MCB 5472 Lecture #5:
Gene Prediction and Annotation
February 24, 2014

## Note on the assignment

- Depending on your settings PSI-BLAST can take a while to run
  - Do not leave this until the last minute!
- Recall from Assignment Lecture #1: `nohup` can allow you to leave a job running on the cluster
  - E.g., `nohup [task] & > nohup.out`

## Do you have a DNA sequence…

- Limited utility by itself
- Annotations describe what the DNA does
  - Structural: what features are present on the DNA?
  - Functional: what do those features do?

## How to annotate: 2 methods

1. From first principles:
   - Experimental data in the literature
   - Algorithmic rules
2. From orthology / homology to previously annotated sequences

## Annotation accuracy

- Manual annotation from experimental data in the literature is highly accurate
  - Although not all experiments are unequivocal
- Annotations using algorithms can be quite accurate
  - Depends on the complexity of the problem the algorithm is trying to solve
- Annotations based on orthology relies on the assumption that function is conserved
  - Depends on how rigorously orthologs are defined
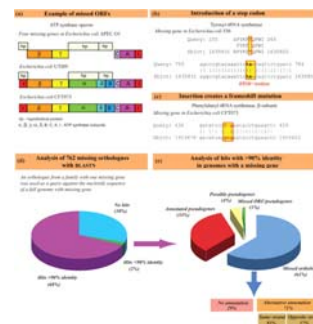  - Depends on functions not changing over time

## Gene annotation

- Gene and protein annotation is typically algorithmic
- Genes and proteins have specific features that algorithms use to define them
- Algorithms for bacteria and archaea work quite well, eukaryotes more difficult because of additional complexity, e.g., splicing

## Prokaryote gene finding

- Glimmer, GeneMark: Markov Models
  - Genes modeled based on differences between coding and non-coding regions
    - E.g., typically start with ATG, end with stop codon
    - E.g., ORF overlap
    - E.g., ribosome binding regions
  - Often have difficulty to decide which strand is coding.
- Prodigal: summed likelihood of finding individual gene features
- Can be challenged by %GC bias
  - Better performance by training on known genome annotations

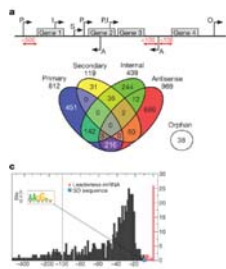Fig. 1. Analysis of ORFs missing in one out of 30 completely annotated Escherichia genomes.



Poptsova M S , and Gogarten J P Microbiology 2010;156:1909-1917

SGM

Microbiology

## Remember: genes are not transcripts!

- 5' mRNA analysis in *Helicobacter pylori* shows much greater transcript diversity than evident from simple gene annotations
- Most NCBI annotations equate genes with transcripts



Sharma et al. 2010 Nature 464:250-255

## Eukaryotic gene finding

- E.g., Augustus, GeneMark-ES
- *ab initio* methods work less well compared to prokaryotic genes
  - More complicated transcripts (e.g., splice variants)
  - Less information at promoter (e.g., Prodigal uses Shine-Delgarno sequences; -35 and -10 regions vs. single TATA box)
- NCBI annotations more clearly separate genes (includes pseudogenes), mRNA (typically spliced) & protein (spliced like mRNA)

## Adding information to gene annotations

1. Combine multiple prediction methods
   - For prokaryotes, typically longest transcript chosen
   - For eukaryotes, typically all splice variants kept
2. Search for homologous genes in related taxa
   - True genes will be evolutionarily conserved
   - Annotation errors can be propagated
     - Annotations do not specify the evidence supporting them
3. Integrate RNAseq
   - Augustus can incorporate into its predictions directly
   - Rare for prokaryotes
   - Requires genes be expressed and detectable

## Metagenomes and single-cell genomes

- Assemblies are typically much more fragmented than those of cultured microbes
- Requires dedicated gene prediction methods
  - Training information often missing/obscured
  - Gene fragments obscure genomic features used for gene prediction

## Non-coding RNAs

- Some HMM-based software
  - RNAMMER (ribosomal RNAs)
  - tRNAscan-SE (tRNAs)
- Rfam: database of non-coding RNA families
  - Curated sequence alignments taking into account secondary structures
  - Infernal: software for searching DNA sequence databases using structured RNA molecule profiles
    - Takes RNA secondary structure into account via "covariance models"
  - Sister project to Pfam (see later)
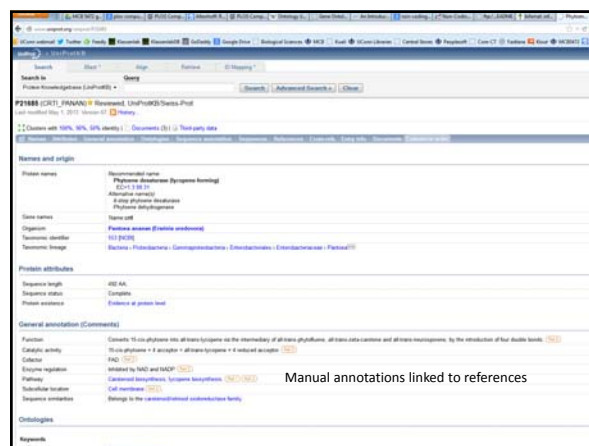
# Functional annotations

## Manual annotation

- Low-throughput
- High accuracy

## SwissProt

- Started 1986 at the Swiss Institute for Bioinformatics, later developed at the European Bioinformatics Institute
- Goal: providing reliable protein sequences having a high level of annotation
  - Directly curated from literature information
  - Contrast to NCBI: a sequence repository with some automated annotation pipelines
- Current version (2014_02): 542,503 sequences annotated from 22,6190 references

## UniProt

- Ultimately manual annotation couldn't keep up, parallel TrEMBL database created using automated annotation
- UniProtKB stores combined SwissProt/TrEMBL databases, incorporates Protein Information Resource (PIR), built on M. Dayhoff's atlas
- Syncs with EMBL/DDBJ/GenBank nucleotide databases
- Hosts several protein annotation schemes
- ExPASy – major proteomics analysis resource
- www.uniprot.org



Manual annotations linked to references

## Ecocyc – an example manually edited model organism database
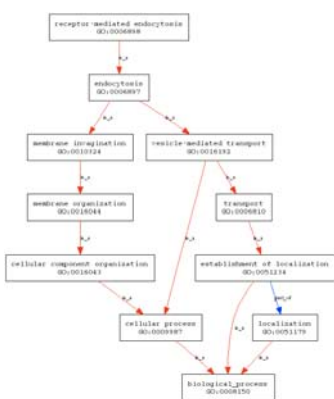
- http://ecocyc.org/



## Ontologies

- Manual annotations originally used free-text labels, not standardized
  - Problem: free text is difficult for computers to make use of
- Ontologies: knowledge representation using standardized terms and interrelationships
  - Amenable to computation

## E.g., GO

- Controlled vocabulary
- Defined relationships
- "Directed acyclic graph"
  - Links are directional
  - No individually circular paths



http://lopasgen677s09.weebly.com/gene-ontology.html

## Gene Ontology (GO)

- http://www.geneontology.org/
- Consortium that defines standardized terms and relationships
- Centered on model organism databases
  - E.g., human, mouse, Drosophila, E.coli
  - Most curation derived from these sources, but do extend more broadly
- Linked and mapped to many other resources
- Used by many computational analysis tools

## GO domains

- GO is divided into three domains, encompassing three separate functional properties
  - Biological process: what it does
  - Molecular function: how it does it
  - Cellular component: where it does it

## GO evidence codes

- GO uniquely has an ontology to describe the evidence supporting annotations



http://www.geneontology.org/GO.evidence.tree.shtml

## GO on UniProt



## GO on UniProt

- Some GO annotations added manually
- Some mapped to other term databases



## Annotation families

- There are many different types of protein annotations, often with different foci and methods
- Hand vs. automatically generated
- Entire vs partial proteins

## Pfam

- http://pfam.sanger.ac.uk/
- Originally constructed in the late 1990's for annotation of the *C. elegans* genome
- Developed & maintained by the Sanger Institute and S. Eddy (now Howard Hughes)
- Purpose: to overcome the % alignment problem inherit to BLAST
  - i.e., BLAST hits may not reflect homology over the entire query and/or reference sequence
- Currently (v27.0) 14,831 manually curated protein domain families

## Pfam

- Pfam-A: manually selected and aligned alignments and HMMs of protein domains
  - v27.0: 14,831 families
  - At least 1 domain in 80% of proteins in UniProt
    - Figure is still scaling with database sizes
    - Represents 58% of total sequence in UniProt
- Pfam-B: automatically-generated families for domains not in Pfam-A
  - Mostly families with only a few members

## Pfam example

## Slide 1: Pfam example

Family: ABC_tran (PF00005)



## Slide 2: Pfam example

Family: ABC_tran (PF00005)



## Slide 3

# Clusters of Orthologus Groups (COG)

- One of the earliest attempts to define protein families by orthology (Tatusov et al 1997 Science)
- Used BLAST between proteins from multiple genomes to define triangles, i.e., triplets where each is a best match to the others



Kristensen et al. 2010 Bioinformatics 26:1481-1487

## Slide 4

# COG triangles

- Allows single-direction best hits
- Start with central triangle and add edges whenever possible
- Causes paralogs to be linked
- Allows distant & fast evolving homologs to be linked through intermediates



Sold lines: RBHs
Dotted lines: single direction

Tatusov et al 1997 Science 278: 631-637

## Slide 5

# COGs

- Bacterial COGs not updated often (last 2003)
- COGs more recently defined for other groups:
  - KOGs (eukaryotes)
  - arCOGs (archaea)
  - POGs (phages)
- Each COG family has a free-text annotation
  - 4873 families total
- Grouped into 24 superfamilies
  - COGs can belong to >1 superfamilies
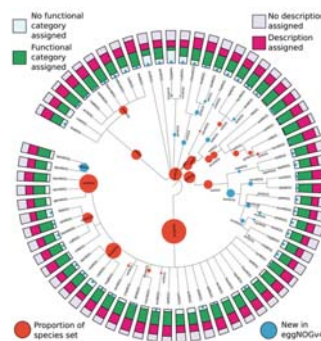
## Slide 6

# eggNOG

- 'evolutionary genealogy of genes: Non-supervised Orthologous Groups'
- Constructed & maintained by EMBL (Peer Bork)
- Attempt to extend and update COG/KOG database annotations without requiring manual annotations (which do not scale)

## eggNOG: method

- Use BLAST/fasta/Smith-Watterman alignments to find best matches
- Represent in-paralogs by single sequences
- Map sequences to COG/KOGs
- Triangle cluster non-matching sequences
- Add single RBH hits to clusters
- Automatically split multi-domain proteins
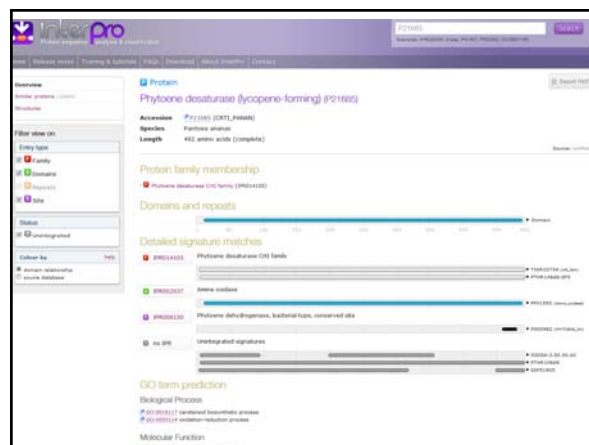- Derive annotations by consensus within groups derived from multiple annotation sources

## eggNOG:

- http://eggnog.embl.de/version_4.0.beta/
- 107 different annotation levels
- 1.7 million ortholog groups
- 7.7 million proteins
- Probably the currently most comprehensive ortholog database
- Can use to construct PSSMs/HMMs



## Interpro

- Classifies proteins according to a combination of multiple protein motifs
- Multiple sources synthesized into single Interpro classification system
  - Four broad annotation types: Family, Domains, Repeats, Sites
- Interpro terms mapped to GO
- InterProScan – resource to annotate proteins using all member databases
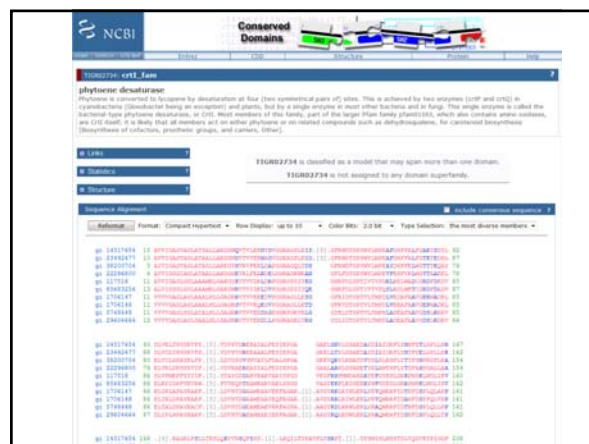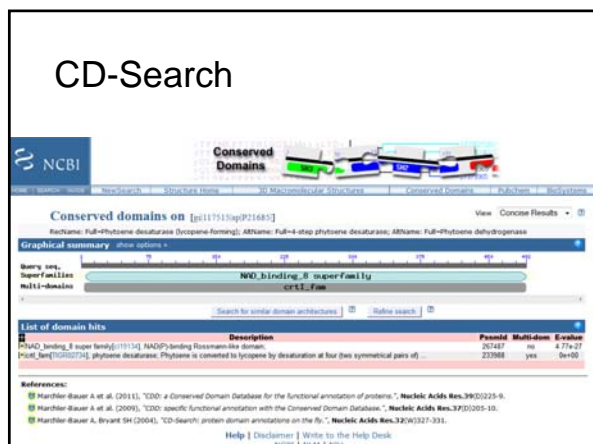  - HMM and regular expression-based classifications



## Interpro: member databases

- Pfam (domains, curated; Sanger)
- PROSITE (diagnostic motifs; SIB)
- HAMAP (homologs, curated; SIB)
- PRINTS (conserved motifs; U. Manchester)
- ProDom (domains, automatic via PSI-BLAST; PRABI Villerubanne)
- SMART (domains and architectures esp. signaling, curated; EMBL)
- TIGRFAMs (homologs, curated; JCVI)
- PIRSF (homologs & domains, ; Georgetown)
- SUPERFAMILY (structures, curated, U Bristol)
- CATH-Gene3D (homologs, mapped to structures, automatic via Markov clustering; University College London)
- PANTHER (functional homologs, curated, USC)

## Conserved Domains Database (CDD)

- Protein classification database maintained by NCBI
- CDD database based on domains curated by NCBI using structural alignments
- Also includes external resources: Pfam, SMART, COG, PRK, TIGRFAM
- Downloadable PSSMs for each CDD family for querying via RPS-BLAST
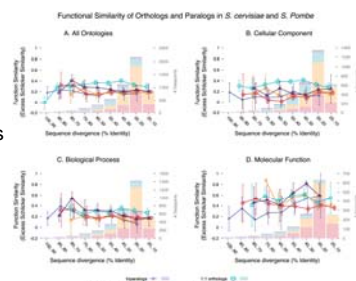
## CD-Search





## Are functions actually conserved?

- All of the protein annotation methods that we have discussed assume the hypothesis that function is evolutionarily conserved
- But we know that this can be confounded by duplication/loss and xenology
  - Can be addressed by better methods of determining orthology
  - Not typically accommodated by annotation databases
- Even orthologous functions can drift and/or be promiscuous
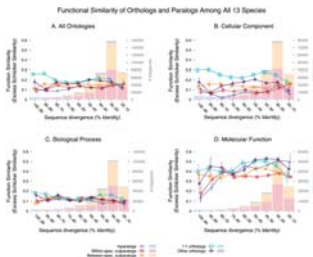
## Are functions actually conserved?

- Compare curated GO annotations of orthologs and paralogs
  - Corrected for annotation biases
- Functions of orthologs more similar than paralogs, but not perfectly



Altenhoff et al. (2012) PLoS Comput. Biol. 8:e1002514

## Are functions actually conserved?

- Same for 13 species instead of just 2
  - Paralogy potentially a greater confounder
  - Differences in GO annotation completeness



Altenhoff et al. (2012) PLoS Comput. Biol. 8:e1002514

## Are functions actually conserved?

- Define function as expression similarity between same human and mice tissues
- Same trend: ortholog function more conserved than paralogs, not absolute



Chen & Zhange (2012) PLoS Comput. Biol. 8:e1002784

## Are functions actually conserved?

- Yes, but not perfectly even for highly conserved sequences
- Likely depends on definition of "function"
- Annotated functions are likely quite broad in most cases

## Protein database vs. pathways and reactions

- Protein databases are based on homology
  - Hypothesis that function is conserved
- Reaction databases classify function without reference to homology
  - Function can be due to evolutionary convergence
  - GO is an example of this we have already seen
- Reaction and pathway annotations are therefore closer to function but further from underlying evolutionary mechanism

## Enyzme Commission

- One of the oldest functional annotation schemes, arising out of biochemistry
- Four part numerical nomenclature having increasing specificity
  - EC 3: hydrolases
  - EC 3.4: hydrolases acting on peptide bonds
  - EC 3.4.11: hydrolases cleaving amino-terminal amino acids from a peptide
  - EC 3.4.11.4: hydrolases cleaving amino-terminal amino acids from a tripeptide
- Database updates are infrequent

## Kyoto Encyclopedia of Genes and Genomes (KEGG)

- Manually edited pathway database
- Orthologs defined in other genomes
- Reactions combined into metabolic maps
  - Pathways are typically quite general
- Individual proteins can be freely queried via web
- Individual genomes can be annotated via KAAS server
- Underlying database NO LONGER FREE

KEGG example

**KAAS – KEGG Automatic Annotation Server**

http://www.genome.jp/kegg/kaas/

# Biocyc

- Collection of curated metabolic pathways
- Typically smaller modules compared to KEGG
- www.biocyc.org

BioCyc Tier 1: Intensively Curated Databases

| DATABASE | SCOPE | HIGHLIGHTS | ORGANIZATION |
|----------|-------|-----------|--------------|
| EcoCyc | Escherichia coli K-12 MG1655 Model Organism Database | • Literature curation of complete genome<br>• Information from 25,406 publications<br>• Transcriptional regulatory network<br>• Flux-balance metabolic model | SRI International |
| MetaCyc | Multiorganism Metabolic Pathway and Enzyme Database | • 2,097 metabolic pathways from 2460 organisms<br>• Extensive commentary<br>• Information from 37,570 publications | SRI International |
| HumanCyc | Homo sapiens | • 291 metabolic pathways | SRI International |
| AraCyc | Arabidopsis thaliana | • 400 metabolic pathways<br>• Information from 3,500 publications | S. Rhee, Department of Plant Biology, Carnegie Institution, USA |
| YeastCyc | Saccharomyces cerevisiae | • 152 metabolic pathways<br>• Information from 1,000 publications | SGD Curators, Stanford U., USA |
| LeishCyc | Leishmania major Friedlin | • 141 metabolic pathways | Bio21 Molecular Science and Biotechnology Institute, University of Melbourne |

BioCyc Tier 2: Computationally-Derived Databases Subject to Moderate Curation

BioCyc Tier 3: Computationally-Derived Databases Subject to No Curation

Create Your Own Pathway/Genome Database

# Metacyc example

# Metacyc example

# Metacyc example

# Uniprot cross-references Interpro, metacyc

## Annotation process

- Use web to find information about particular proteins
- Use individual tools separately on your genome
  - Allows most customization, proofchecking
  - Standard for eukaryotic genomes
- Use automatic prediction servers
  - Common for prokaryotes
  - E.g., NCBI, IMG (JGI), RAST, Megan, MAGE
  - Each vary slightly in algorithm, user engagement and proofchecking, visualization
- Transfer homology from previously-annotated sequences
  - Can propagate incorrect annotations
  - Can limit coverage