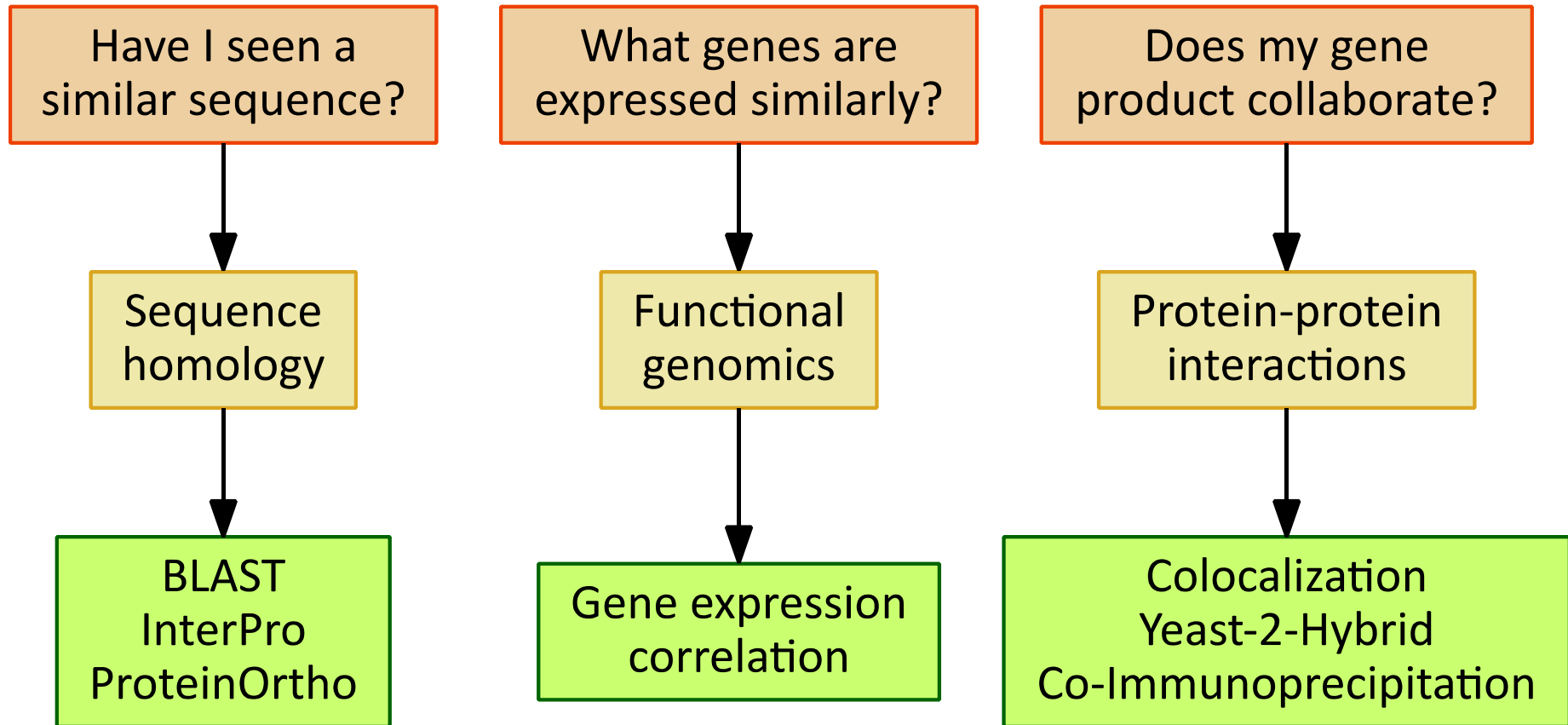# What does this gene do? Bioinformatics for function

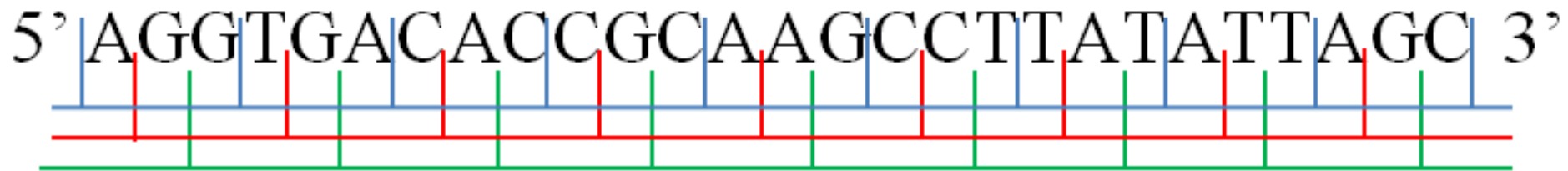DAVID L. TABB, PH.D.

MOLECULAR BIOLOGY AND HUMAN GENETICS

STELLENBOSCH UNIVERSITY

# Outline

| Have I seen a similar sequence? | What genes are expressed similarly? | Does my gene product collaborate? |
|---|---|---|

↓ ↓ ↓

| Sequence homology | Functional genomics | Protein-protein interactions |
|---|---|---|

↓ ↓ ↓

| BLAST InterPro ProteinOrtho | Gene expression correlation | Colocalization Yeast-2-Hybrid Co-Immunoprecipitation |
|---|---|---|

# If you have transcripts, generate polypeptides, too.

5' AGGTGACACCGCAAGCCTTATATTAGC 3'

- EMBOSS `transeq` can 6-frame translate polypeptide sequence from assembled mRNAs.

- EMBOSS `checktrans` can eliminate any ORF below 100 amino acids.

- Downstream tools are faster if you translate.

# From putative transcript to a translation

```
CTGGTCGGCTATTAGAAAAGAAAGATCGGCTAAGTCCTTCGGACCTGATCAACTTAATCC
GGGAGCTACTGATTTCAACTACTTCGACCTAACTCTCTGAAAAGATGAATTACACAAGTT
TTATCTTCGCCTTTCAGCTTTGCATAATTTTGTGTTCGTCTGGCTATTACTGTCAGTCCA
TAATTTTTAGGGAAATAGAAAACCTAAGGGACTATTTTAACGCAAGTAATCCAGATGTAG
CAGATGGTGGGTCGCTTTCATAGATATTTTGAAGAATTGGAGAGAGGAGAGTGATAAAA
CAGTAATTCAAAGCCAAATTGTCTNNNNNNACTTGAAAATGTTTGAAAACTTGAAAGATA
ACCAGCTCATTCAAAGGAGCATGATACCATCAAGGAAGACATGCTTGATAAGCTGTTAA
ATAGCAGCTCTGATAAACGGAATGACTTCCTCAAGCTGACTCAAATTCCTGTAAATGATC
TGCAGGTCCAGCGCAAAGCGATAAATGAACTCTTCAAAGTGATGAATGATCTCTCACCAA
GATCTAACCTCAGAAAGAGGAAAAGGAGCCAGAATCTGTTTCAAGGCCGGAGAGCATCGA
AATAATGGTCATCCTGCCTGCAATATTTGAATTTTTTATATAAATCTATTTATTAATATT
TAATATTTTACATTATTTATATGAAGAATATATTTTTAGACTCATCAATCAAAGTATTTAT
```

GRLLEKKDRLSPSDLINLIR
ELLISTTST*LSEKMNYTSF
IFAFQLCIILCSSGYYCQSI
IFREIENLRDYFNASNPDVA
DGGSLFIDILKNWREESDKT
VIQSQIVXXXLKMFENLKDN
QLIQRSMDTIKEDMLDKLLN
SSSDKRNDFLKLTQIPVNDL
QVQRKAINELFKVMNDLSPR
SNLRKRKRSQNLFQGRRASK
*WSSCLQYLNFLYKSIY*YL
IFYIIYMKNIFLDSSIKVF

N: indeterminate base call
X: indeterminate amino acid
*: stop codon

*This reading frame produces longest ORF.*

# InterPro can be run on each sequence or in batch (Linux).



*Horizontal position is location within sequence.*

Family matches reflect hits to consensus for orthologous IFN-Gammas.
Superfamilies reflect 4-helical cytokines.
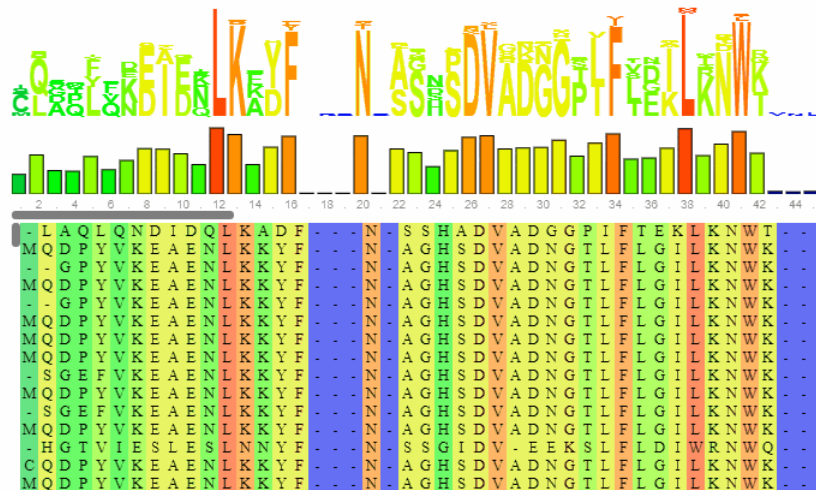Unintegrated Gene3D reflects structure availability!

https://www.ebi.ac.uk/interpro/

5

# Gene3D and CATH
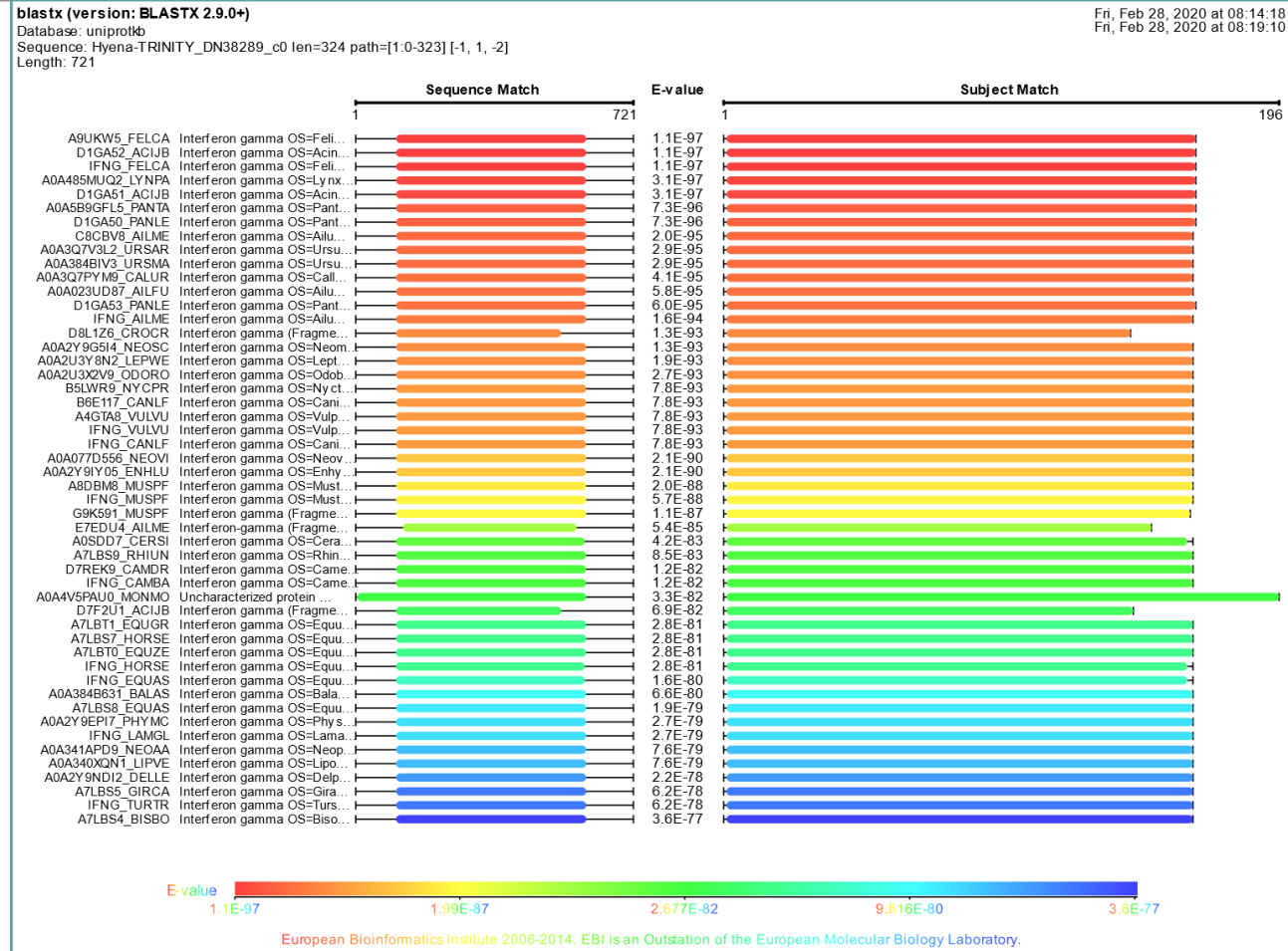


- This subsequence matches sequences with known structure.

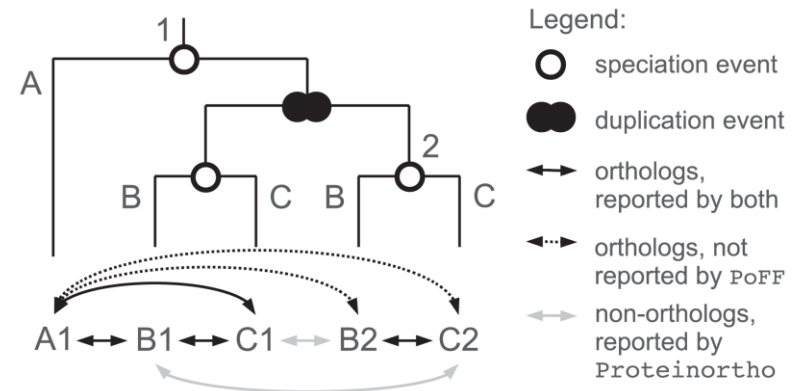- We can overlay those structures, and we can discern sequence motif commonalities.

# BLAST can be run on each sequence or in batch.

- blastx matches mRNA to protein.

- E-values show strength of match.

- Taxonomy shown in accession suffix:
  FELCA: housecat
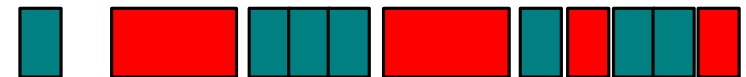  ACIJB: cheetah
  LYNPA: lynx
  PANTA: tiger



blastx (version: BLASTX 2.9.0+)
Database: uniprotkb
Sequence: Hyena-TRINITY_DN38289_c0 len=324 path=[1:0-323] [-1, 1, -2]
Length: 721

Fri, Feb 28, 2020 at 08:14:18
Fri, Feb 28, 2020 at 08:19:10

European Bioinformatics Institute 2006-2014. EBI is an Outstation of the European Molecular Biology Laboratory.

# Orthologs are the "same gene across species boundaries"

- ProteinOrtho uses speciation and gene duplication models to relate sequences.



Legend:
○ speciation event
● duplication event
↔ orthologs, reported by both
←--→ orthologs, not reported by PoFF
↔ non-orthologs, reported by Proteinortho

- Diamond indexing accelerates all protein vs. all protein homology detection.

**PFERPEAEAMCTSFKENPT**

**RLVRPEVDVMCTAFHDNEE**

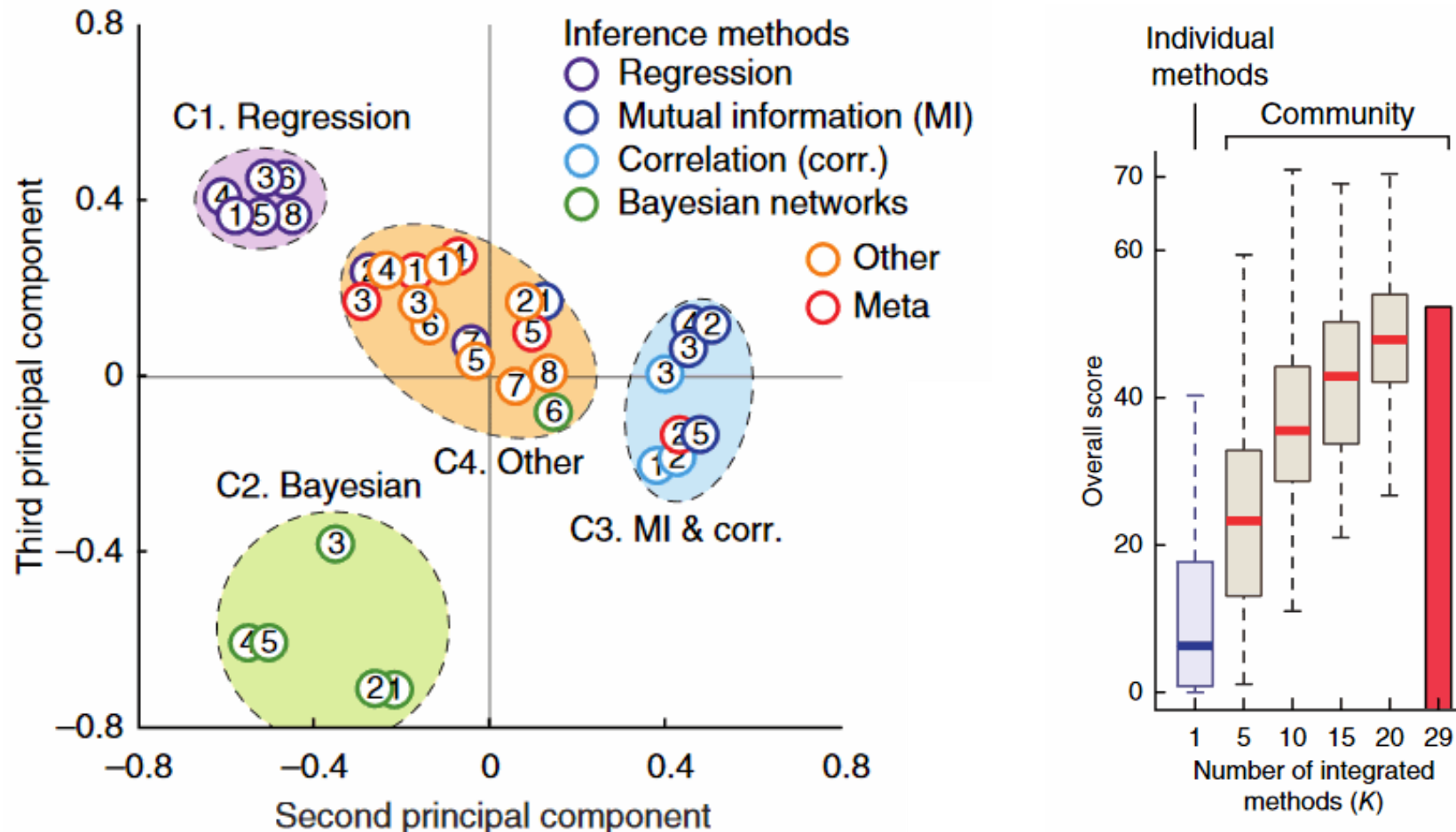M. Lechner et al. *PLOS One* (2014) 9: e105015
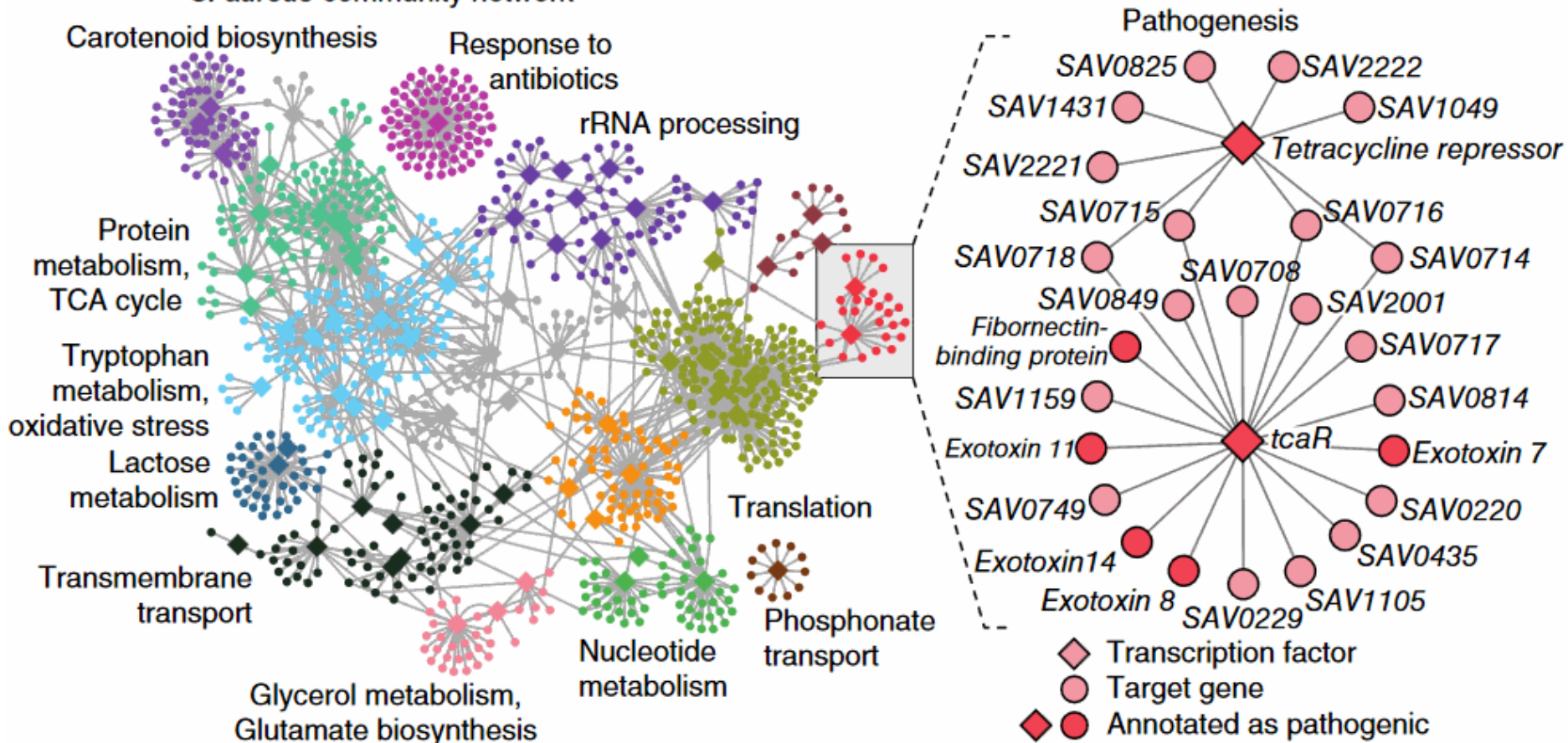
# Infer function by coexpression



"Gene-expression data can be used to define groups of genes that show similar patterns of expression, or co-variation, across multiple conditions."

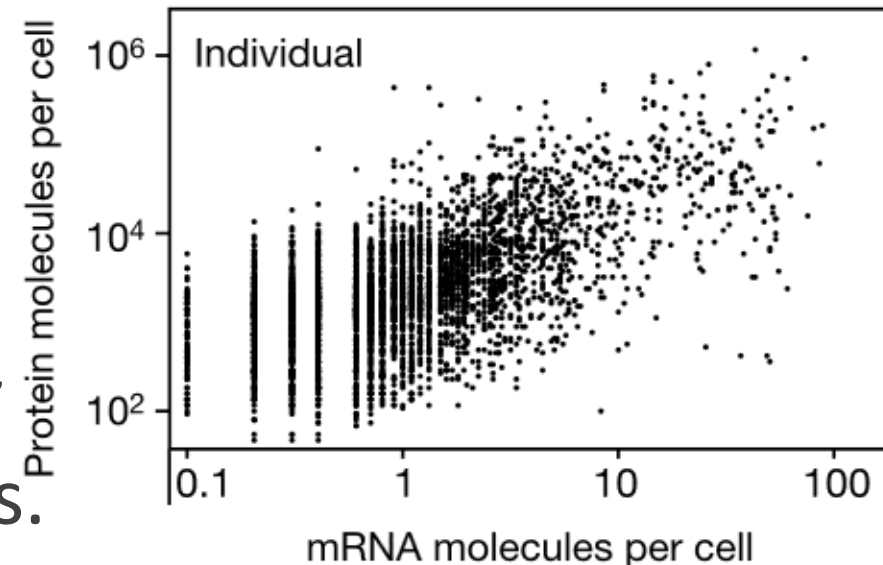# Correlation of gene expression benefits from multiple methods

# Networks of hierarchical, modular relationships
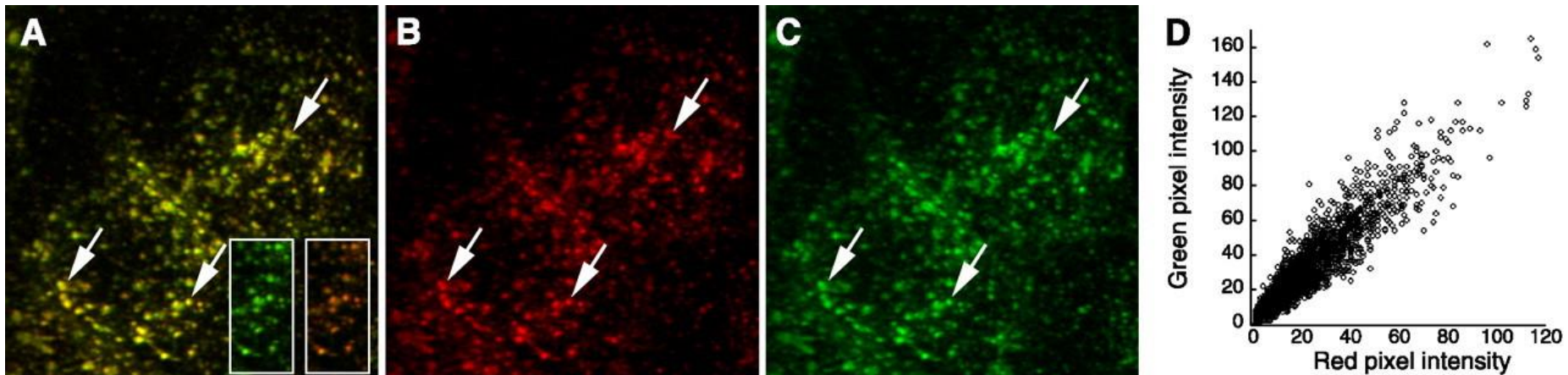
# mRNA and protein quant give different results

- Proteins in *S. cerevisiae* ranged from 50 to more than a million molecules per cell.

- Western blots quantified proteins, while microarrays quantified messenger RNA.

- Spearman rank corr. yields r=0.57. *Turnover* uses separate processes.

# Colocalization: aspects of spatial association

- **Co-occur** overlaps in space

- **Correlate** scales similarly across structures

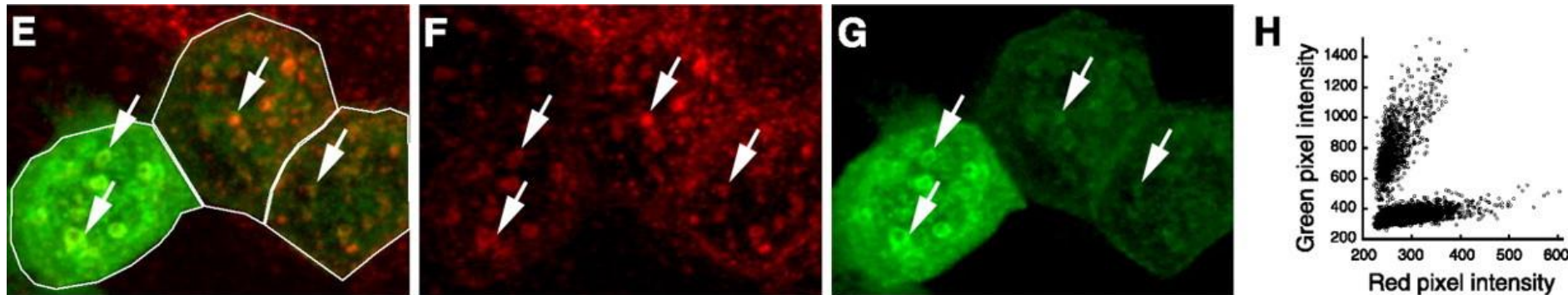- **Interact** requires FRET or EM resolution

*Endosomes in MDCK cells: Pearson CC= 0.944*



KW Dunn et al. *Am. J. Physiol. Cell Physiol*. (2011) 300: C723-C742

# Regions of Interest and single cell analysis

- Each cell may reveal different relationship.

- *Segmenting* cells from background boosts signal-to-noise.

- Specifying ROIs can separate cellular info.
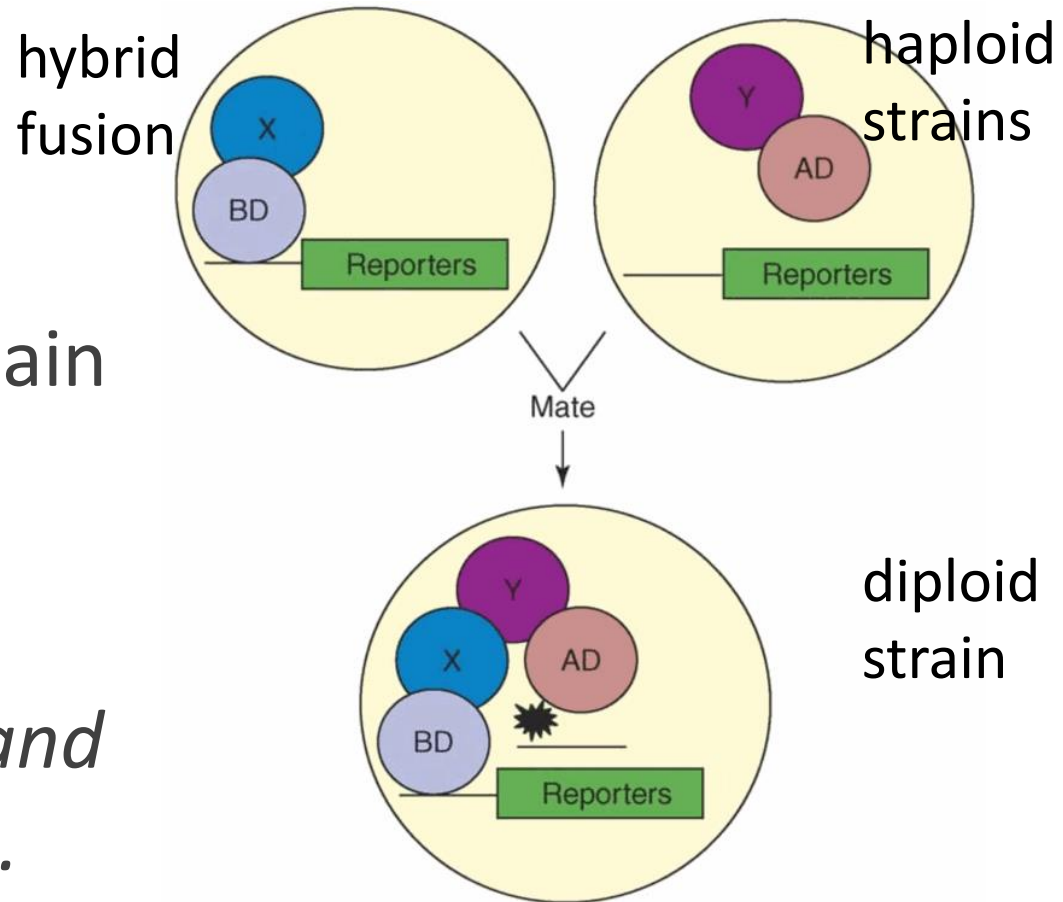
# Defining protein-protein interactions (PPIs)

"The physical contact considered in PPIs should be specific, not just all proteins that bump into each other by chance. It also should exclude interactions that a protein experiences when it is being made, folded, quality checked, or degraded."

# Yeast two-hybrid establishes binary relationship

- X: first gene

- Y: second gene

- BD: DNA-binding domain

- AD: Transcription activating domain

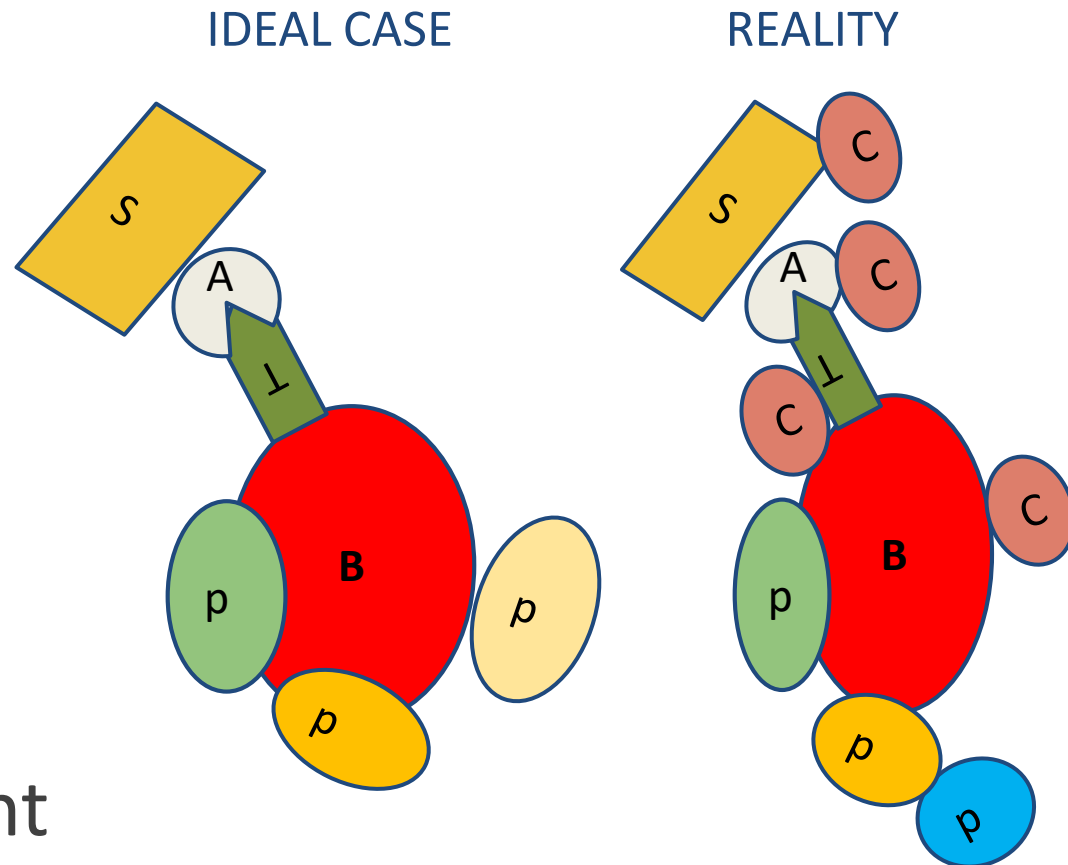*If X and Y interact, BD and AD cooperate to report.*



hybrid fusion

haploid strains
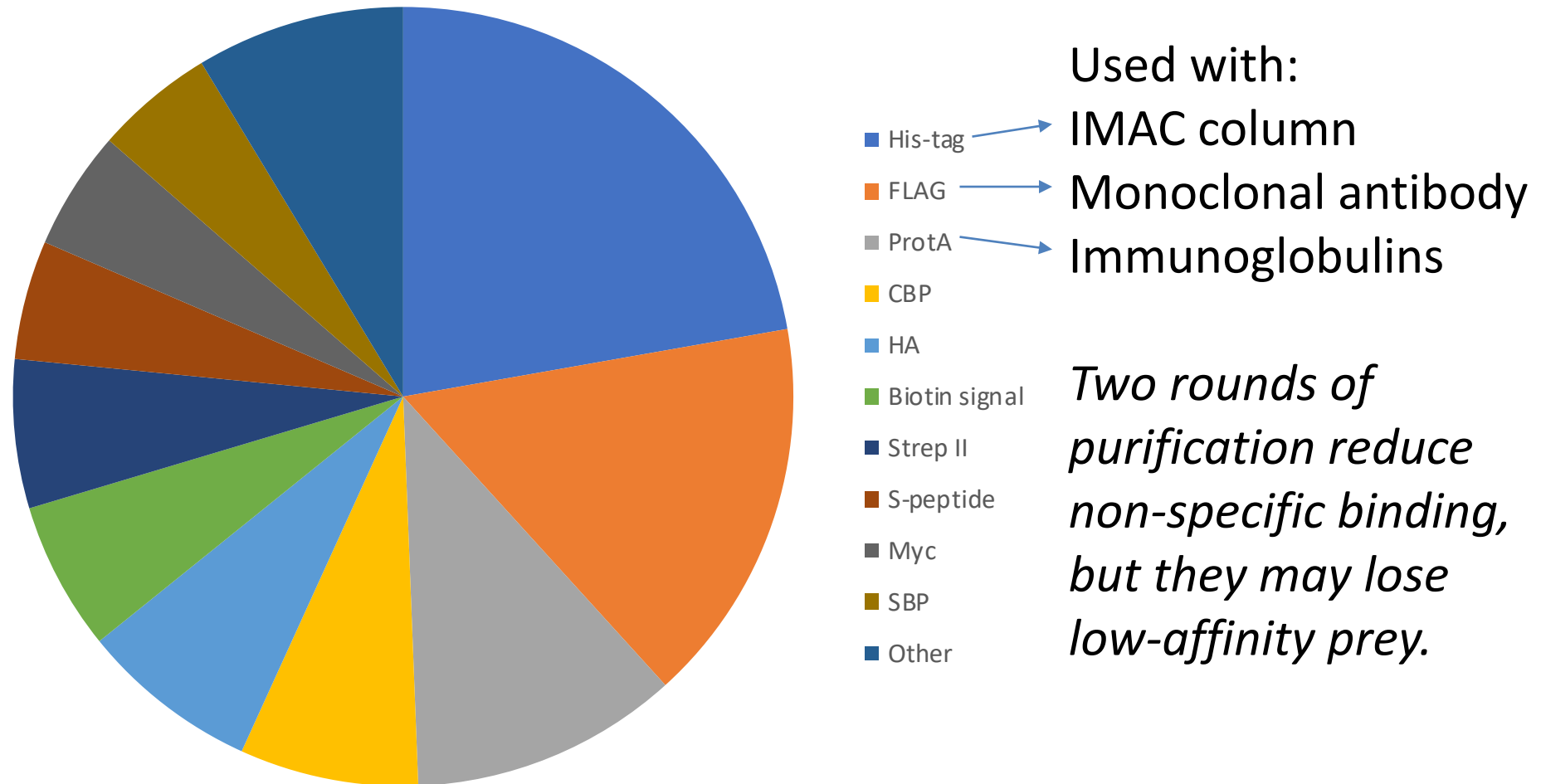
diploid strain

Current Opinion in Biotechnology

JR Parrish et al. *Curr. Opinion Biotech*. (2006) 17: 387-393.

# Co-immunoprecipitation probes multiple interactions

- Support
- Antibody
- Tag
- Bait
- Prey
- Contaminant

IDEAL CASE

REALITY

# Many possible combinations for tandem affinity purification



Legend:
- His-tag → Used with: IMAC column
- FLAG → Monoclonal antibody
- ProtA → Immunoglobulins
- CBP
- HA
- Biotin signal
- Strep II
- S-peptide
- Myc
- SBP
- Other

Used with:
IMAC column
Monoclonal antibody
Immunoglobulins

*Two rounds of purification reduce non-specific binding, but they may lose low-affinity prey.*
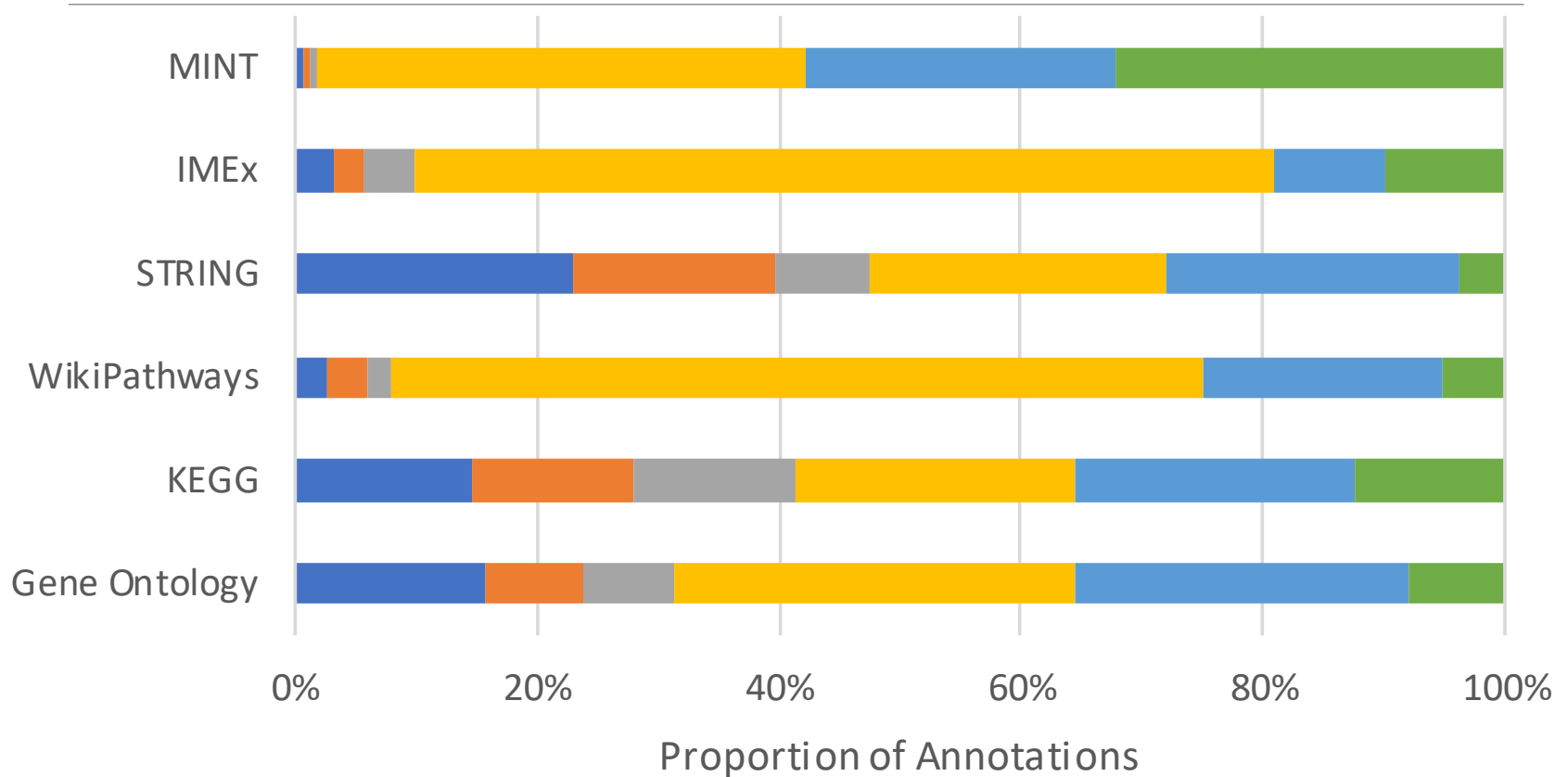
Yifeng Li. Biotechnol. *Appl. Biochem*. (2010) 55: 73-83.

# When will structural data enable protein-protein inference?

- "Localized regions on protein surfaces [are] conserved among structural neighbors that participate in protein-protein interactions."

- "As long as structural information is available for a given pair of proteins... the set of 'template complexes' available in the current structural databases can be used to generate coarse-grained models of protein-protein interactions."

# Pathway and network annotation vary considerably



**Proportion of Annotations**

Legend: ■ A. thaliana ■ C. elegans ■ D. melanogaster ■ H. sapiens ■ M. musculus ■ S. cerevisiae

# Takeaway Messages

▪Learning a transcript sequence allows you to leverage considerable sequence resources.

▪Expression data may reveal co-expression partners besides responsiveness to stimuli.

▪Colocalization and protein-protein interactions require experimentation.

▪Pathway and network data skew by organism.