## Functional Annotation / Protein Annotation

- Protein structure
- Structure/Function prediction
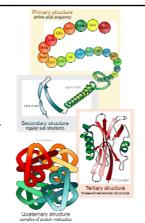- Functional Classification

## Functional prediction

- In absence of a proven function or mutant phenotype, a gene is no more than a transcribed piece of DNA

- If we can predict gene function in a genome, we can create a "parts list" of molecular functions that allow to make assumptions about the organism

- Actual experiments to elucidate or even just validate gene function usually take many years. Per gene!

## From sequence to function

DNA sequence (ATGAAGTTGATGGCAGCG...)

↓ *simple rule*

protein sequence (MKLMAA...)

protein structure

*prediction*
secondary structure    *ab initio* folding
sequence alignment / domain assignment

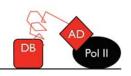protein function

*prediction*
sequence alignment / domain assignment

- Primary
  - AA sequence
  - Post-translational modifications
- Secondary
  - α-helix
  - β-sheet
- Tertiary
  - 3d folding driven by non-specific hydrophobic interactions
  - Stabilised by specific tertiary interactions
- Quaternary
  - Like tertiary but with multiple protein chains



# Protein domains

- Many proteins consist of several structural and functional entities called "domain"
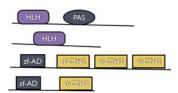


E.g. Transcription factor

DNA-binding and trans-activating domain structurally and functionally different.

- Often, proteins of the same family are 'mixtures' of a set of standard domains

# Domain architectures



Descriptions of domains:
http://pfam.sanger.ac.uk/

## *Ab initio* prediction of secondary structure from primary structure

- learning directly from X-ray structures
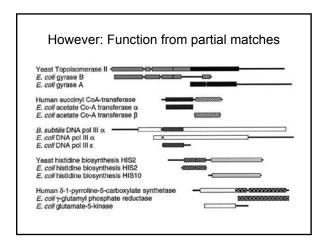- consideration of environment
- neural network-based training

e.g. Dor *et al.* - Achieving 80% ten-fold cross-validated accuracy for secondary structure prediction by large-scale training – Proteins 2007

---

## Prediction by alignment

- Primary sequence similarity >30% can be assumed to have the same 3D structure (but not necessarily function - beware of details!)

- Any available structural or functional data on orthologues (the "same" protein in a different organism) can be of great relevance.

---

## The 30% sequence identity rule



Sander & Schneider, Proteins 9(1):56-68, 1991

**BLAST**
bi-directional best hit

is often the most simple tool to establish orthology

BLAST of A$_{species 1}$ identifies B$_{species 2}$
BLAST of B$_{species 2}$ identifies A$_{species 1}$
as best hits
and sequence similarity > 30%.
BINGO.

## Intermediate sequences increase the detection of homology between sequences

1. 

2. 

3. 

1 and 2 are homologous
2 and 3 are homologous
1 and 3 are **not** homologous

- Multi-domain proteins introduce errors

Park *et al.* - Intermediate sequences increase the detection of homology between sequences – J Mol Bio. 1997

---

## However: Function from partial matches



Yeast Topoisomerase II
*E. coli* gyrase B
*E. coli* gyrase A

Human succinyl CoA-transferase
*E. coli* acetate Co-A transferase α
*E. coli* acetate Co-A transferase β

*B. subtilis* DNA pol III α
*E. coli* DNA pol III α
*E. coli* DNA pol III ε

Yeast histidine biosynthesis HIS2
*E. coli* histidine biosynthesis HIS2
*E. coli* histidine biosynthesis HIS10

Human δ-1-pyrroline-5-carboxylate synthetase
*E. coli* γ-glutamyl phosphate reductase
*E. coli* glutamate-5-kinase

---

## Identifying domains

- Structural domains show a high degree of conservation.

- Multiple sequence alignments [of the amino acids found in the domain] from multiple proteins (from different or the same species) can be used for a prediction.



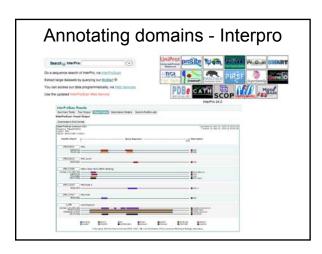sequence alignment          hidden Markov model

## Learning a profile



HMM profiles have a higher sensitivity than alignment-based approaches AND effectively compare the candidate sequence against all sequences in the alignment profile at the same time.

RKMAHHARERRRR...
RGLMHNELEKRRR...
KRARANELEKQMV...
RTAAHKQNERKMR...

## Limitations of profile HMMs

Higher-order relationships are not preserved (but this is also true for sequence alignments).

RKMAHHARERRRR...
RGLMHNELEKRRR...
KRARANELEKQMV...
RTAAHKQNERKMR...

. . .

Why are amino acids linked?



## Annotating domains - Interpro
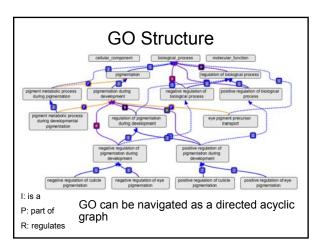
## Functional annotation

- For enzymes: EC number
- 6 groupings – Oxidoreductases, Transferases, Hydrolases, Lyases, Isomerases, Ligases



## Functional annotation

For other proteins: Gene Ontology, Reactome, KEGG

- GO is the *de facto* standard used by all major model organism databases

- Founded in 2000 by the GO Consortium lead by Michael Ashburner

- Database curators read the scientific literature and assign functional classifications along with an evidence code to proteins

- These annotations follow a controlled vocabulary that is organized into an ontology.

## GO Structure



I: is a
P: part of
R: regulates

GO can be navigated as a directed acyclic graph

## Inferences and logic



Annotation granularity and inconsistencies are a big issue:
1: Index finger *part of* Extremity
2: Ring finger *is a* Finger *part of* Hand *part of* Extremity



## Evidence codes



**Introduction**
Experimental Evidence Codes
  EXP: Inferred from Experiment — **the best!**
  IDA: Inferred from Direct Assay
  IPI: Inferred from Physical Interaction
  IMP: Inferred from Mutant Phenotype
  IGI: Inferred from Genetic Interaction
  IEP: Inferred from Expression Pattern
Computational Analysis Evidence Codes
  ISS: Inferred from Sequence or Structural Similarity
  ISO: Inferred from Sequence Orthology
  ISA: Inferred from Sequence Alignment
  ISM: Inferred from Sequence Model
  IGC: Inferred from Genomic Context
  RCA: Inferred from Reviewed Computational Analysis
Author Statement Evidence Codes
  TAS: Traceable Author Statement
  NAS: Non-traceable Author Statement
Curator Statement Evidence Codes
  IC: Inferred by Curator
  ND: No biological Data available
**Automatically-assigned Evidence Codes**
  IEA: Inferred from Electronic Annotation — **don't trust!**

| Species | Genes | Annot. |
|---|---|---|
| *Drosophila melanogaster*<br>FlyBase | 12517 | 72277<br>(16660 non-IEA) |
| *Escherichia coli*<br>EcoCyc & EcoliHub | 3543 | 38310<br>(1898 non-IEA) |
| *Ehrlichia chaffeensis Arkansas 3CVI* | 1091 | 2861<br>(2861 non-IEA) |
| *Gallus gallus*<br>GO Annotations @ EBI | 16306 | 70671<br>(2035 non-IEA) |

## References

- Dor *et al.* - Achieving 80% ten-fold cross-validated accuracy for secondary structure prediction by large-scale training – Proteins 2007
- Sander & Schneider - Database of homology-derived protein structures and the structural meaning of sequence alignment – Proteins 1991
- Park *et al.* - Intermediate sequences increase the detection of homology between sequences – J Mol Bio. 1997
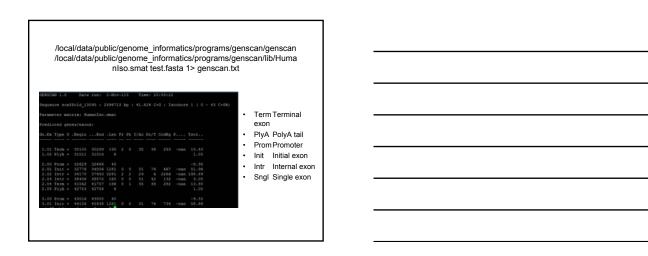
## References

- Must read: The documentation at
  – www.geneontology.org

## Assignment 2

- 20 min group presentations
- Annotate *Drosophila* species, compare annotations across species
- Groups on Moodle – either Gene or Protein

- ALL submit a copy of the presentation on Moodle; PER GROUP print out a copy of slides and give to me on day of presentations

**Online resources to explore in your own time - I**

**Resource centres:**

The EBI: www.ebi.ac.uk
The NCBI: www.ncbi.nlm.nih.gov
UCSC: genome.ucsc.edu

**Model organism databases:**
Budding yeast: www.yeastgenome.org
Worm: www.wormbase.org
Fly: www.flybase.org
Mouse: www.informatics.jax.org
Rat: rgd.mcw.edu
Zebrafish: zfin.org

**Database building tool:** www.intermine.org
**InterMine databases:**
www.flymine.org
intermine.modencode.org
ratmine.mcw.edu
yeastmine.yeastgenome.org
www.metabolicmine.org
www.flytf.org
mitominer.mrc-mbu.cam.ac.uk
targetmine.nibio.go.jp

**Ontologies:**
Gene ontology:
www.geneontology.org
Sequence ontology:
www.sequenceontology.org

**DNA sequence/ genomes:**
Ensembl: www.ensembl.org
Short read archive:
www.ncbi.nlm.nih.gov/sra

**Proteins:**
www.uniprot.org
www.ebi.ac.uk/interpro

**RNA:**
rfam.sanger.ac.uk

**Pathways:**
www.reactome.org
www.genome.jp/kegg/

---

/local/data/public/genome_informatics/programs/ncbi-blast-2.5.0+/bin/blastp -query dana-all-translation-r1.04.fasta -db /local/data/public/genome_informatics/assignment_2/swissprot_database/uniprot_sprot.fasta -out d_ana_results -outfmt 6 -num_threads 60



- Query ID
- Subject ID
- Percentage Identity of alignment
- Alignment length
- Mismatch count
- Gap open count
- Query start
- Query end
- Subject start
- Subject end
- E-value
- bitScore

---

/local/data/public/genome_informatics/programs/genscan/genscan /local/data/public/genome_informatics/programs/genscan/lib/HumanIso.smat test.fasta 1> genscan.txt



- Term Terminal exon
- PlyA PolyA tail
- Prom Promoter
- Init Initial exon
- Intr Internal exon
- Sngl Single exon

/local/data/public/genome_informatics/programs/hmmer-3.1b2-linux-intel-
x86_64/binaries/hmmscan --cpu 60
/local/data/public/genome_informatics/assignment_2/Pfam-A.hmm dana-all-translation-
r1.04.fasta 1> dana_test.txt

--tblout <f> : table of per-sequence hits          --domtblout <f> : table of per-domain hits