.28 | H | 6 |
| D | 3.65 | K | 10.53 |
| E | 4.25 | R | 12.43 |

# UV Absorptivity*

- **$OD_{280}$ = (5690 x #W + 1280 x #Y)/MW x Conc.**
- **Conc. = $OD_{280}$ x MW/(5690 X #W + 1280 x #Y)**



Very useful for measuring protein concentration

# Hydrophobicity*

- **Average Hphob calculation: $H_{ave} = (\Sigma nAA_i \times Hphob_i)/N$**

- **Indicates Solubility, stability, location**

- **If $H_{ave} < 1$ the protein is soluble**

- **If $H_{ave} > 1$ it is likely a membrane protein**

| Kyte / Doolittle Hyrophobicity Scale | | | |
|---|---|---|---|
| Residue | Hphob | Residue | Hphob |
| A | 1.8 | M | 1.9 |
| C | 2.5 | N | -3.5 |
| D | -3.5 | P | -1.6 |
| E | -3.5 | Q | -3.5 |
| F | 2.8 | R | -4.5 |
| G | -0.4 | S | -0.8 |
| H | -3.2 | T | -0.7 |
| I | 4.5 | V | 4.2 |
| K | -3.9 | W | -0.9 |
| L | 3.8 | Y | -1.3 |

# Annotation Methods

- **Annotation by homology (BLAST)**
  - **requires a large, well annotated database of protein sequences**
- **Annotation by sequence composition**
  - **simple statistical/mathematical methods**
- **Annotation by sequence features, profiles or motifs**
  - **requires sophisticated sequence analysis tools**

# Where To Go



http://www.expasy.ch/tools/#proteome

# Sequence Feature Databases

- **PROSITE - http://www.expasy.ch/prosite/**

- **InterPro - http://www.ebi.ac.uk/interpro/**

- **PPT-DB - http://www.pptdb.ca/**

**To use these databases just submit your PROTEIN sequence to the database and download the output.  They provide domain information, predicted disulfides, functional sites, active sites, secondary structure – IF THERE IS A MATCH**

# Using Prosite

# Prosite Output

# What if your Sequence doesn't match to Something in the Database?

- **Don't worry**

- **You can use prediction programs and freely available web servers that use machine learning, neural networks, HMMs and other cool bioinformatic tricks to predict some of the same things that your database matching tools try to identify**

# What Can Be Predicted?*

- **O-Glycosylation Sites**
- **Phosphorylation Sites**
- **Protease Cut Sites**
- **Nuclear Targeting Sites**
- **Mitochondrial Targ Sites**
- **Chloroplast Targ Sites**
- **Signal Sequences**
- **Signal Sequence Cleav.**
- **Peroxisome Targ Sites**

- **ER Targeting Sites**
- **Transmembrane Sites**
- **Tyrosine Sulfation Sites**
- **GPInositol Anchor Sites**
- **PEST sites**
- **Coil-Coil Sites**
- **T-Cell/MHC Epitopes**
- **Protein Lifetime**
- **A whole lot more….**

# Cutting Edge Sequence Feature Servers*

- **Membrane Helix Prediction**
  - **http://www.cbs.dtu.dk/services/TMHMM-2.0/**

- **T-Cell Epitope Prediction**
  - **http://www.syfpeithi.de/home.htm**

- **O-Glycosylation Prediction**
  - **http://www.cbs.dtu.dk/services/NetOGlyc/**

- **Phosphorylation Prediction**
  - **http://www.cbs.dtu.dk/services/NetPhos/**

- **Protein Localization Prediction**
  - **http://psort.ims.u-tokyo.ac.jp/**

# 2º Structure Prediction*

- **PredictProtein-PHD (72%)**

  – **http://www.predictprotein.org**

- **Jpred (73-75%)**

  – **http://www.compbio.dundee.ac.uk/~www-jpred/**

- **PSIpred (77%)**

  – **http://bioinf.cs.ucl.ac.uk/psipred/**

- **Proteus2 (78-90%)**

  – **http://www.proteus2.ca/proteus2/**

# Putting It All Together

**http://basys.ca/basys/cgi/submit.pl**

# BASys

- **BASys (Bacterial Annotation System) is a web server that performs automated, in-depth annotation of bacterial genomic sequences**

- **It accepts raw DNA sequence data and an optional list of gene identification information and provides extensive textual and hyperlinked image output**

# BASys

- **BASys uses more than 30 programs to determine nearly 60 annotation subfields for each gene, including:**

- **Gene/protein name, GO function, COG function, possible paralogues and orthologues, molecular weight, isoelectric point, operon structure, subcellular localization, signal peptides, transmembrane regions, secondary structure, 3-D structure and reactions**

# Submitting to BASys

# Wait…

# BASys Output

# BASys Output (Map)

# BASys Output (Map)

# BASys Output (Gene Link)

# Conclusion

- **Genome annotation is the same as proteome annotation – required after any gene sequencing and gene ID effort**
- **Can be done either manually or automatically**
- **Need for high throughput, automated "pipelines" to keep up with the volume of genome sequence data**
- **Area of active research and development with about ½ of all bioinformaticians working on some aspect of this process**