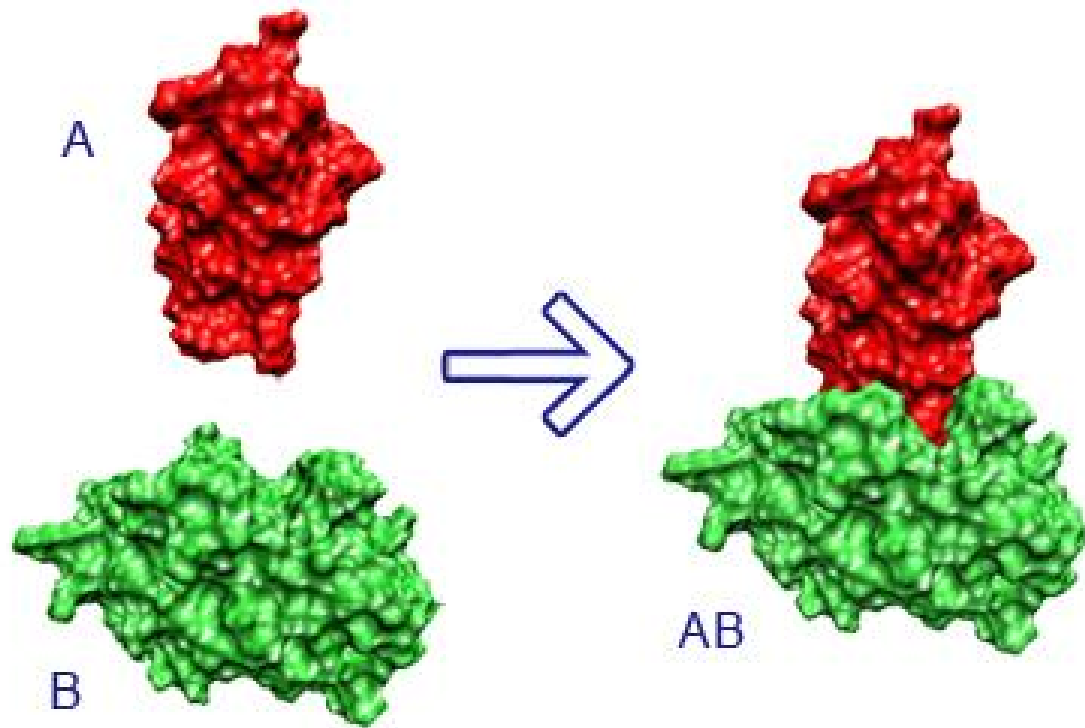# Structure-based prediction of protein-protein interactions on a genome-wide scale.

**Qiangfeng Cliff Zhang, Donald Petrey, Lei Deng**

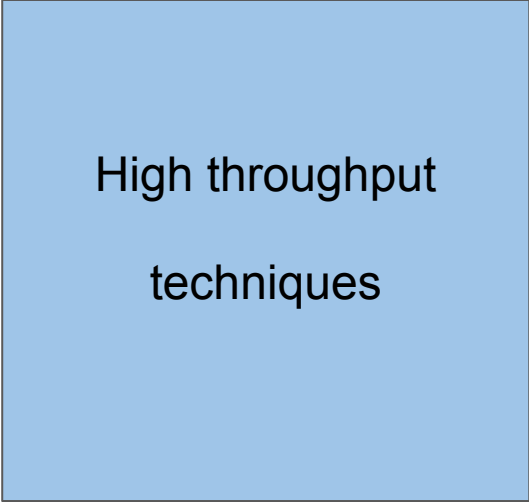# Structure-based prediction of protein-protein interactions on a genome-wide scale.

**Qiangfeng Cliff Zhang, Donald Petrey, Lei Deng**

# How do we know what we know?

# How do we know what we know?

High throughput

techniques

# How do we know what we know?

High throughput

techniques

Computational
methods based on
non-structural
evidence

# Structure-based prediction of protein-protein interactions on a genome-wide scale.

**Qiangfeng Cliff Zhang**

# Structure-based prediction of protein-protein interactions on a genome-wide scale.

Homology Models

**Qiangfeng Cliff Zhang**
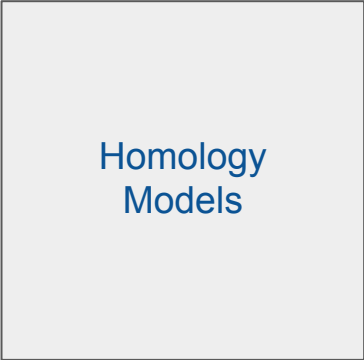
# Structure-based prediction of protein-protein interactions on a genome-wide scale.

**Qiangfeng Cliff Zhang**

Homology Models

Structural Neighbors

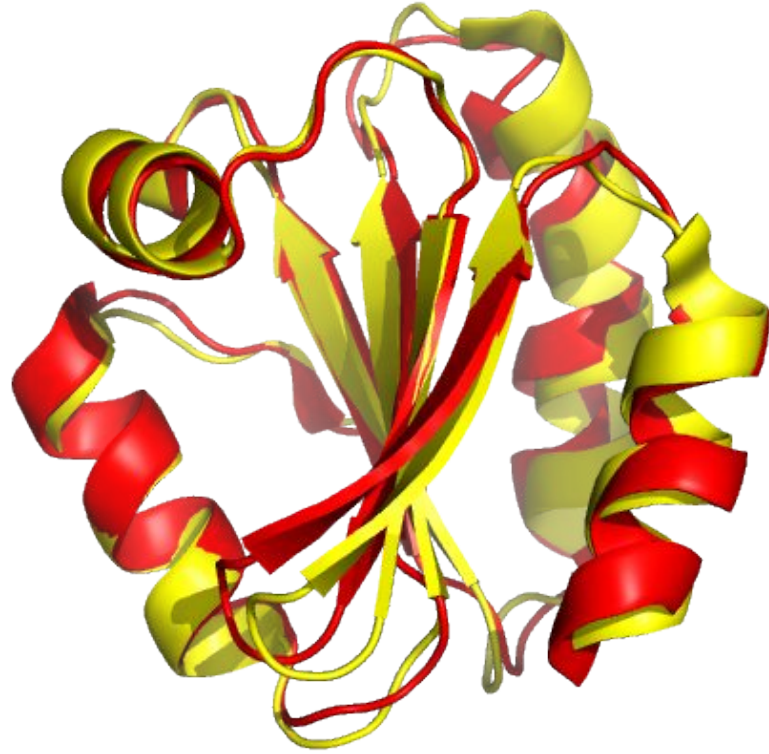Homology
Models

Homology
Models

**Sequence
Alignment**

Structural
Neighbors

Structural
Neighbors

**Structural
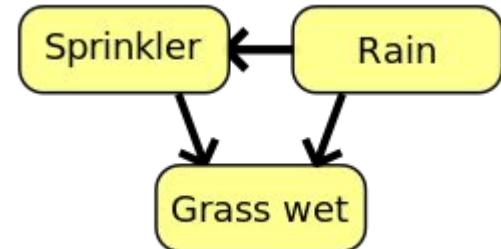Alignment**

Structural Neighbors

**Structural Alignment**

# Structural information
# +
# Non-structural information

**Structural information**
**+**
**Non-structural information**

**=Pre-PPI**

Sprinkler ← Rain

Sprinkler → Grass wet

Rain → Grass wet
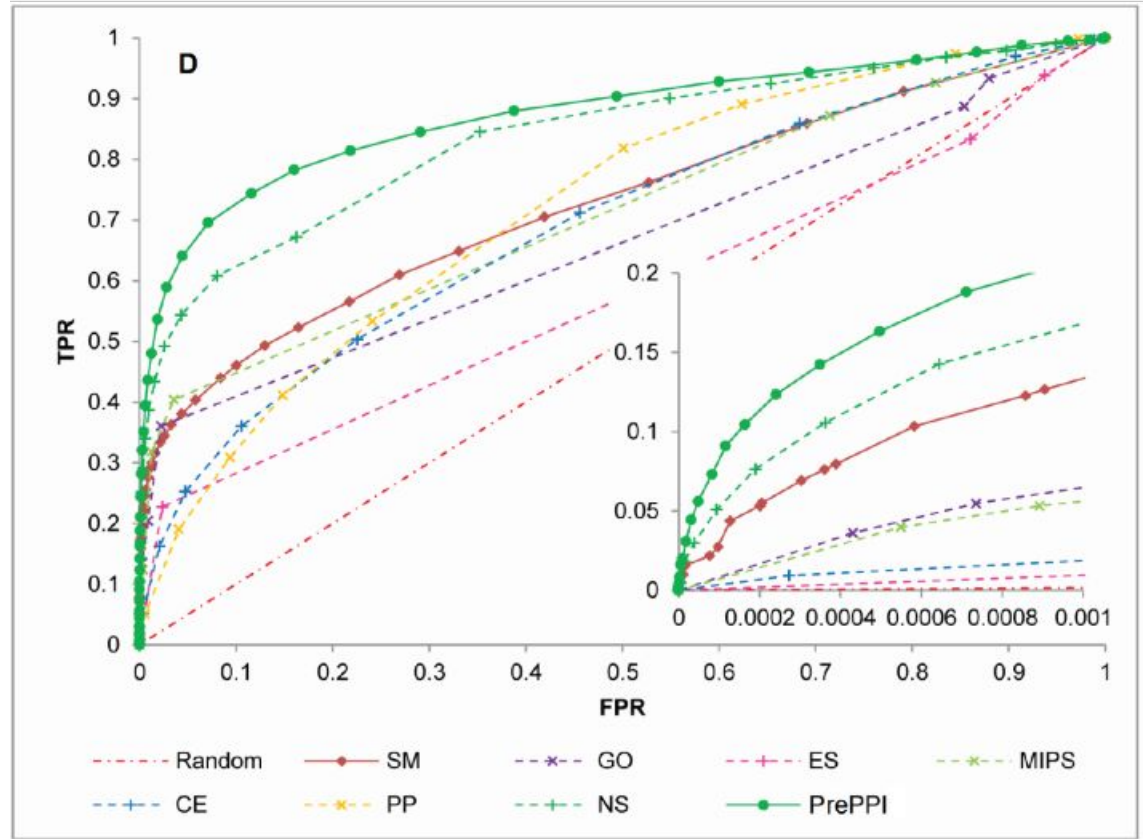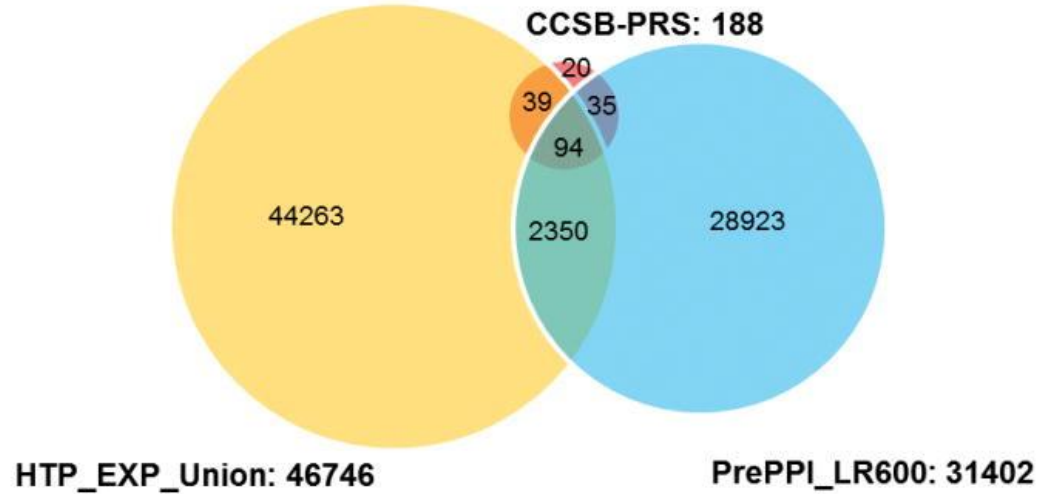
# How good is it?

# How good is it?

Pre-PPI

VS.

Structural/
Non-structural info.
alone

# How good is it?

Pre-PPI

VS.

High Throughput
Techniques



CCSB-PRS: 188

20

39  35

94

44263        2350        28923

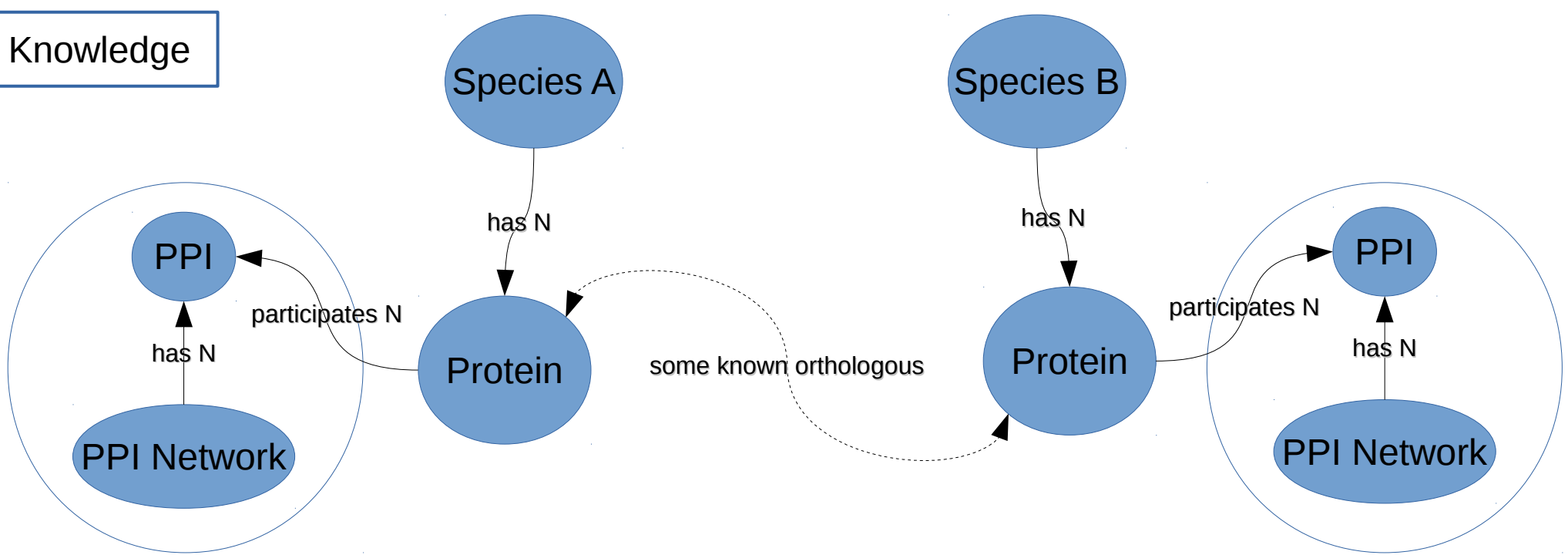HTP_EXP_Union: 46746                PrePPI_LR600: 31402

# Panorama of ancient metazoan macromolecular complexes

A paper with two primary authors
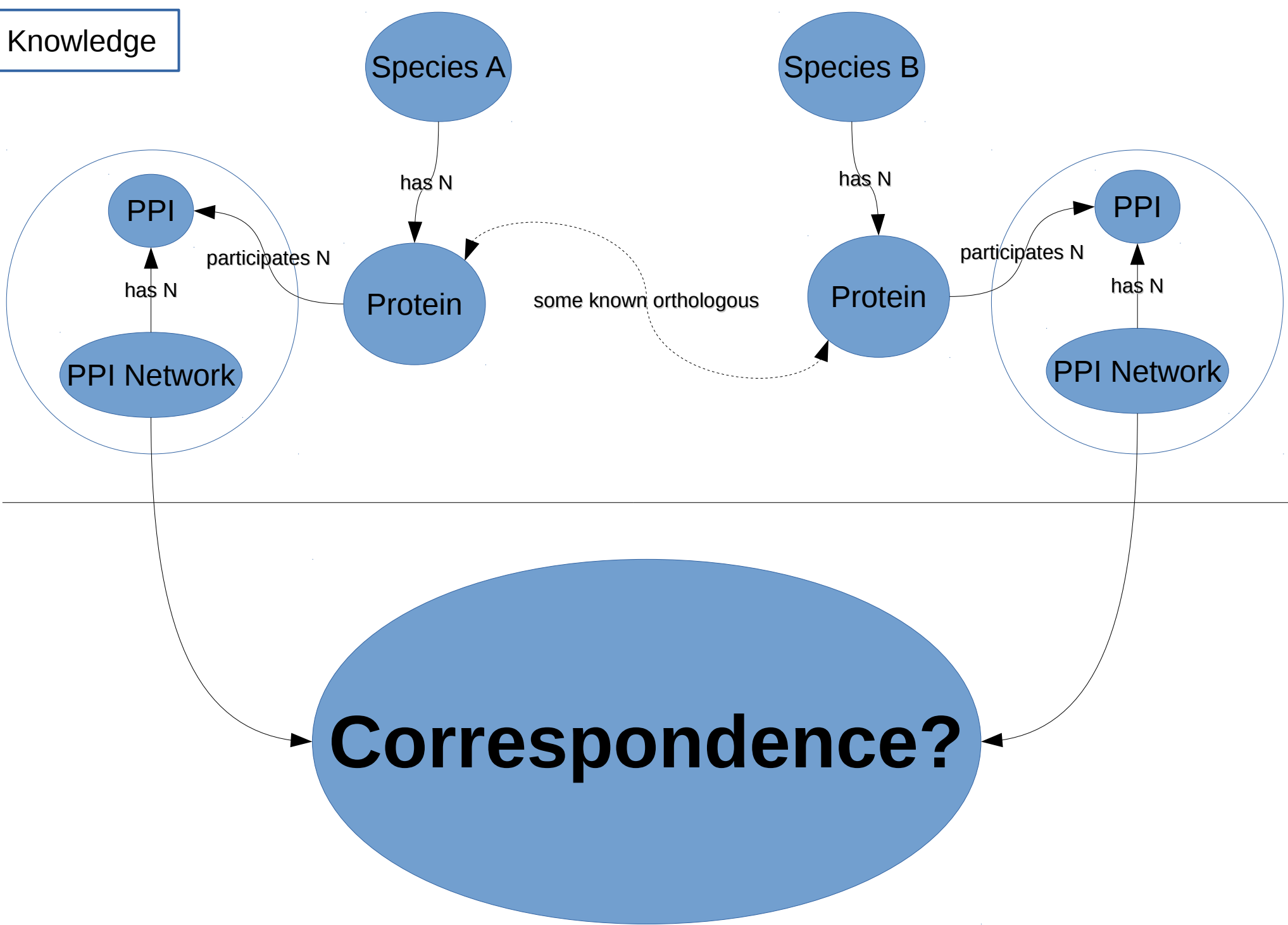and too many others to list

# Who should care?
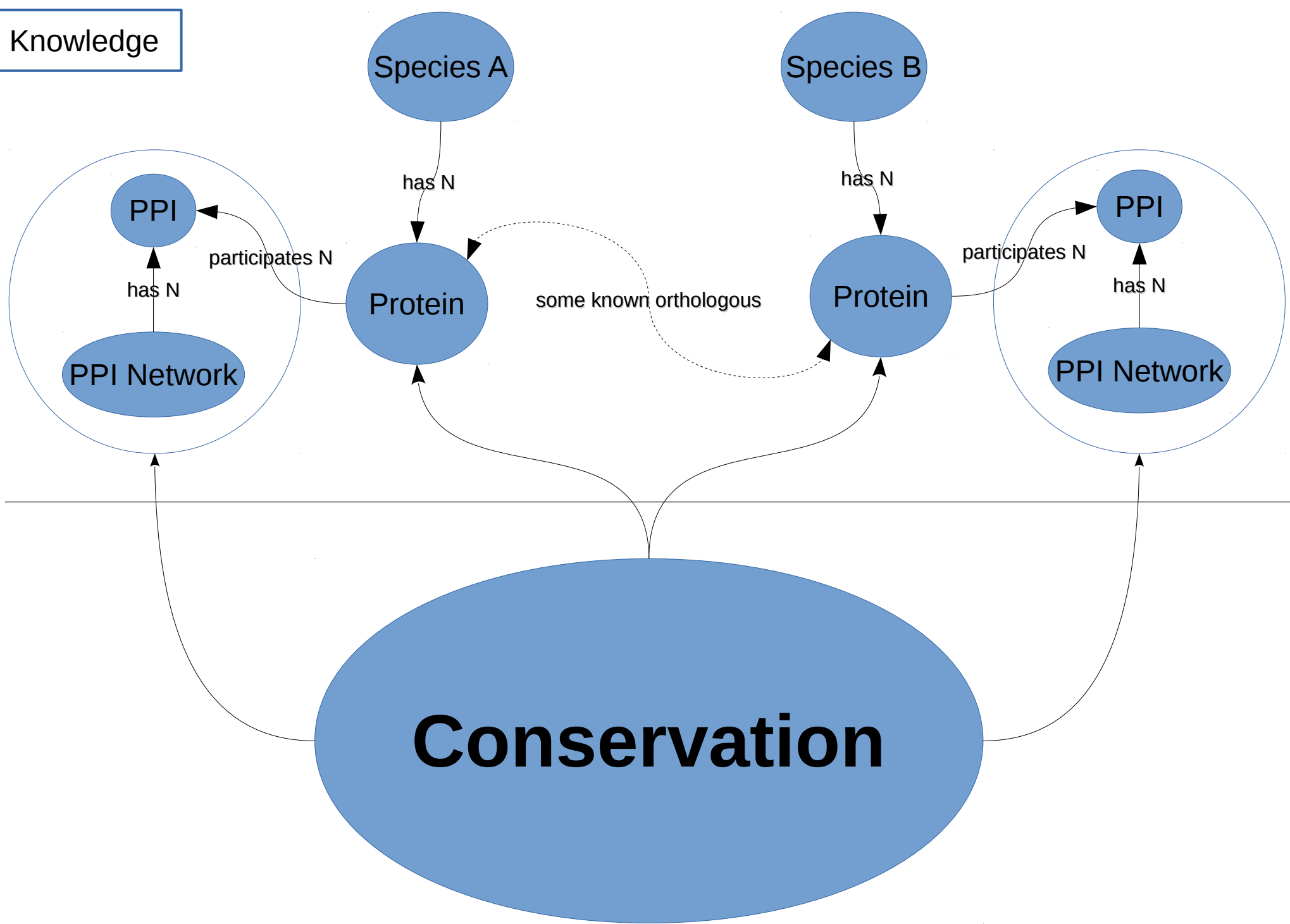
folks who want to
**understand**
*cellular processes*

Species A
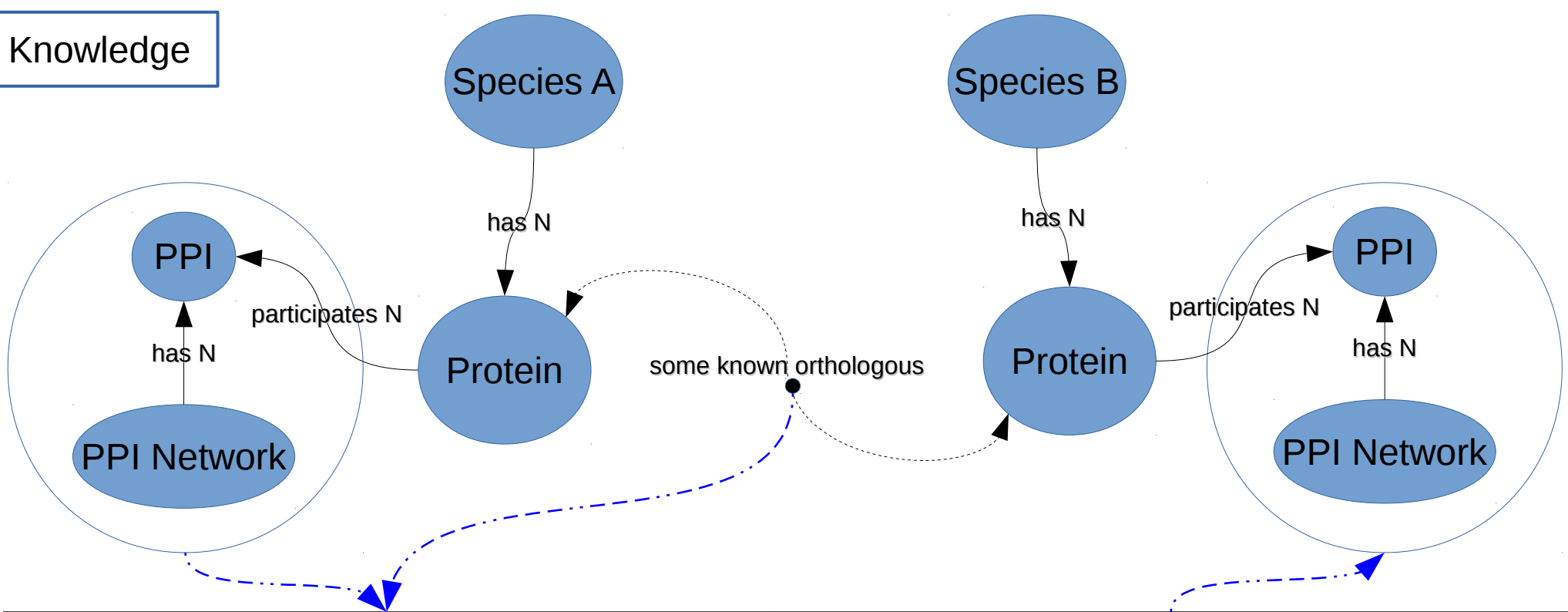
Species B

has N

has N

PPI

PPI

participates N

participates N

Protein

Protein

some known orthologous

has N

has N

PPI Network

PPI Network

**Conservation**
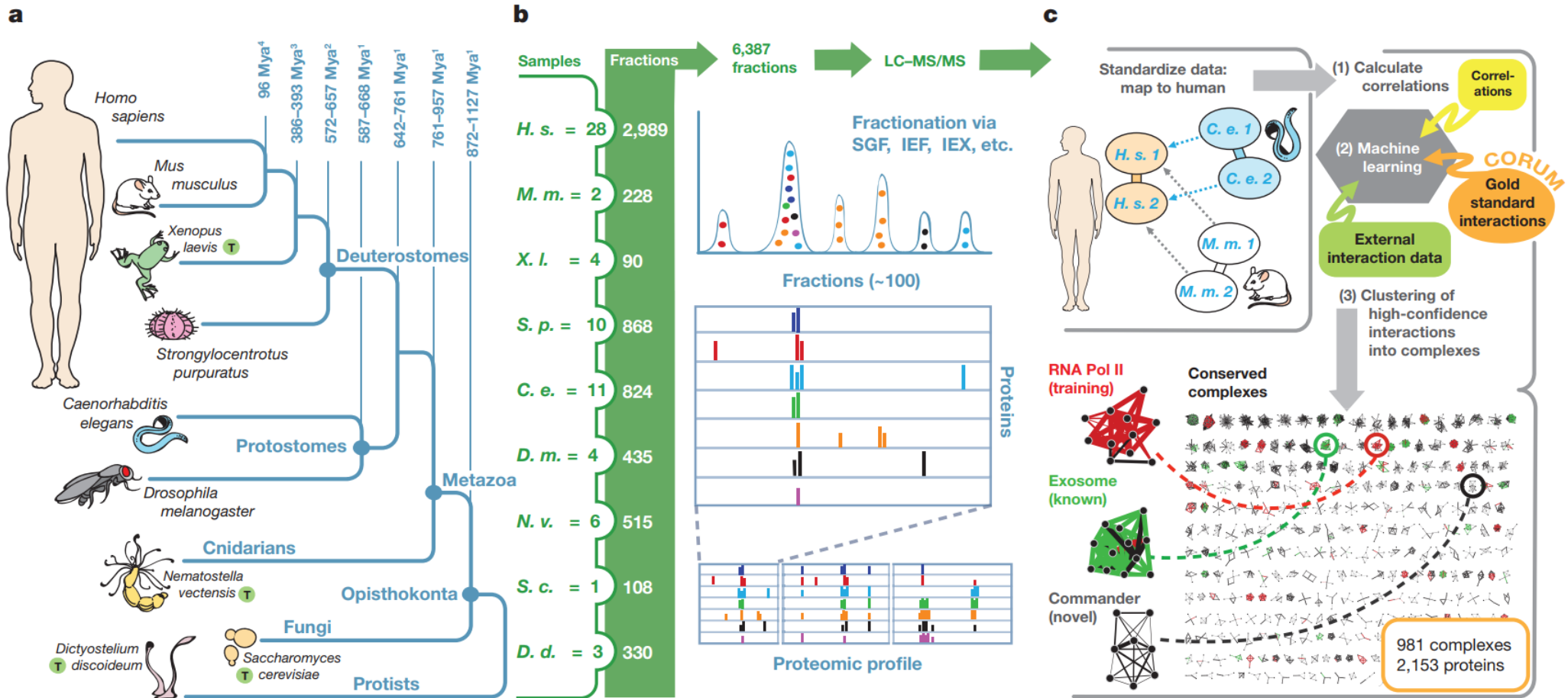
# 10,000ft Objective

- Describe extent of PPI network conservation

- Identify conserved PPI networks

- Characterize conserved PPI networks

# 10,000ft Objective

- Describe extent of PPI network conservation

- Identify conserved PPI networks

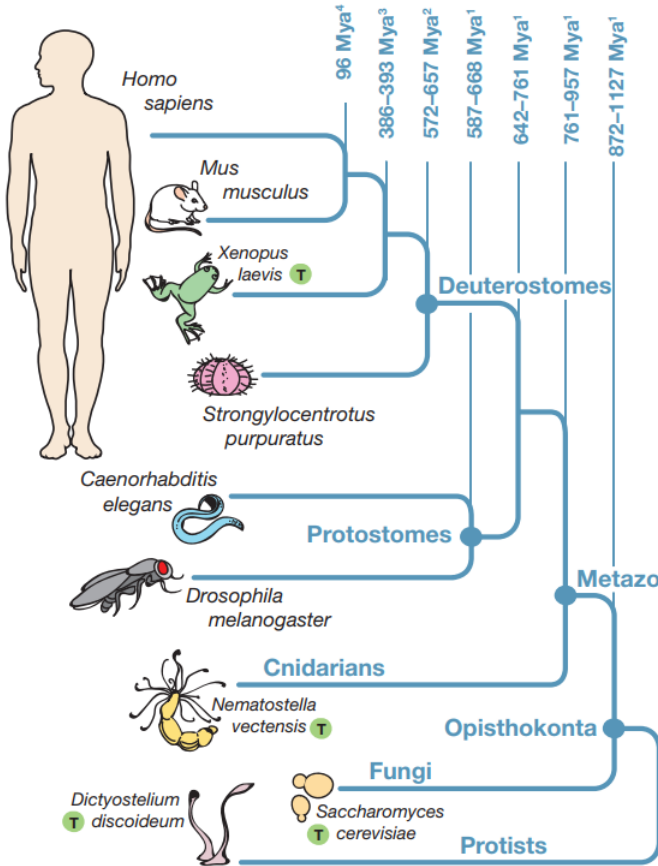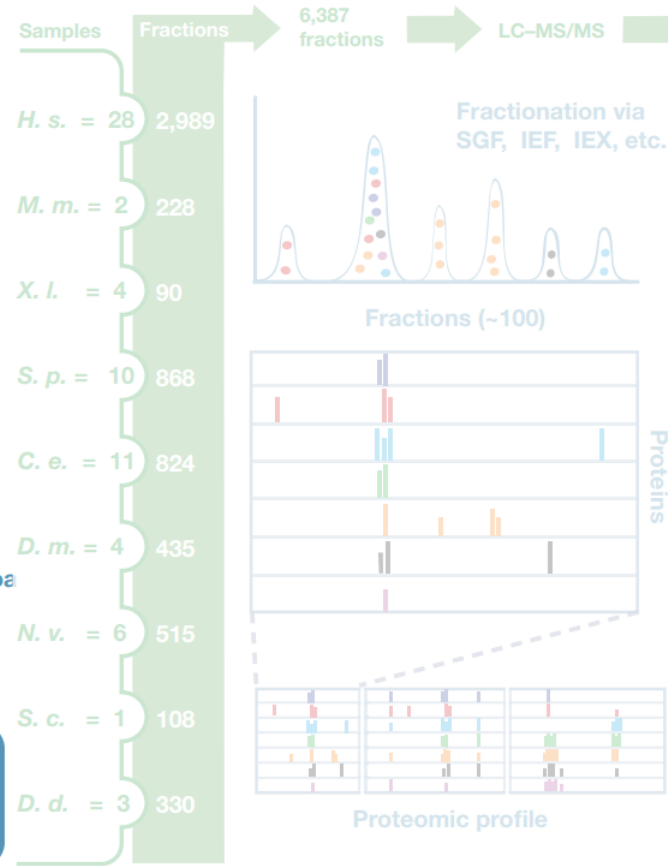- Characterize conserved PPI networks
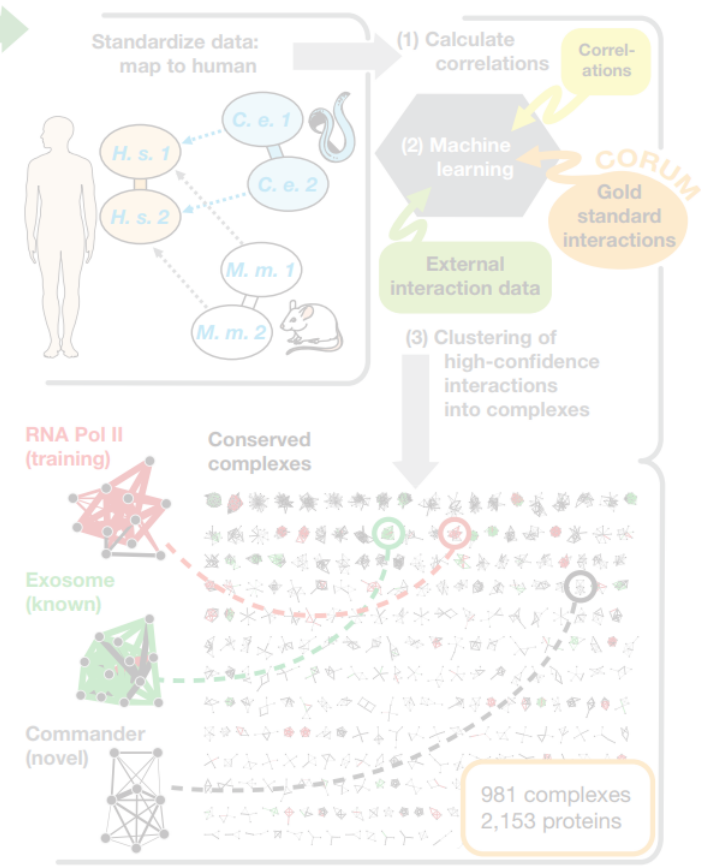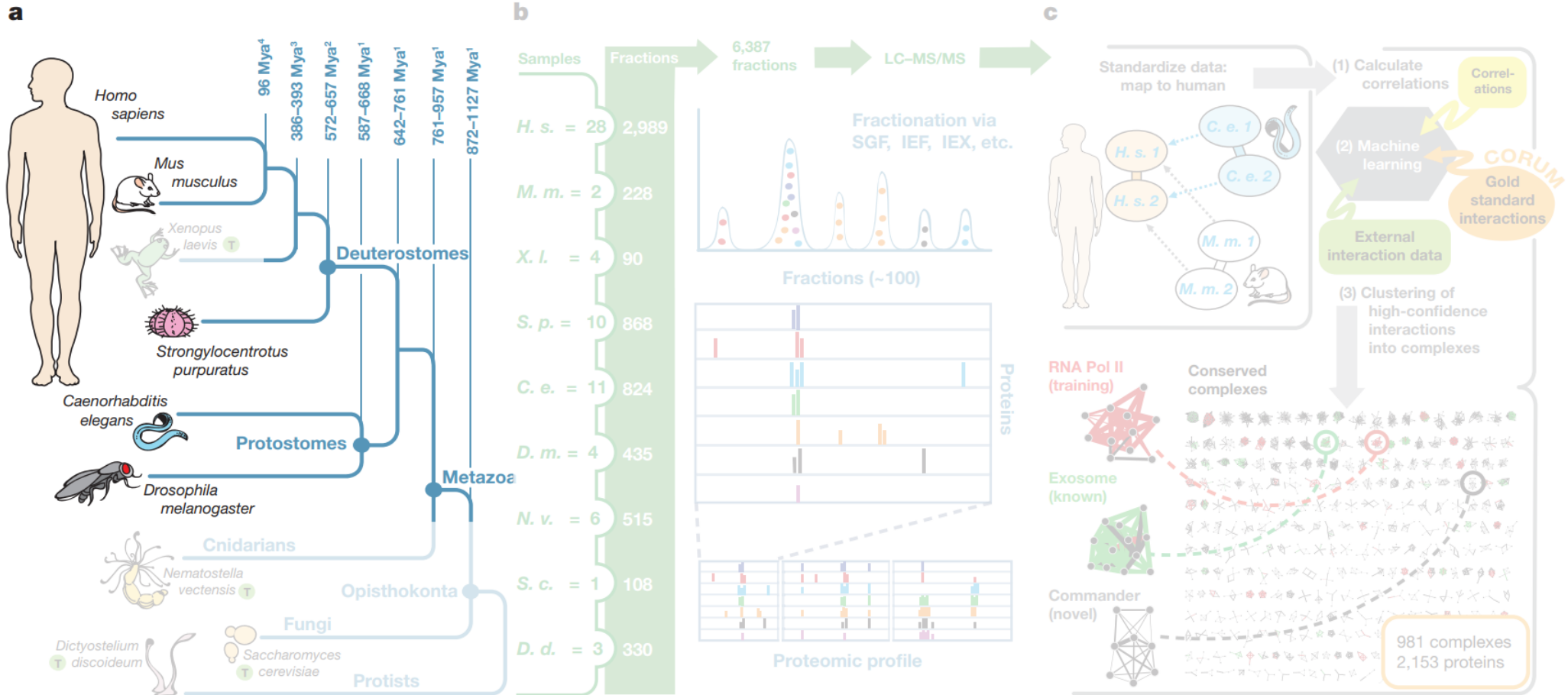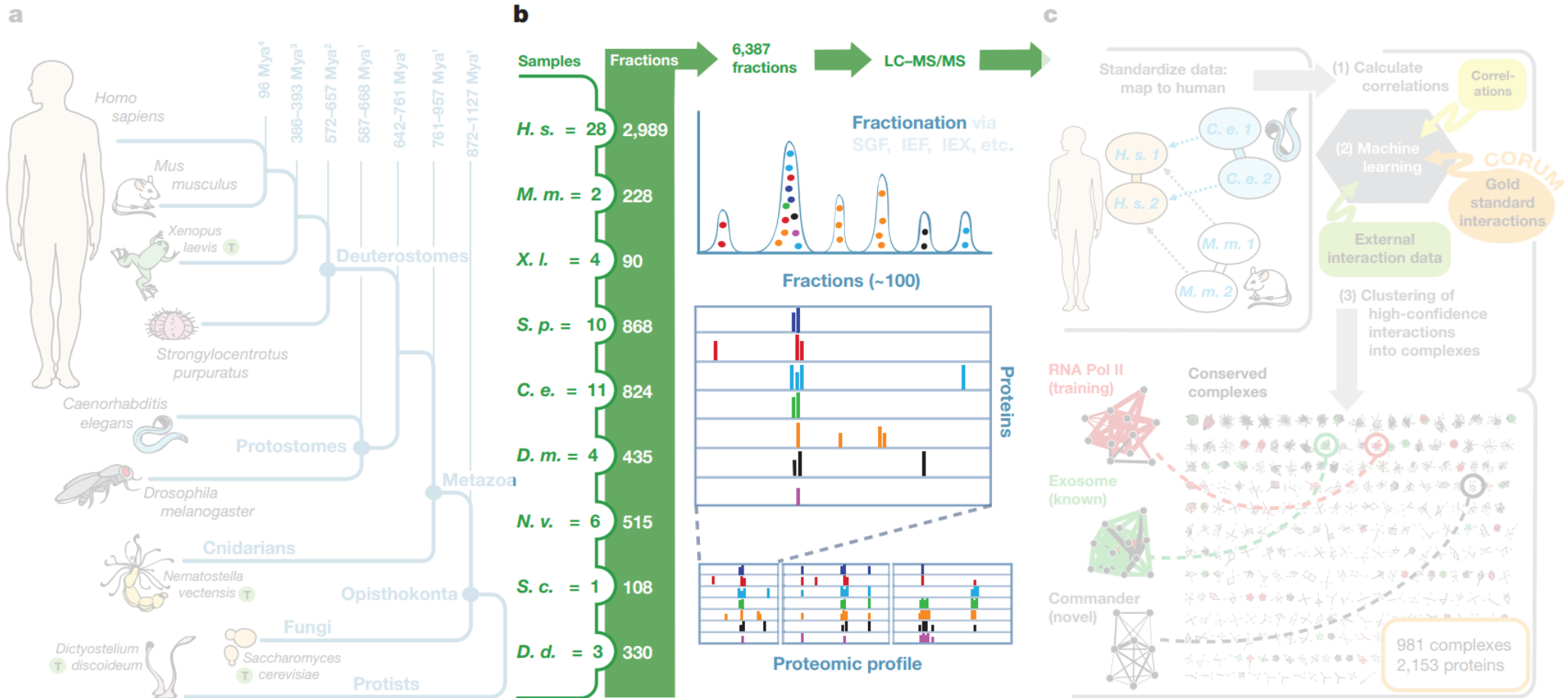
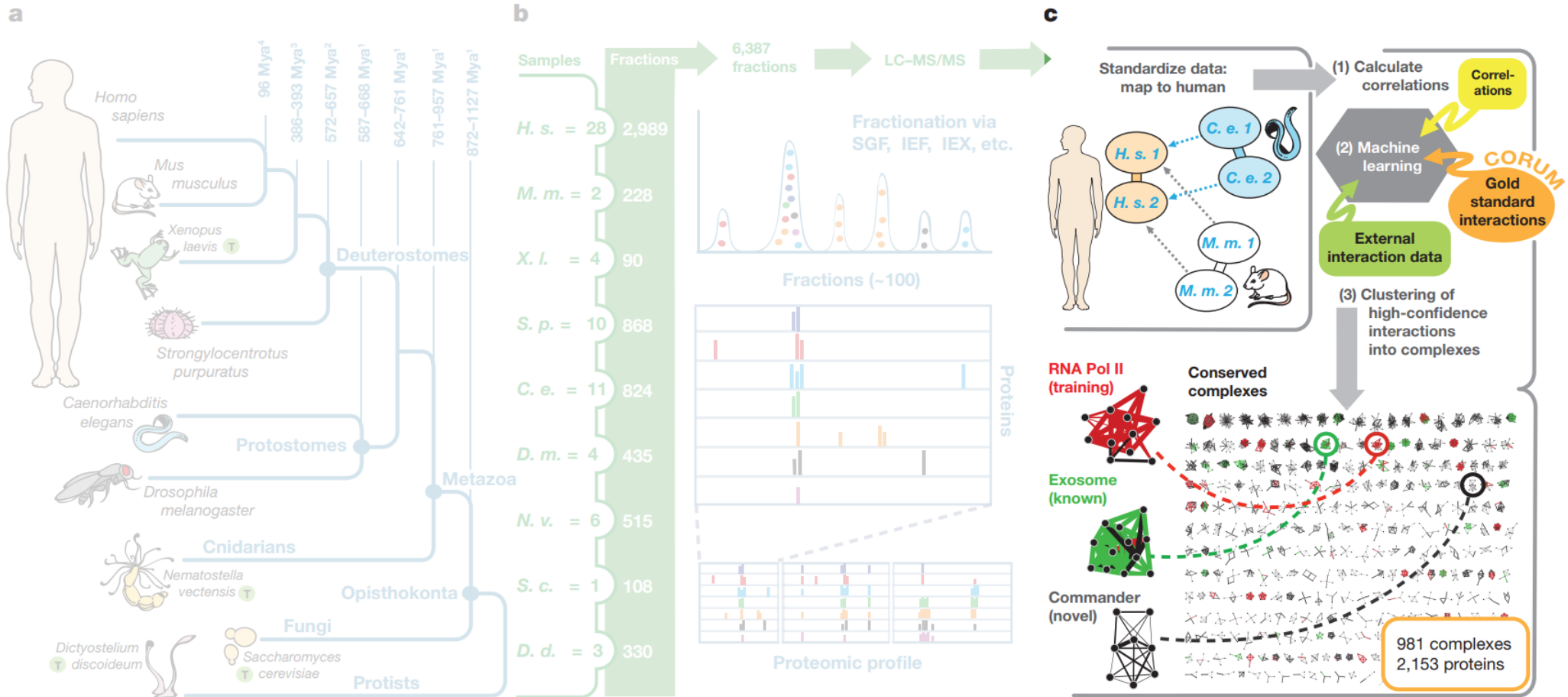# Methods: Overview

# Population

# Population

Training

# Data Acquisition

# Data Analysis

The data analysis is the hard part.

Naturally that's what we're gonna focus on.

# Basis for Comparison

***Have***

– Interspecies PPI pairs

***Want***

– Comparable PPI pairs

***Solutions***

– PPI pairs over orthogroups

– PPI pairs normalized to a single species

– ???

# Basis for Comparison
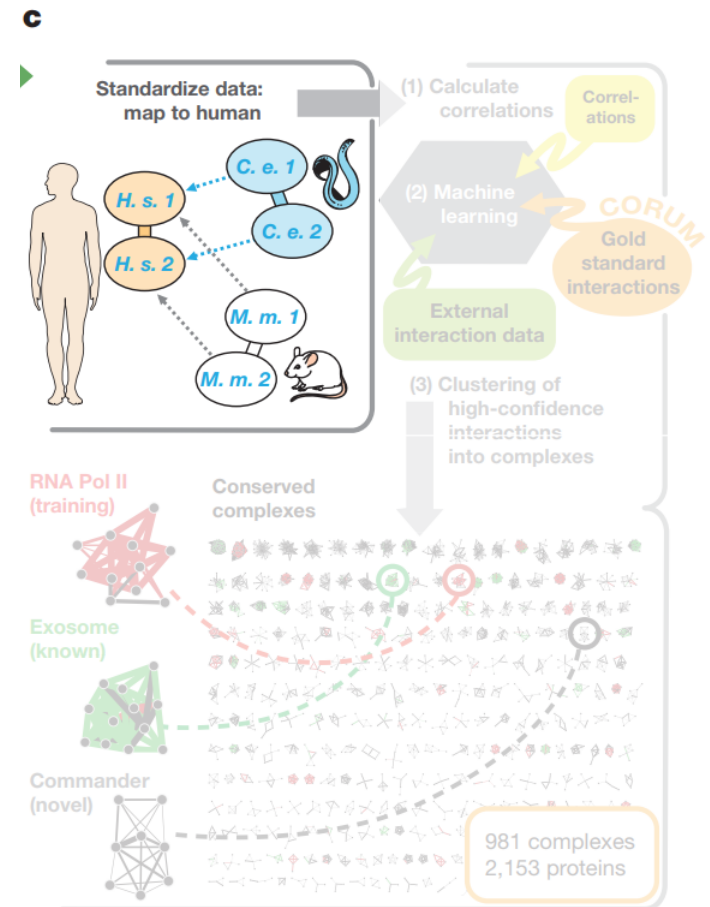
**_Have_**

– Interspecies PPI pairs

**_Want_**

– Comparable PPI pairs

**_Solution_**s

– PPI pairs over orthogroups

– PPI pairs normalized to a single species *(human)*

– ???

arbitrary

# Pause for a thought

***Have***

– Comparable PPI pairs w/ co-fractionation data
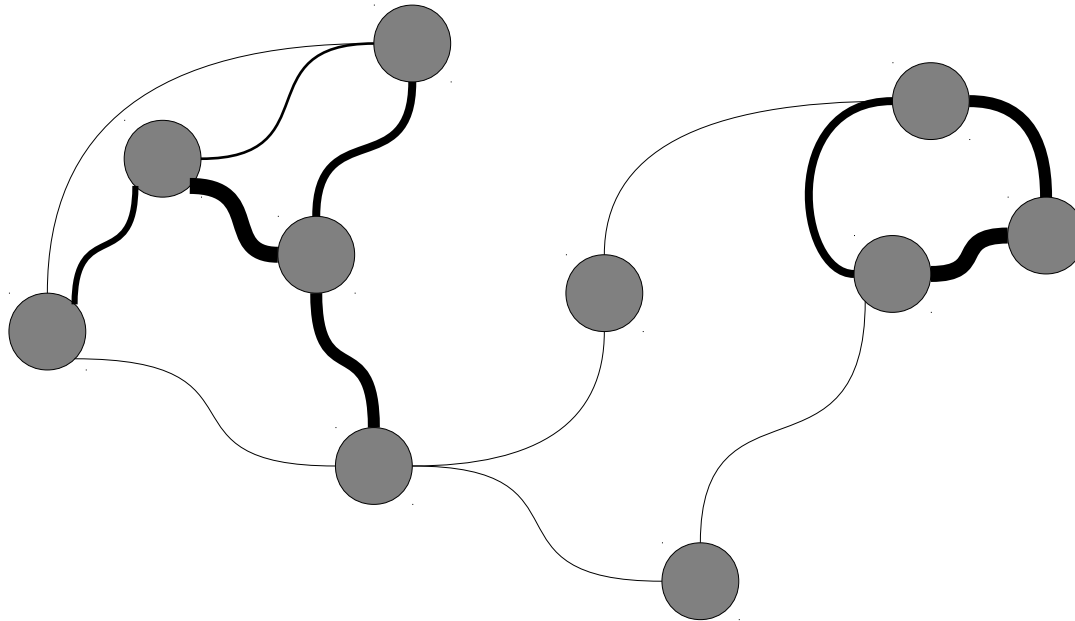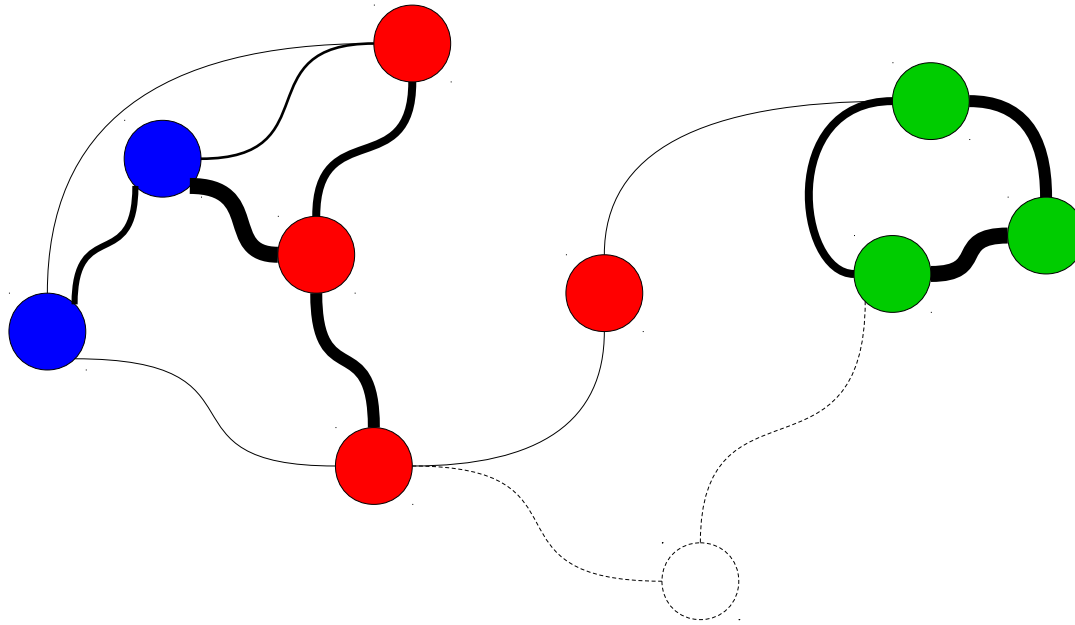
***Want***

– Real PPI networks

***A Solution***

– Score PPI pairs for being co-complex

  - Informed by gold-standard adjacencies

– Transform PPI pairs into a global adjacency network

– Split out subnetworks from global PPI network

  - Tuned for meaningful PPI networks
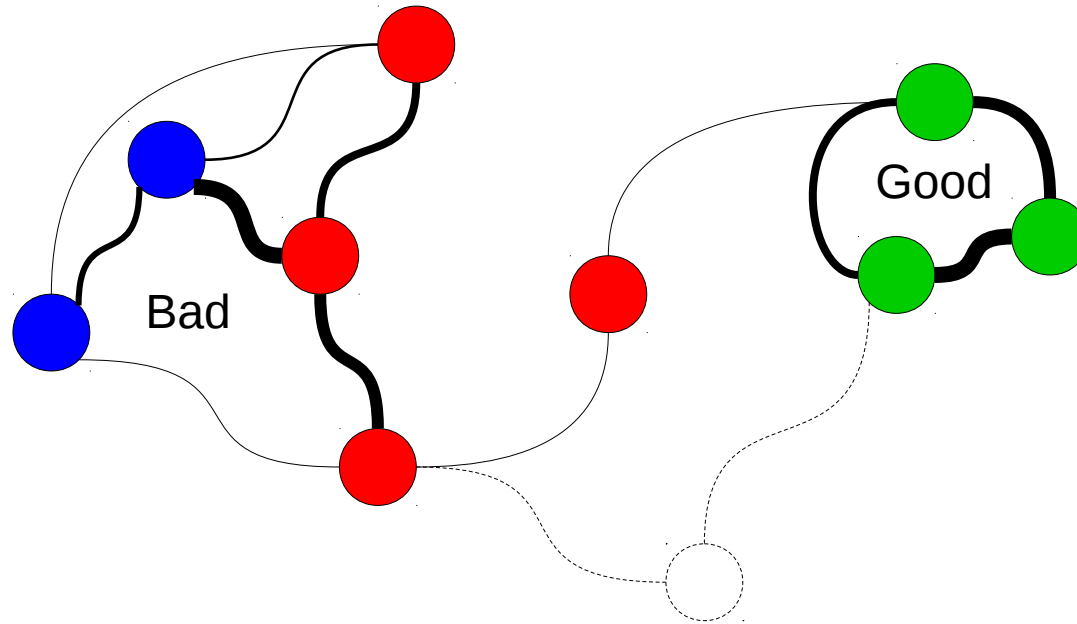
# Graph Clustering Primer
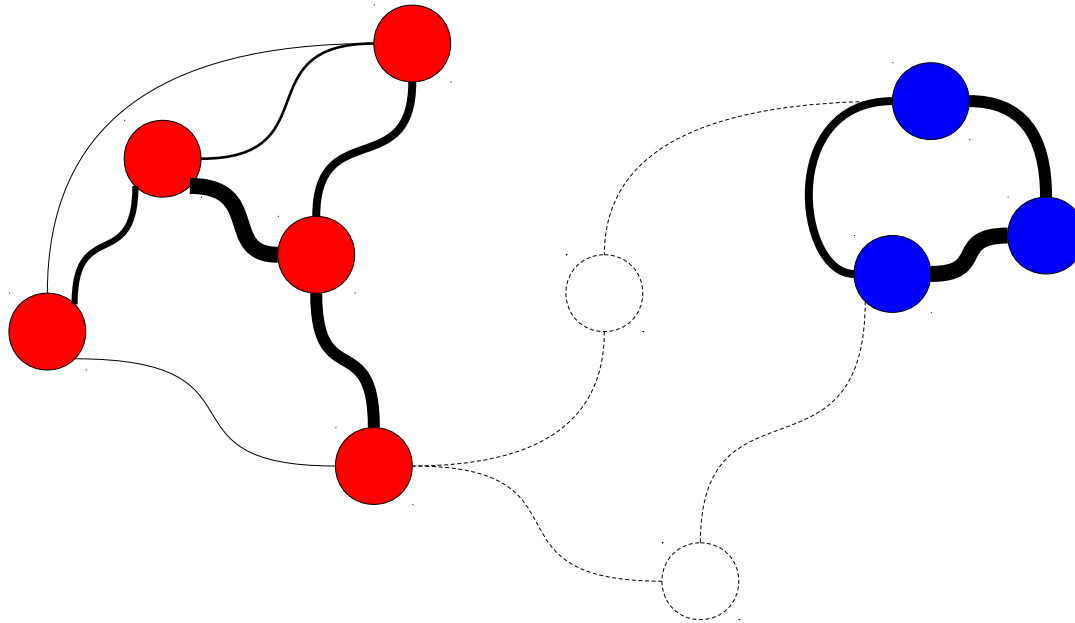
# Graph Clustering Primer



example of a 'clustering'

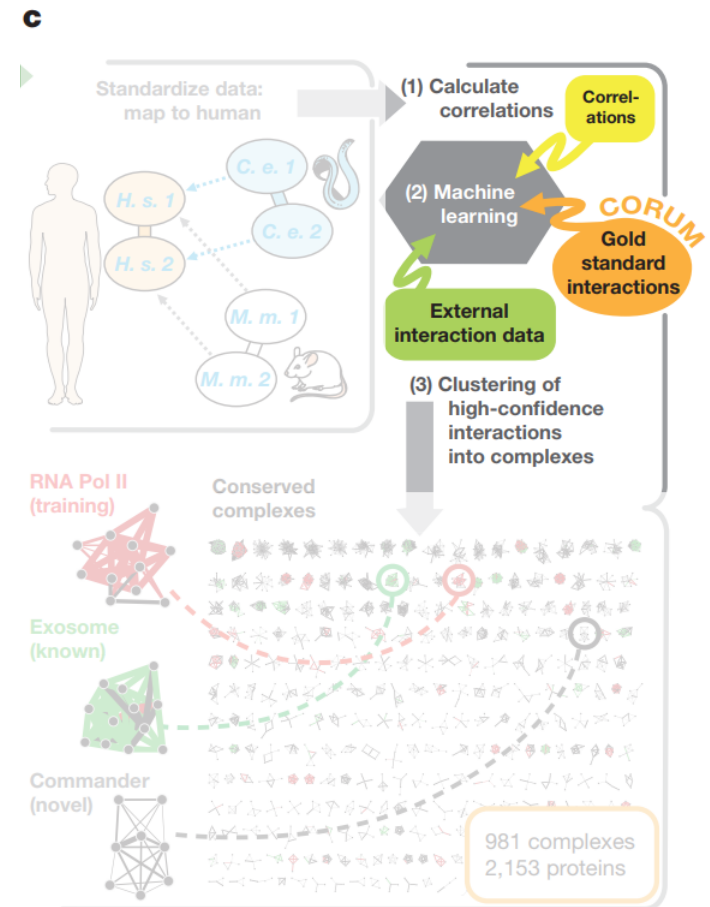# Graph Clustering Primer

# Graph Clustering Primer



probably want something like this

We now return to your scheduled programming

# Gluing together a data pipeline

- Generate features
  - Co-fractionation correlation measures
  - Interaction data from databases, literature
- ML co-complex +/- against CORUM
- Network clustering w/ ClusterOne, MCL
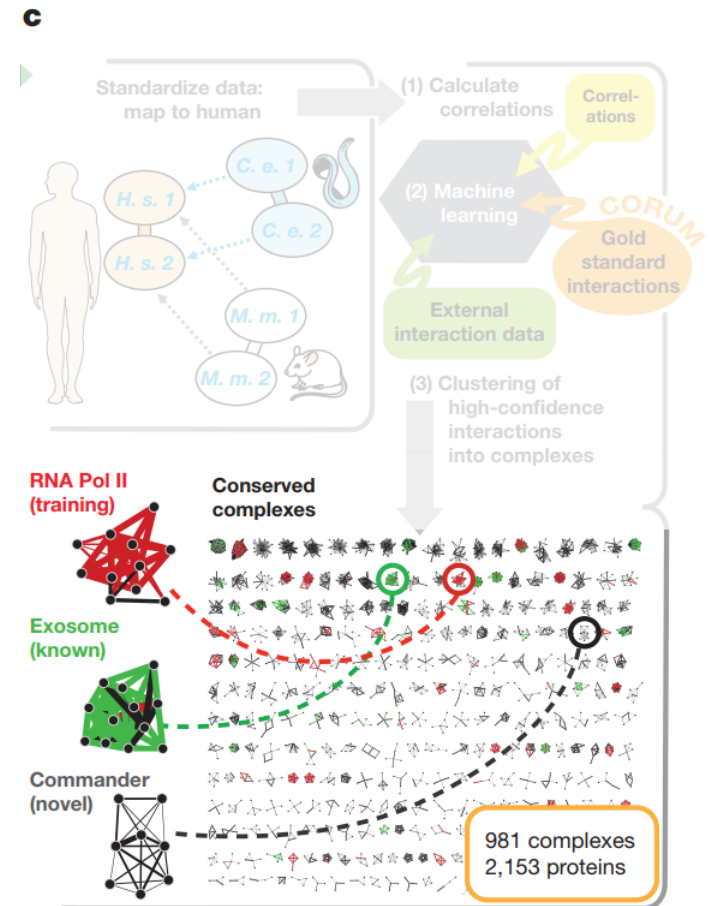
# Characterization of Clusters

**Have**

- PPI networks embedded in global network

**Want**

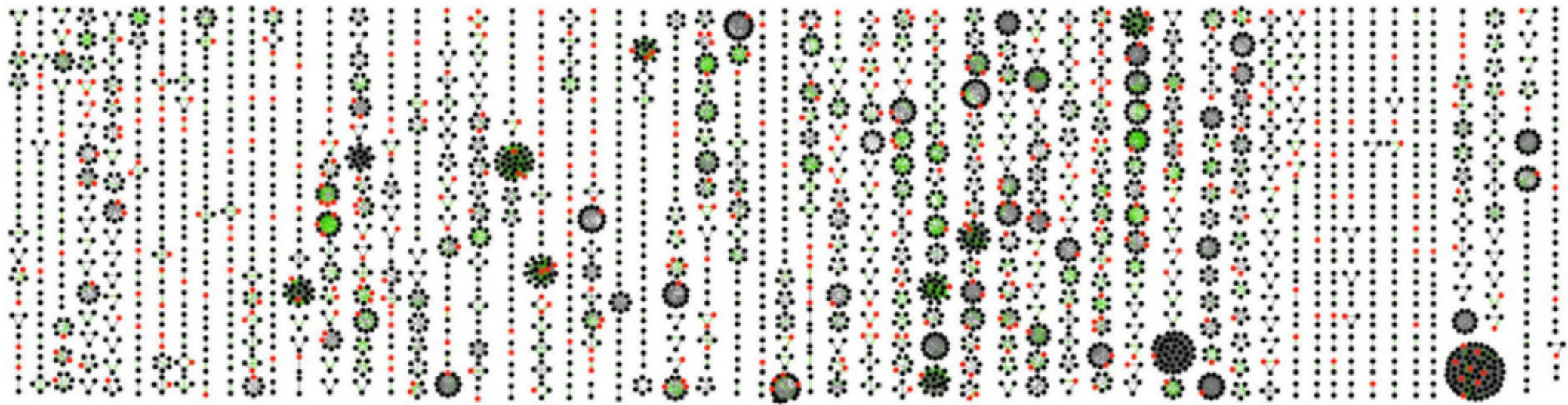- Biologically relevant and 'interesting' features

**A Solution**

- Enrichment analysis for higher-order attributes

# Smattering of Results

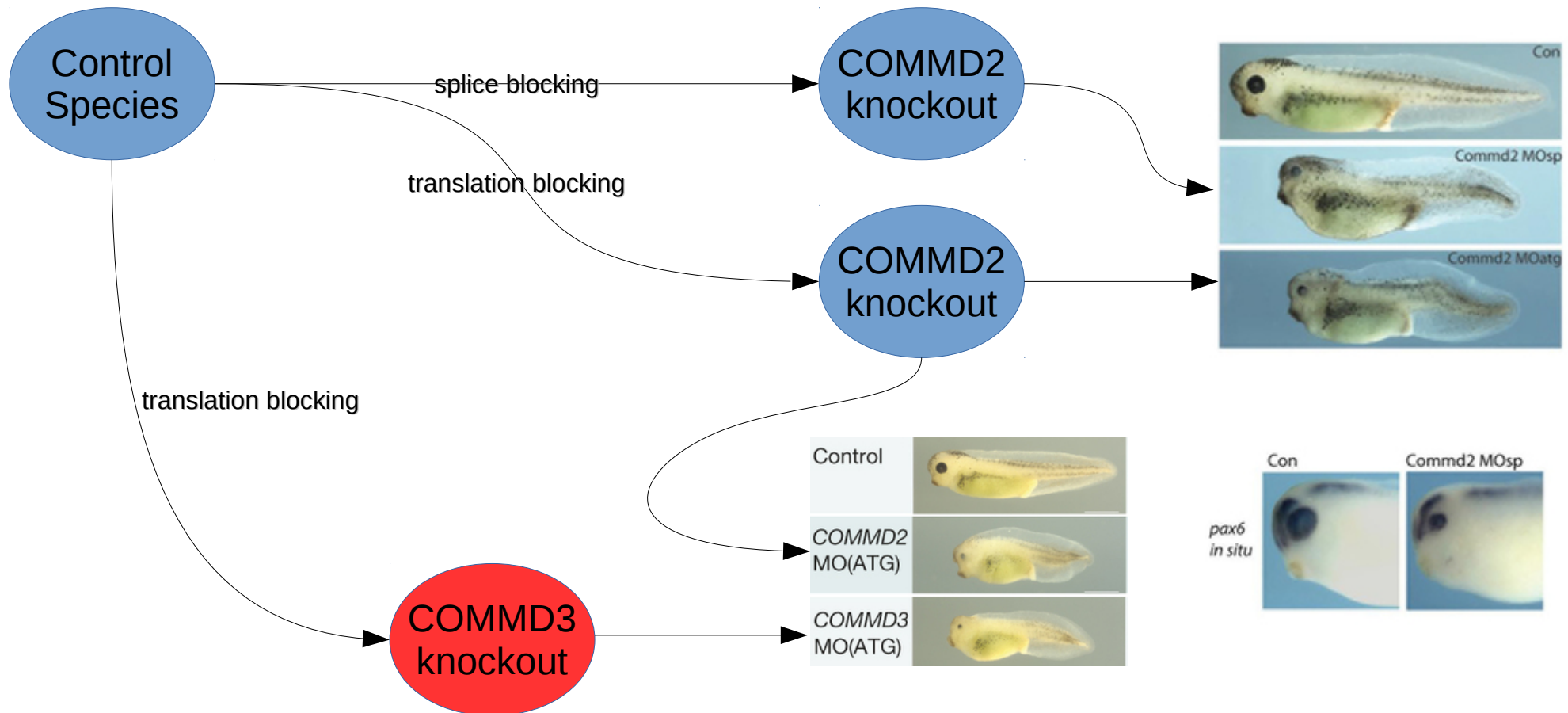## found many neglected proteins/PPIs of interest



**green** links are novel co-complex interactions          **red** dots are unannotated proteins

# Smattering of Results

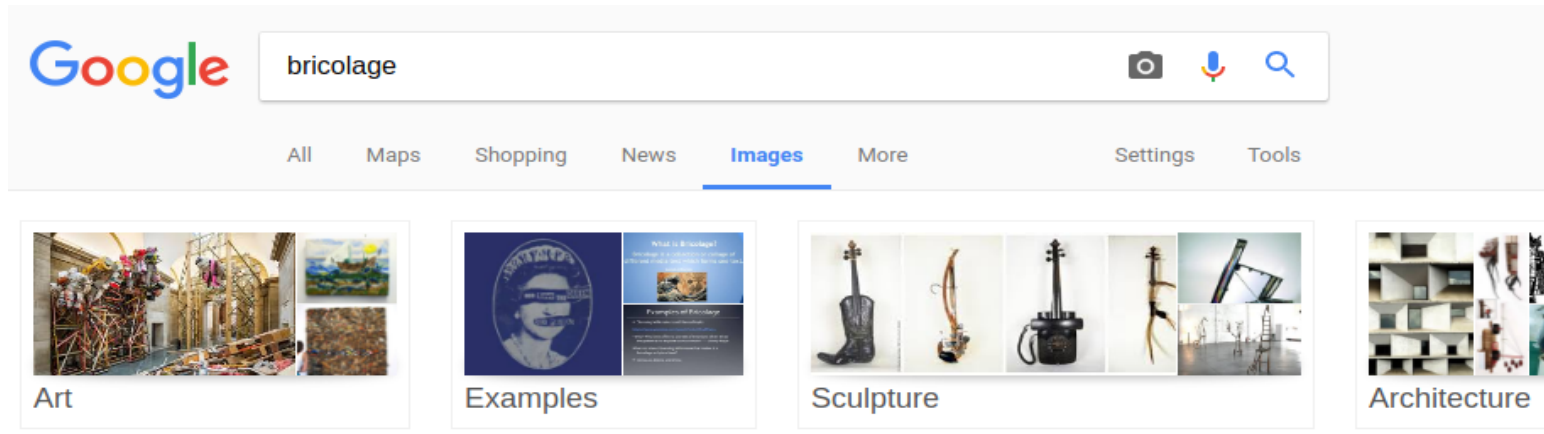## directly demonstrated utility of approach

# Things I wish the authors mentioned

- Were multiple co-fractionation correlation measures really necessary? Which features were relevant?

- Did they account for the false-discovery rate of clusters?

- Was there a specific reason for choosing humans as the ortholog basis? How much do results change when using a different species?

# Final Comments

The candy is a lie