# The function of proteins

* Structural: the organelles of the cell

* Signaling: pass information from the environment and between different parts of the cell; turn genes on & off.
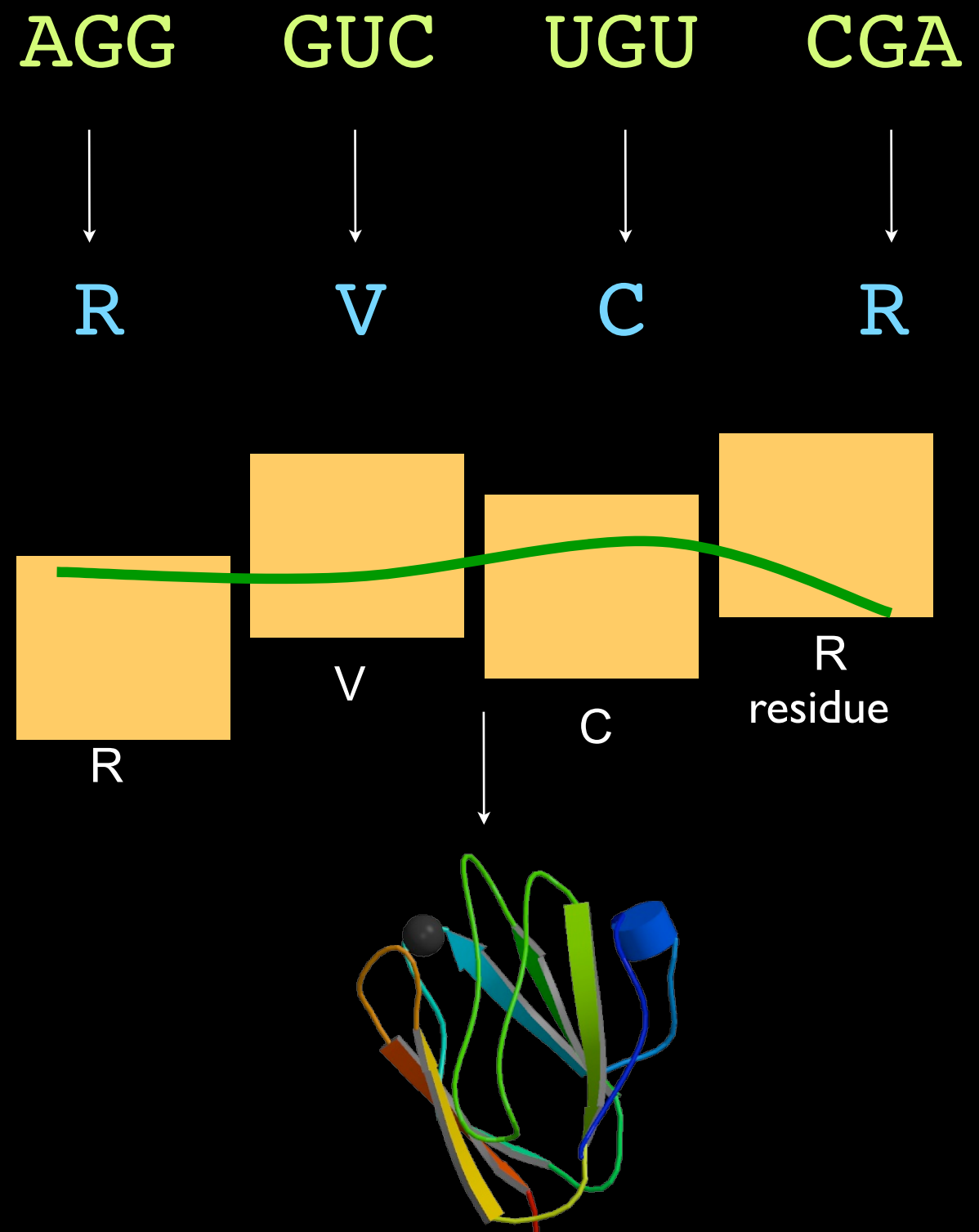
* Catalyze reactions (act as *enzymes*).

# Proteins

mRNA
$\Sigma = \{A,C,G,U\}$

protein
$|\Sigma| = 20$ amino acids

Amino acids with flexible side chains strung together on a backbone

**Proteins are the Building Blocks of Life**

Their shape is instrumental in determining their function.

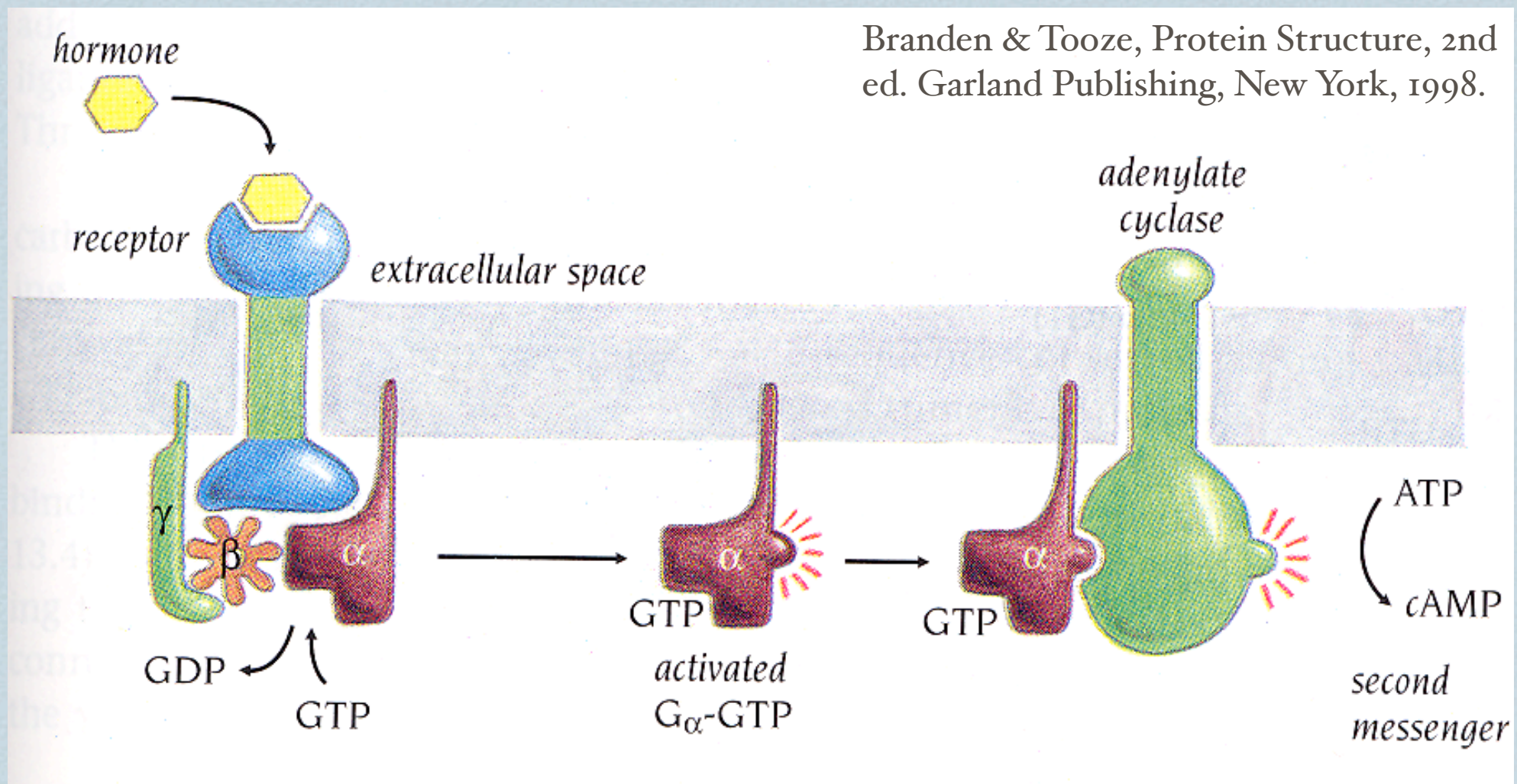AGG    GUC    UGU    CGA

R        V        C        R

R        V        C        R residue

- **Central dogma: DNA → mRNA → Proteins**
- **Proteins are building blocks of many cellular processes**
- **Conservation ⇒ functional importance**

- **Whole-genome (noisy) protein-protein interaction networks and other networks becoming available:**
  - **function annotation**
  - **combining graphs, assigning confidence, predicting edges, eliminating noise**
  - **comparing, searching graphs**
  - **figuring out how they evolved**

- **Start with experimental techniques for generating the graphs; then move on to network clustering.**

- **Central dogma: DNA → mRNA → Proteins**
- **Proteins are building blocks of many cellular processes**

- **Networks:**

| Network | Nodes | Edges |
|---|---|---|
| Transcription (aka regulatory) | proteins/genes | A "regulates" B |
| Metabolic | Metabolites / small molecules | Reactions |
| **Protein-Protein** | **Proteins** | **Physical Interactions** |

# Experimental Techniques

## CMSC 858L

# Proteins Interact



Branden & Tooze, Protein Structure, 2nd ed. Garland Publishing, New York, 1998.

"Why" proteins interact:

Bring chains of enzymes together

Signal Transduction

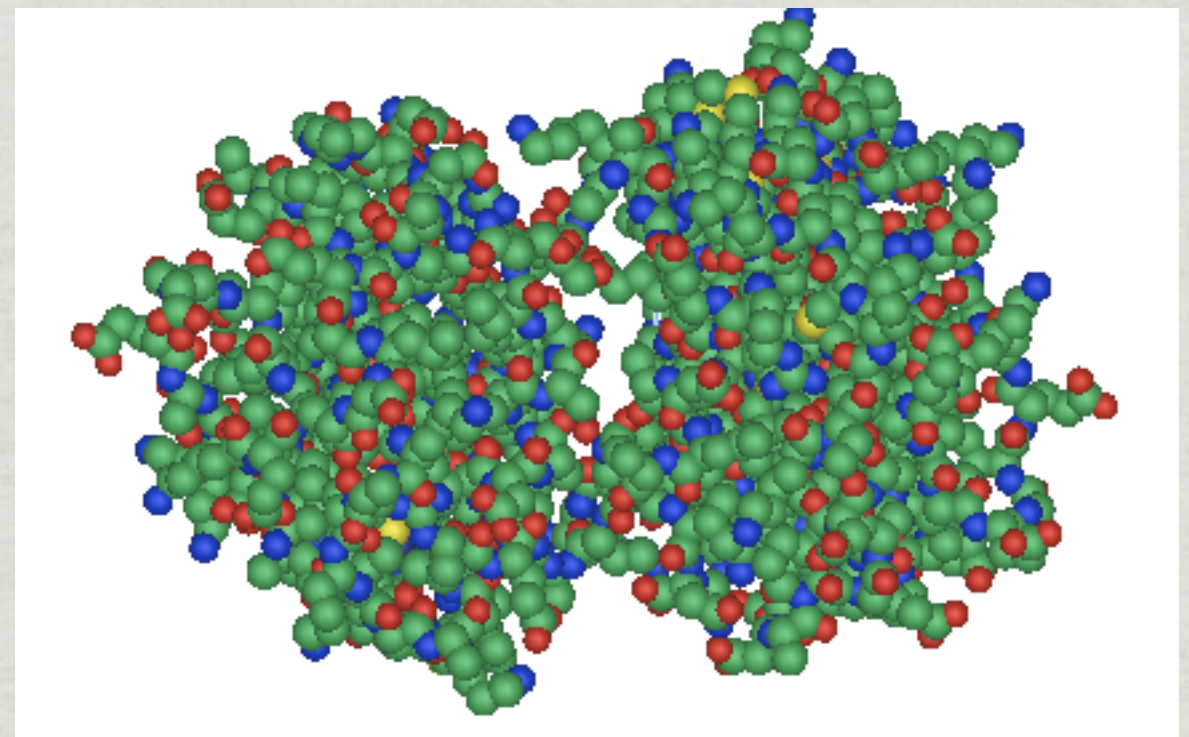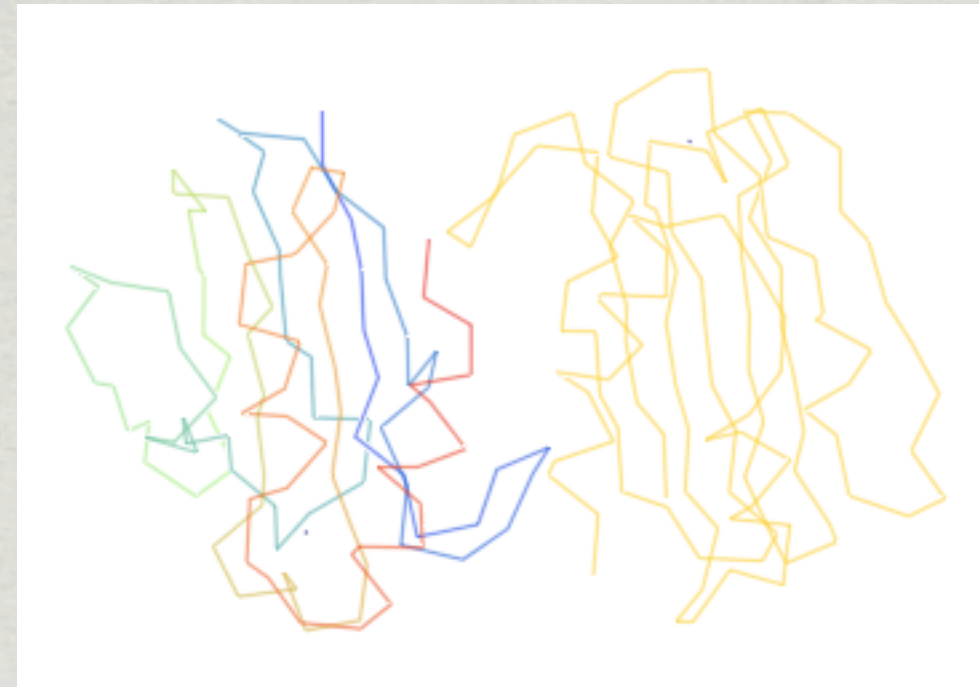"Tethered" Signal Transduction

Form structures

From "Analysis of Biological Networks" Junker and Schreiber, eds

# Experimental Techniques to Determine Protein Interactions

✳ Slow, accurate, costly:

  ✳ X-ray crystallography

  ✳ NMR

✳ High throughput, but noisy:

  ✳ Yeast Two-Hybrid

  ✳ TAP-MS (tandem affinity purification / mass spec)

# Determining protein structure:

* X-ray crystallography

* NMR

  * If you **can** determine structure of a complex, you know the position of each of its atoms.

  * Slow, costly techniques.

  * Don't always work.

* More recent: high-throughput techniques

# PDB

# RCSB PDB
PROTEIN DATA BANK

## An Information Portal to Biological Macromolecular Structures

As of **Tuesday Sep 01, 2009** 🔊 there are 59939 Structures ❓ | PDB Statistics ❓

**WHAT'S NEW** | **HELP** | **PRINT**

PDB ID or keyword ⏷     [ Search ] | **Advanced Search**

| Summary | Derived Data | Sequence | Seq. Similarity | Literature | Biol. & Chem. | Methods | Geometry | Links |

**Home**
- News & Publications
- Policies
- FAQ
- Contact
- Feedback
- About Us

**Deposition**
- All Deposit Services
- Electron Microscopy
- NMR
- Validation Server
- BioSync Beamline
- Related Tools

**Search**
- Advanced Search
- Latest Release
- Latest Publications
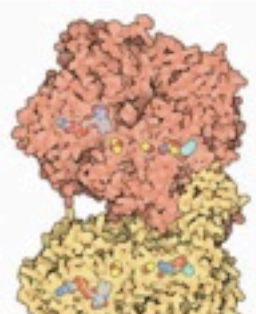- Sequence Search
- Ligand Search
- Unreleased Entries
- Browse Database
- Histograms

**Explorer:**

**Last Structure:** 3HTO

**Tools**
- File Downloads
- File Formats
- Services: RESTful | SOAP
- Widgets
- Compare Structures

## the hemagglutinin structure of an avian H1N1 influenza A virus

# 3hto

📄 **Display Files ▼**
📥 **Download Files ▼**
📄 **Print this Page**
➕ **Share this Page**

DOI:**10.2210/pdb3hto/pdb**

### Primary Citation

### Molecular Description   **Hide**

Classification:   **Viral Protein**🔍
Structure Weight: 55101.73

Molecule: Hemagglutinin HA1 chain
Polymer: 1   Type: polypeptide(L)   Length: 324
Chains: A

Molecule: Hemagglutinin HA2 chain
Polymer: 2   Type: polypeptide(L)   Length: 160
Chains: B

### Source   **Hide**

Polymer: 1

Scientific Name:   **Influenza a virus**🔍

Polymer: 2

Scientific Name:   **Influenza a virus**🔍

Polymer: 3

### ◄ Biological Molecule 1 ❓ ►

More Images...

☒ **View in Jmol**   SimpleViewer ❓
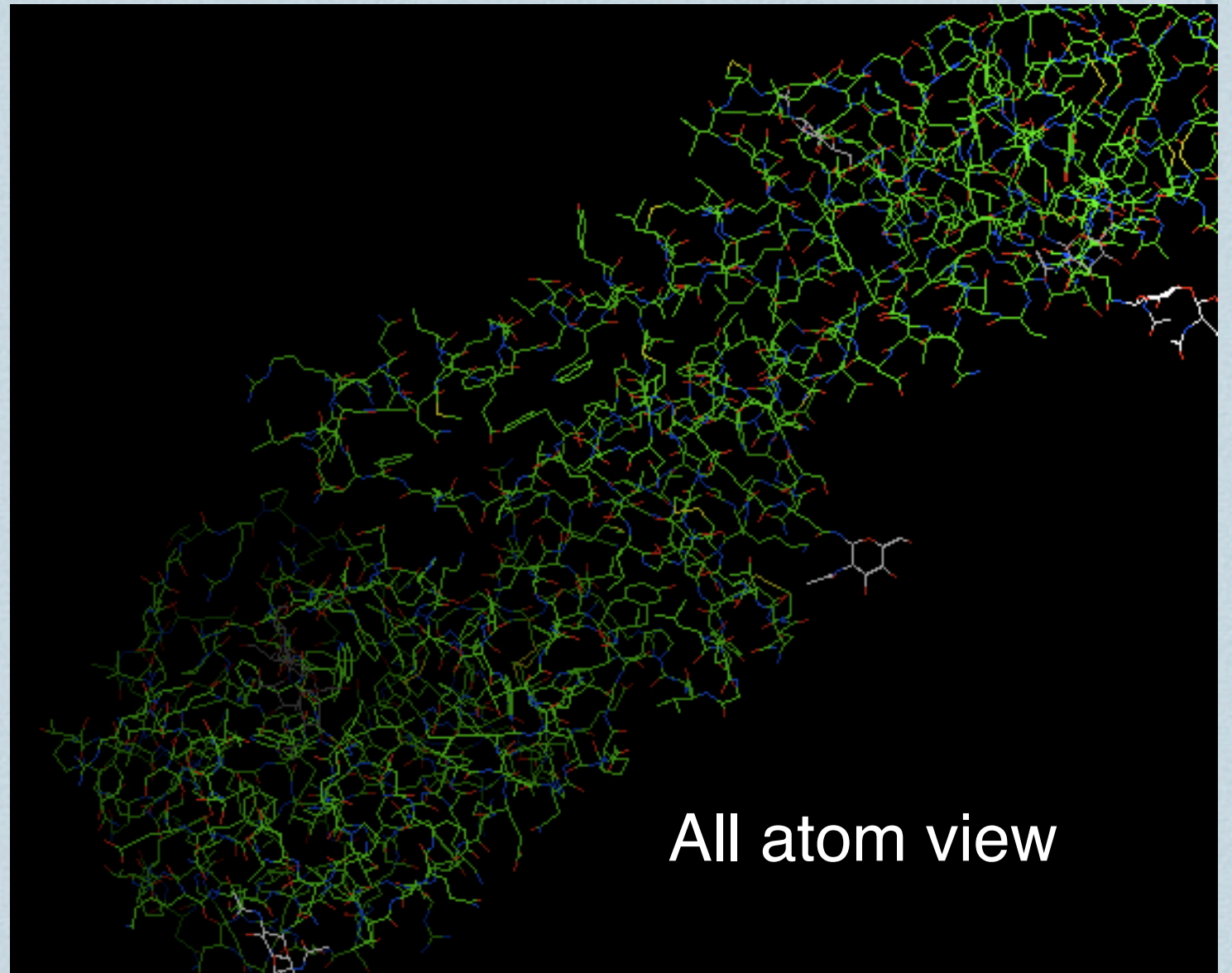Protein Workshop
Other Viewers ▼

Oligomeric State: TRIMERIC
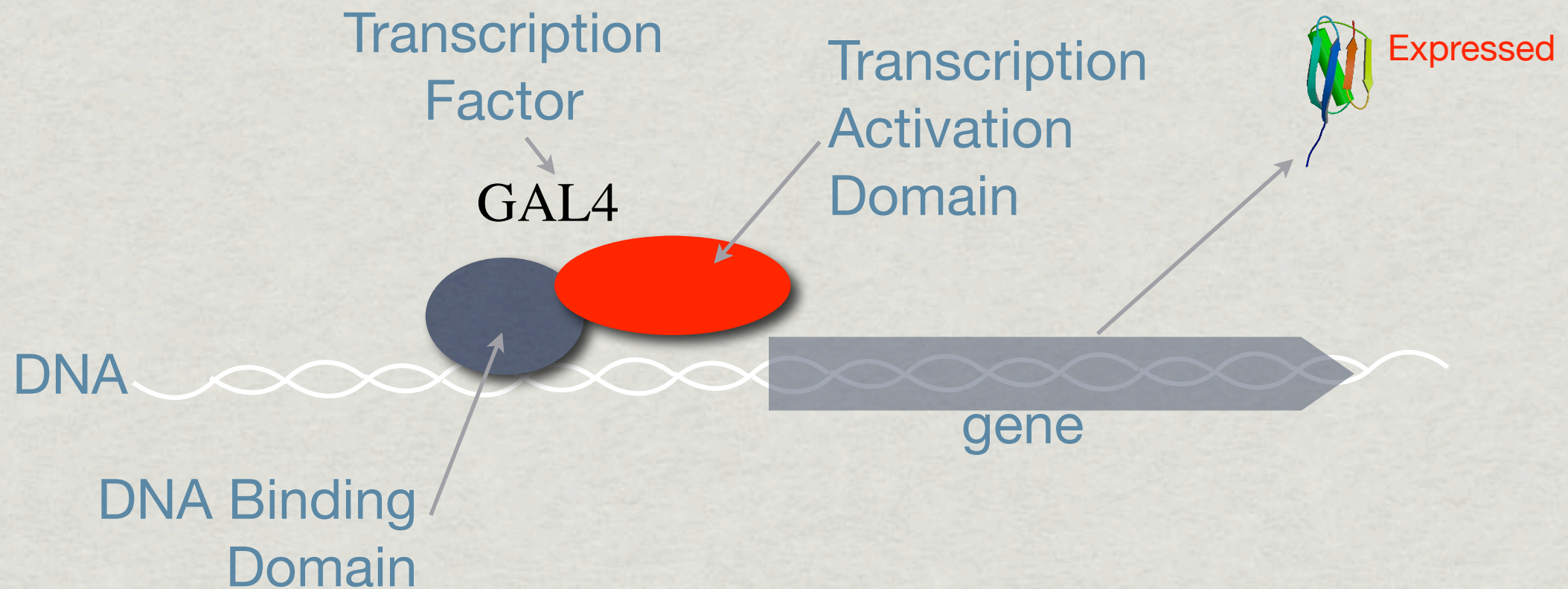
# Other View of a Protein

AVIAN H5 HAEMAGGLUTININ



"Cartoon" drawing, showing major features such as alpha helices and beta sheets
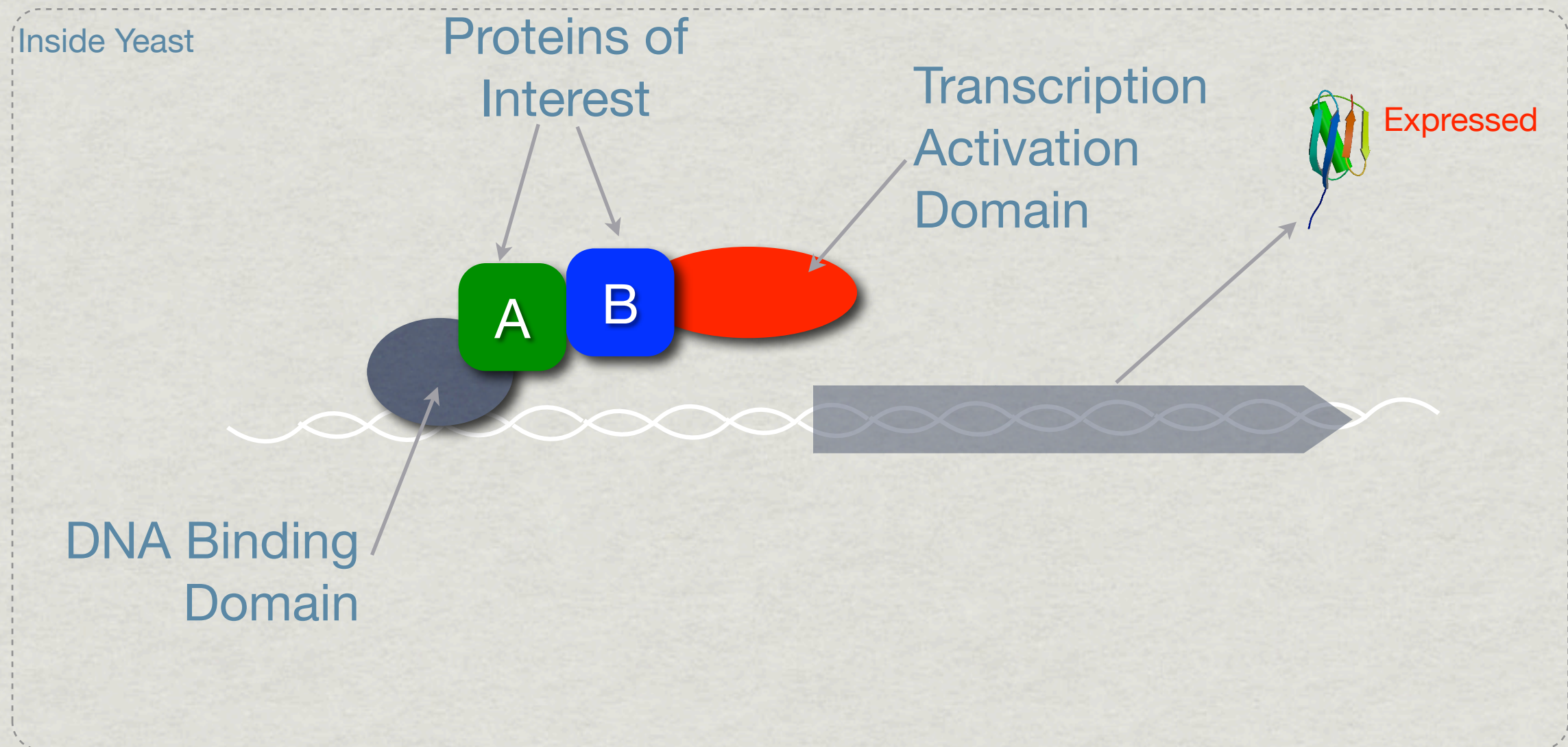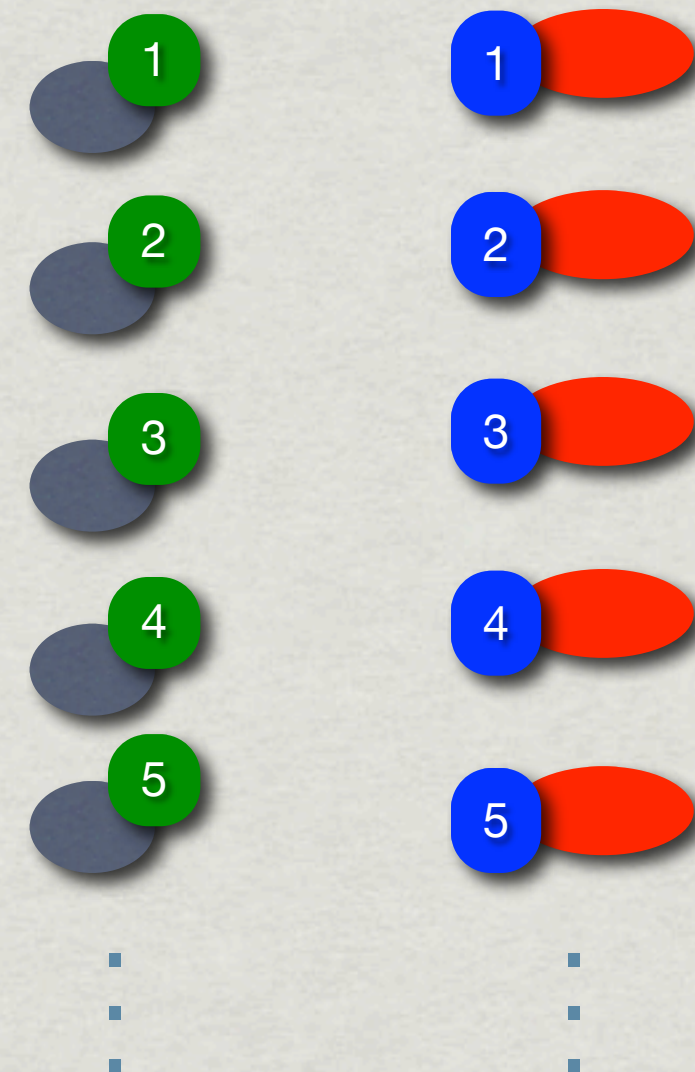


All atom view

# Yeast Two-Hybrid



"Domain" = functional, evolutionary conserved unit of a protein

# Scaling Up (Ito et al, 2001)
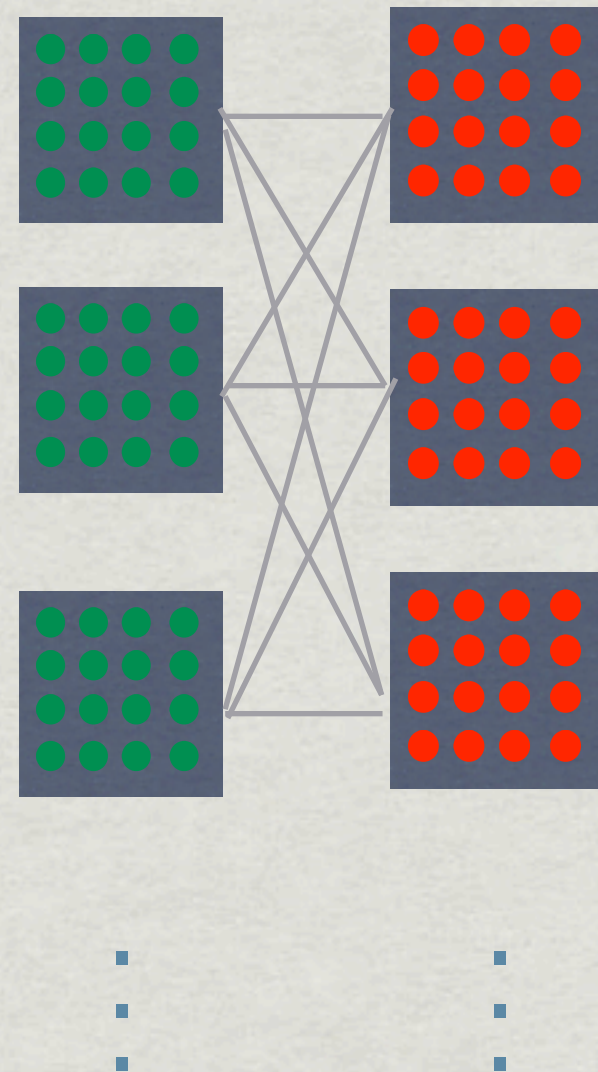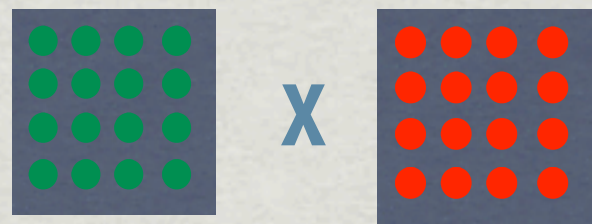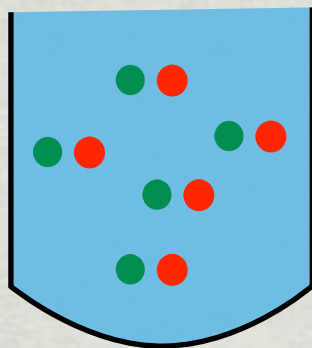


**96-well plates**
Each well contains a yeast strain with a different hybrid

**BAIT**

**PREY**

~ 6,000 genes / 96
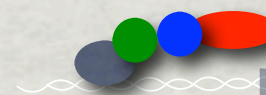= 62 plates
= 3,844 crosses between plates

# Ito et al, 2001 Results:

~ 18 million gene pairs; (prey,bait) & (bait, prey)

# of colonies that passed all 4 tests

Involve 3278 proteins out of ~ 6,000

**Table 1. Summary of the comprehensive two-hybrid screening**

| | |
|---|---|
| Mating reactions | 3,844 |
| Combinations to be examined | ~$3.5 \times 10^7$ |
| Positive colonies | 15,523 |
| ISTs  INTERACTION SEQUENCE TAGS | 13,754 |
| Independent two-hybrid interactions | 4,549 |
|    More than 2 IST hits | 1,533 |
|    More than 3 IST hits (core data) | 841 |

**Table 2. Comparison between the two genomewide IST projects**

| Dataset | Total interactions | Known interactions* (%) |
|---|---|---|
| Uetz et al. (11) | 691† | 88 (12.7) |
| This study | | |
|    More than 2 IST hits | 1,533 | 128 (8.3) |
|    More than 3 IST hits (core data) | 841 | 105 (12.5) |

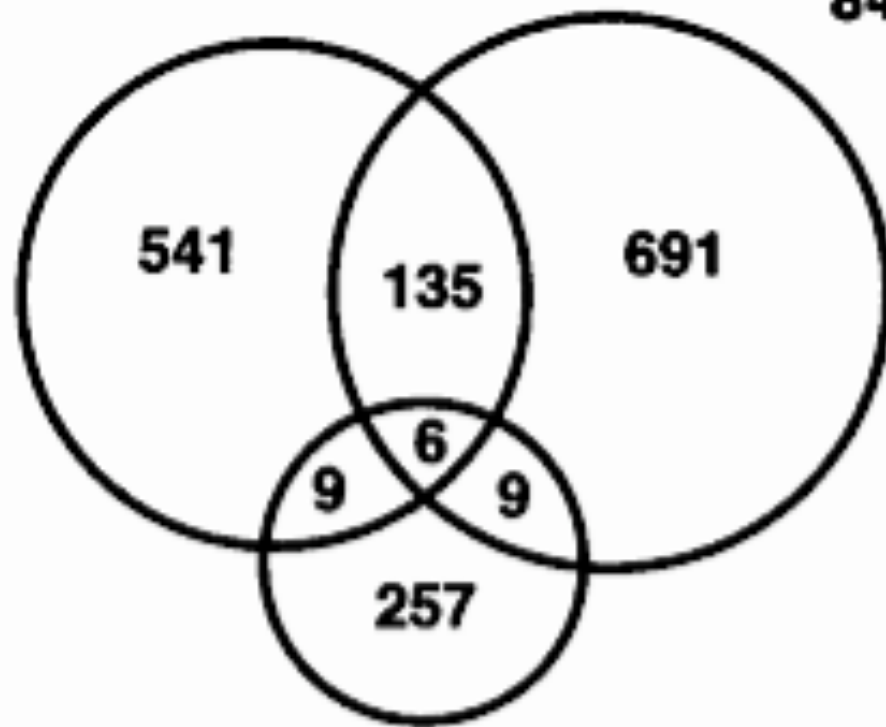*Those described in the YPD (14) as previously known to associate or to occur in the same complex.
†In Uetz et al. (11), total number of interactions revealed by IST approach was claimed to be 692, whereas their list contained 691 interactions.

Comparable overlaps between known interactions

# What could go wrong with Yeast 2-Hybrid?

Uetz *et al.*
IST approach
691

this study
core data
841

541   135   691

6

9   9

257

(Ito et al, 2001)

Uetz *et al.*
protein array
281

Low overlap!
Why?
- different experimental protocols

- different ways of making the hybrid genes (some fold correctly in Ito et al, but not in Uetz et al)

- actual randomness in binding

- Test takes place in nucleus, so proteins that never enter the nucleus won't be tested.

- Will always be using yeast: so required post-translational modifications might not happen.

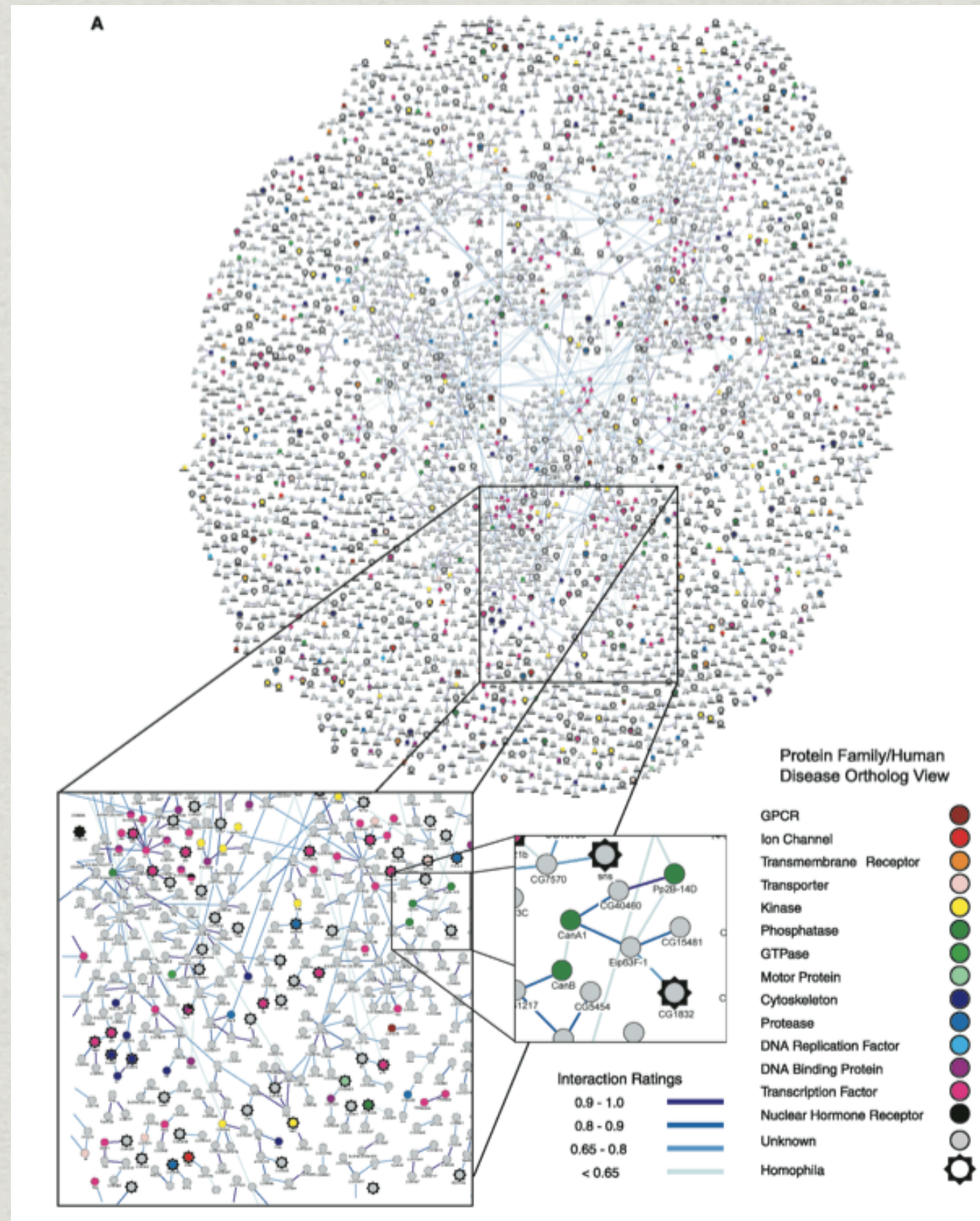- Triple interactions: A - X - B

- Both proteins may not meet in vivo.

- Transcription factors can be hard to test (b/c they may activate the reporter gene w/o binding)

- hydrophobic / membrane proteins may not fold correctly.
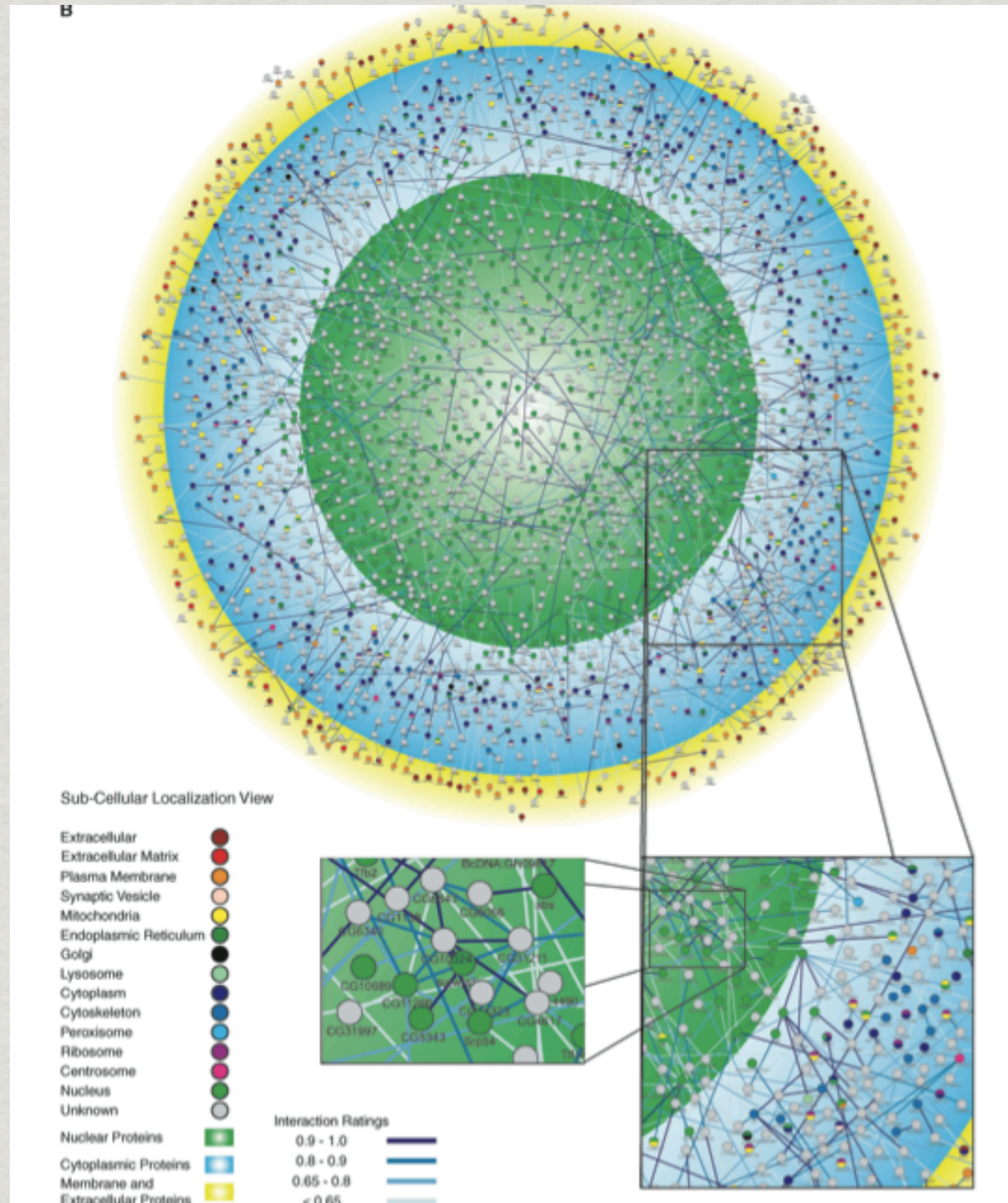
# Fruit Fly (*Drosophila melanogaster*)

# Fly (*Drosophila melanogaster*) (Giot et al, 2003)

7,048 proteins 20,405 interactions

High-confidence: 4,679 proteins 4,780 interactions

Colored and placed by sub-cellular location

(Giot et al, 2003)

## Sub-Cellular Localization View

Extracellular
Extracellular Matrix
Plasma Membrane
Synaptic Vesicle
Mitochondria
Endoplasmic Reticulum
Golgi
Lysosome
Cytoplasm
Cytoskeleton
Peroxisome
Ribosome
Centrosome
Nucleus
Unknown

Nuclear Proteins

Cytoplasmic Proteins

Membrane and
Extracellular Proteins

Unknown cellular location

Interaction Ratings

0.9 - 1.0
0.8 - 0.9
0.65 - 0.8
< 0.65

4 of the 6 highly connected proteins in fact are predicted by other means to be in the nucleus.

# General Topological Properties



**A** — Degree — $k_i$ = number of links connected to node $i$

**B** — Distance — $d_{ij}$ = shortest path length between node $i$ and $j$

**C** — Diameter — $D = \max\{d_{ij} \mid i, j \in N\}$    $N$ : all nodes in the network

**D** — Clustering Coefficient — $c_i = \dfrac{2e_i}{k_i(k_i-1)}$    $e_i$ : number of existing links (labeled in red) among the $k_i$ nodes that connect to node $i$

**E** — Betweenness — $b_l = \sum_{ij} p_{ij}(l) / p_{ij}$    $p_{ij}$ : number of shortest paths between $i$ and $j$
   $p_{ij}(l)$ : number of shortest paths between $i$ and $j$ going through node $l$

Zhu et al, 2007.

**Characteristic:** Huge # of low-degree nodes, with a few very high-degree nodes.

Giant connected
component
    = 3659 edges
    = 3039 nodes

Avg shortest path =
9.4 links, longer than
expected in a
random network (7.7
links)

Node disjoint loops

Real network contains more triangles than random



**B**

Legend:
- Actual
- Random
- 1-level fit
- 2-level fit

Y-axis: Count (1 e+00, 1 e+02, 1 e+04, 1 e+06)

X-axis: Loop perimeter (4, 6, 8, 10, 12, 14)

*C. elegans*

# Li et al, 2004 Results:

| Interaction Set | # interactions | # proteins |
|---|---|---|
| Core | 2157 | 502 baits 1039 preys |
| Non-core | 1892 | 531 baits 1395 preys |



- "Core" means they observed an interaction ≥ 3 times.
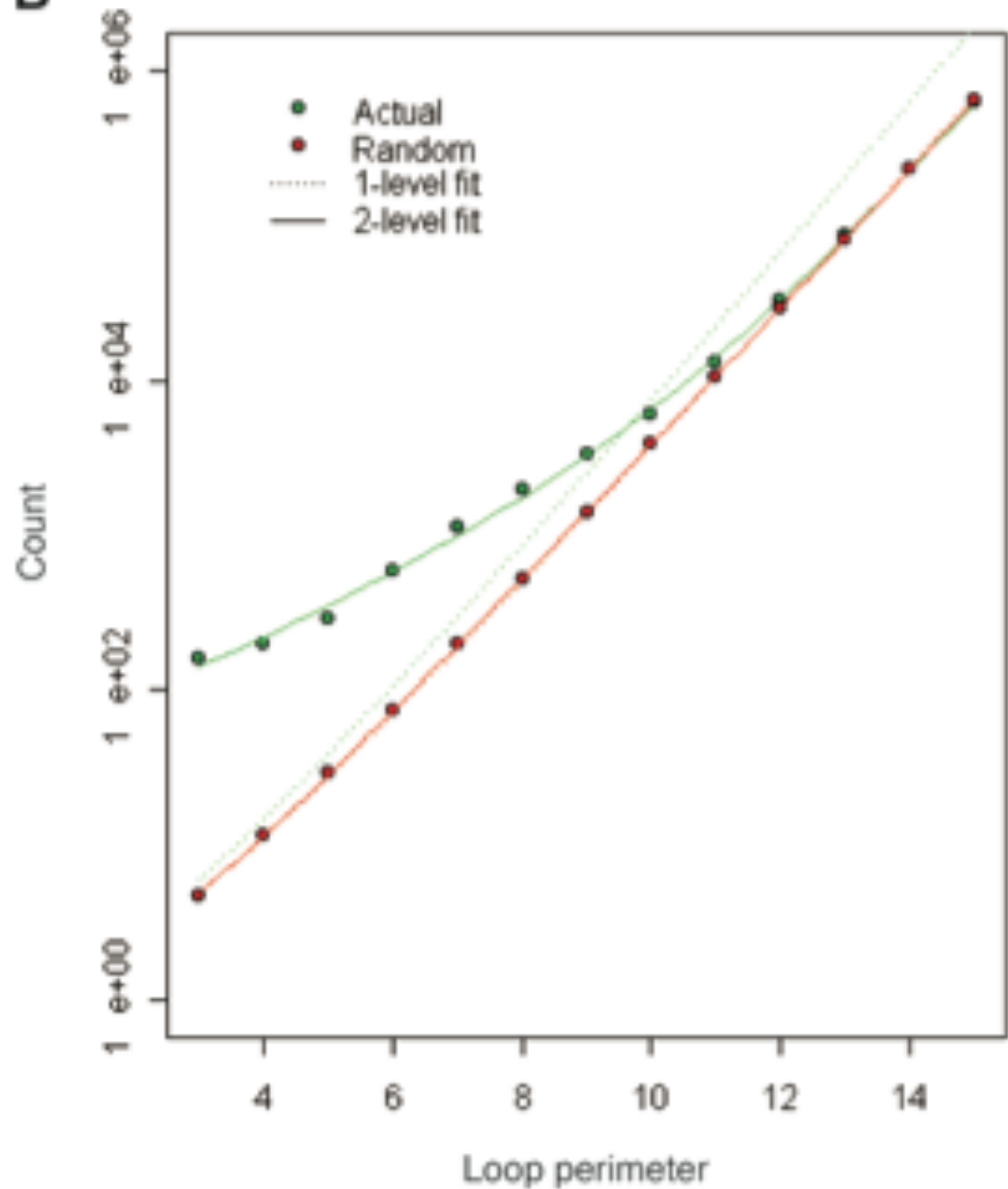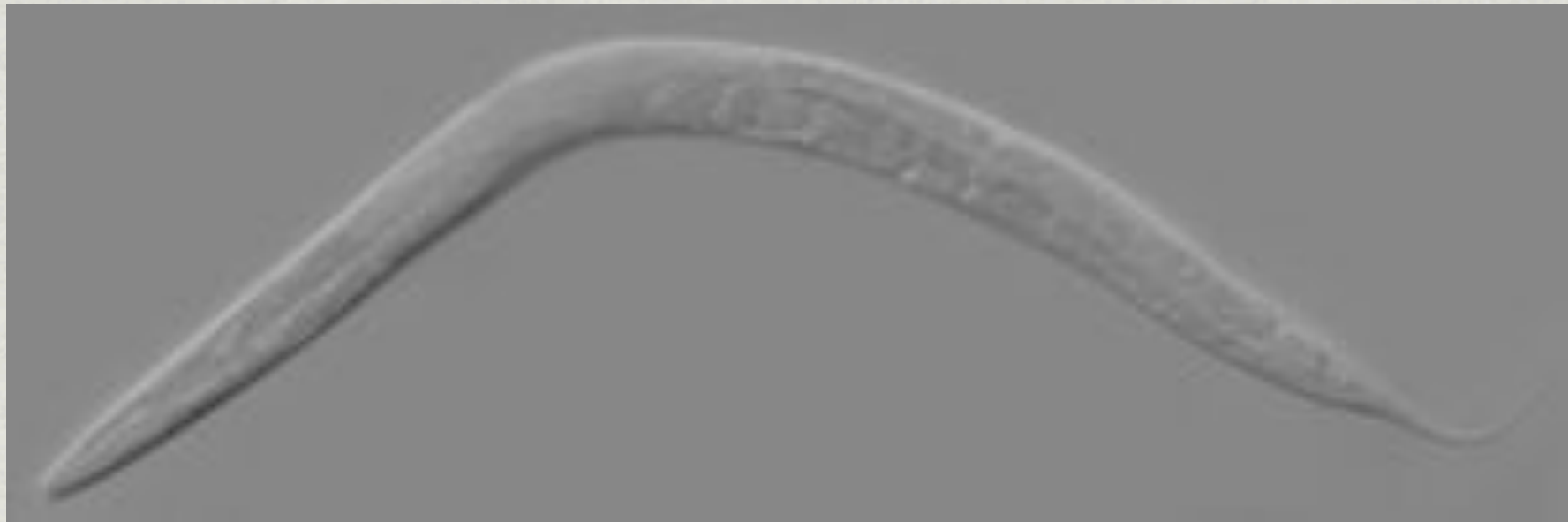
- Asymmetric (bait, prey) vs (prey, bait)

- Out of 2157 core pairs, only **22** were observed in both orientations

- 108 interactions in WormPD involved the tested proteins

- Core contain 8 of these interactions; Non-core contained 2

- Coverage = (8+2)/108 ≈ 10%

# Estimating Reliability

Sprinzak et al, 2003

fraction co-localized
pairs in true edges

fraction of true edges
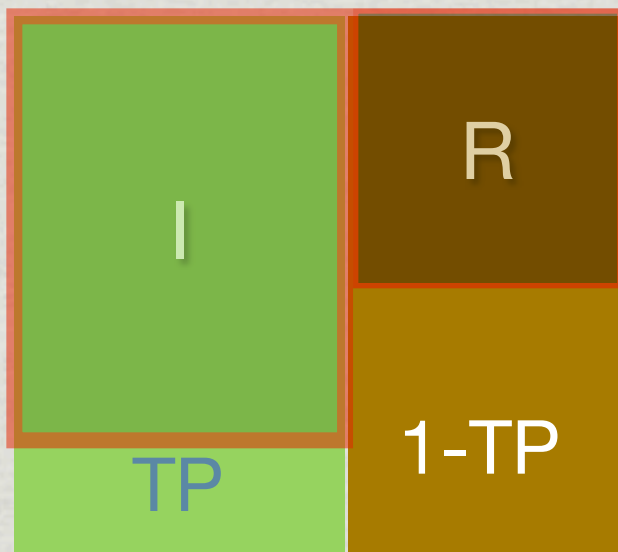
fraction of co-
localized pairs in
random edges

$$D = TP \times I + (1-TP) \times R$$

fraction of **co-localized**
pairs in predicted edges

fraction of false
positives

Assuming FP are
'random-like'

$$TP = (D - R) / (I - R)$$

## Estimating TP

$$TP = (D - R) / (I - R)$$

✳ D = fraction co-localized predicted edges

✳ R ≈ fraction of **all** pairs that are co-localized (~0.36)

✳ I ≈ 1 or 0.95   [assumed to be very high]

**Table 1.** Data sets of pairs of interacting proteins

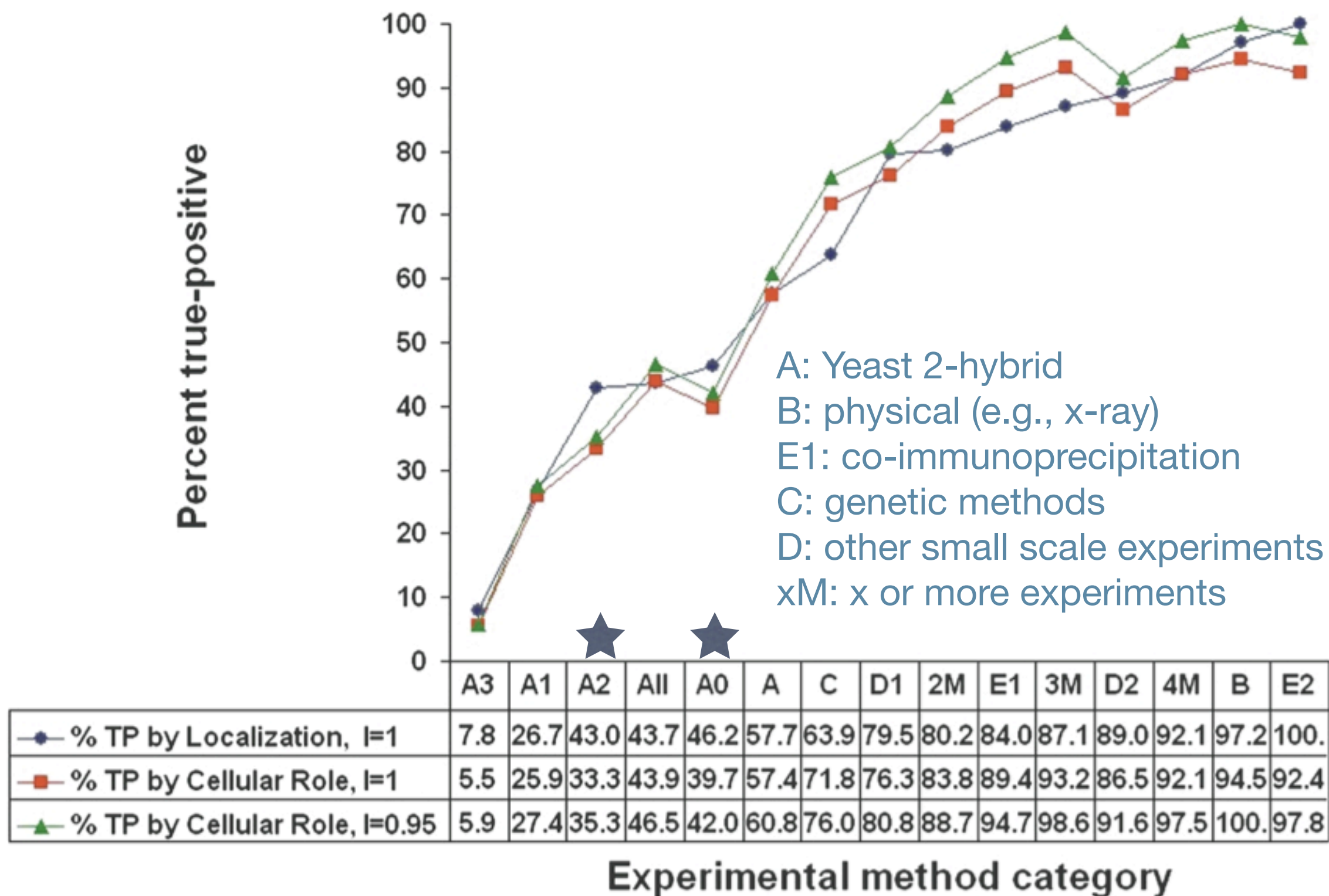| Experimental method category[a] | Number of interacting pairs | Co-localization[b] (%) | Co-cellular-role[b] (%) |
|---|---|---|---|
| All: All methods | 9347 | 64 | 49 |
| A: Small scale Y2H | 1861 | 73 | 62 |
| A0: GY2H Uetz *et al.* (published results) | 956 | 66 | 45 |
| A1: GY2H Uetz *et al.* (unpublished results) | 516 | 53 | 33 |
| A2: GY2H Ito *et al.* (core) | 798 | 64 | 40 |
| A3: GY2H Ito *et al.* (all) | 3655 | 41 | 15 |
| B: Physical methods | 71 | 98 | 95 |
| C: Genetic methods | 1052 | 77 | 75 |
| D1: Biochemical, *in vitro* | 614 | 87 | 79 |
| D2: Biochemical, chromatography | 648 | 93 | 88 |
| E1: Immunological, direct | 1025 | 90 | 90 |
| E2: Immunological, indirect | 34 | 100 | 93 |
| 2M: Two different methods | 2360 | 87 | 85 |
| 3M: Three different methods | 1212 | 92 | 94 |
| 4M: Four different methods | 570 | 95 | 93 |

**Figure 1.** True positive rates in various data sets that are distinguished by the experimental method for determining protein–protein interaction. Blue, percentage of TP based on co-localization ($I = 1$); red, green, percentage of TP based on shared cellular-role for $I = 1$ (red) and $I = 0.95$ (green). $I$ is the fraction of pairs with co-localized pair-mates in true interacting pairs (see the text). For the method categories see the legend to Table 1.

- High-throughput interaction detection
- Yeast two-hybrid - pairwise
  - organisms as machines to learn about organisms
  - yeast, worm, fly, human,...
  - low intersection between repeated experiments
  - *in vivo*, but takes place inside the nucleus.
  - Estimated 50% FP rate
  - statistics: shortest path distribution, degree distribution, # triangles, etc. show that Y2H graphs ≠ random graphs.

- TAP-MS (co-immunoprecipitation) - complexes: Simultaneous interactions between several proteins.

# PubMed: Where nearly all the relevant papers can be found.

Relevant Journals: Science, Nature, PLoS Biology, PLoS Comp. Biology, Bioinformatics, BMC Bioinformatics, Journal Computational Biology, Genome Biology

http://www.ncbi.nlm.nih.gov/PubMed/

| Amino Acid | Etymology | Amino Acid | Etymology |
|---|---|---|---|
| glycine | greek, "Sweet" b/c it tastes sweet | asparagine | first found in asparagus |
| alanine | nonsense, euphonic | aspartic acid | similar to asparagine |
| leucine | greek, "white", first isolated as white crystals | glutamine | first found in wheat gluten |
| isoleucine | isomer of leucine: same atoms, different arrangement | glutamic acid | similar to glutamine |
| proline | shorten "pyrrolidine" | lysine | greek, "a breaking up", b/c first isolated in broken up molecules |
| phenylalanine | alanine + phenyl group | histidine | greek, "tissue" b/c first isolated from tissue protein |
| tyrosine | greek, "cheese" from which it was first isolated | arginine | latin "silver", first isolated in combination with silver atom |
| tryptophan | greek, "trypsin-appearing" b/c first discovered in after action of trypsin | methionine | methyl group attached to sulfur atom (called *theion* in greek) |
| serine | latin, "silk", from which it was first isolated | cystine | greek "bladder" b/c first isolated in bladderstone |
| threonine | related to sugar called 'threose' | valine | related to valeric acid |

Asimov, The Human Brain, 1965