# Hidden in plain sight: The eukaryotically conserved unstudied proteins and a framework for their classification and characterisation.
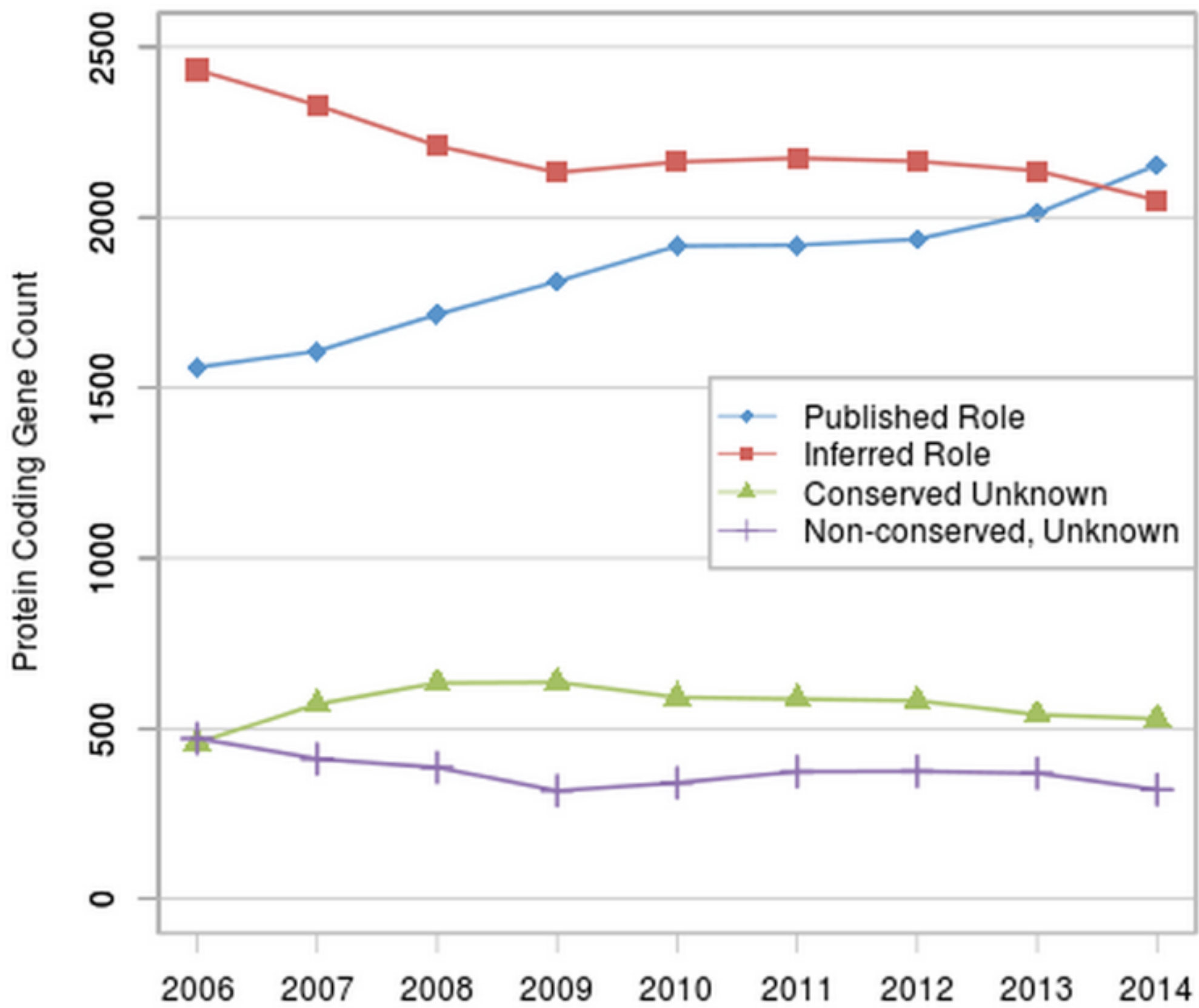
Wood V [1], Bahler J [2], Harris M [1], Lock A [2], Oliver SG [1]

## ABSTRACT

Proteins conserved widely among eukaryotes play fundamentally important roles in the shared, basic mechanisms of life. The roles of many broadly conserved proteins remain unknown, however, despite almost a century of gene- and gene product-specific genetic and biochemical investigation. Even the recent emergence of genome-wide experimental techniques and the availability of near-complete protein inventories for many intensively studied eukaryotic model species have shed light on the functions of few previously uncharacterised conserved proteins. Because the success of many endeavours in basic and translational research, including drug discovery, metabolomics, and systems biology, depends critically on comprehensive representation of conserved functions, a more complete understanding of protein components conserved throughout eukaryotes would have far-reaching benefits for biological research in many species and on a wide range of scales.
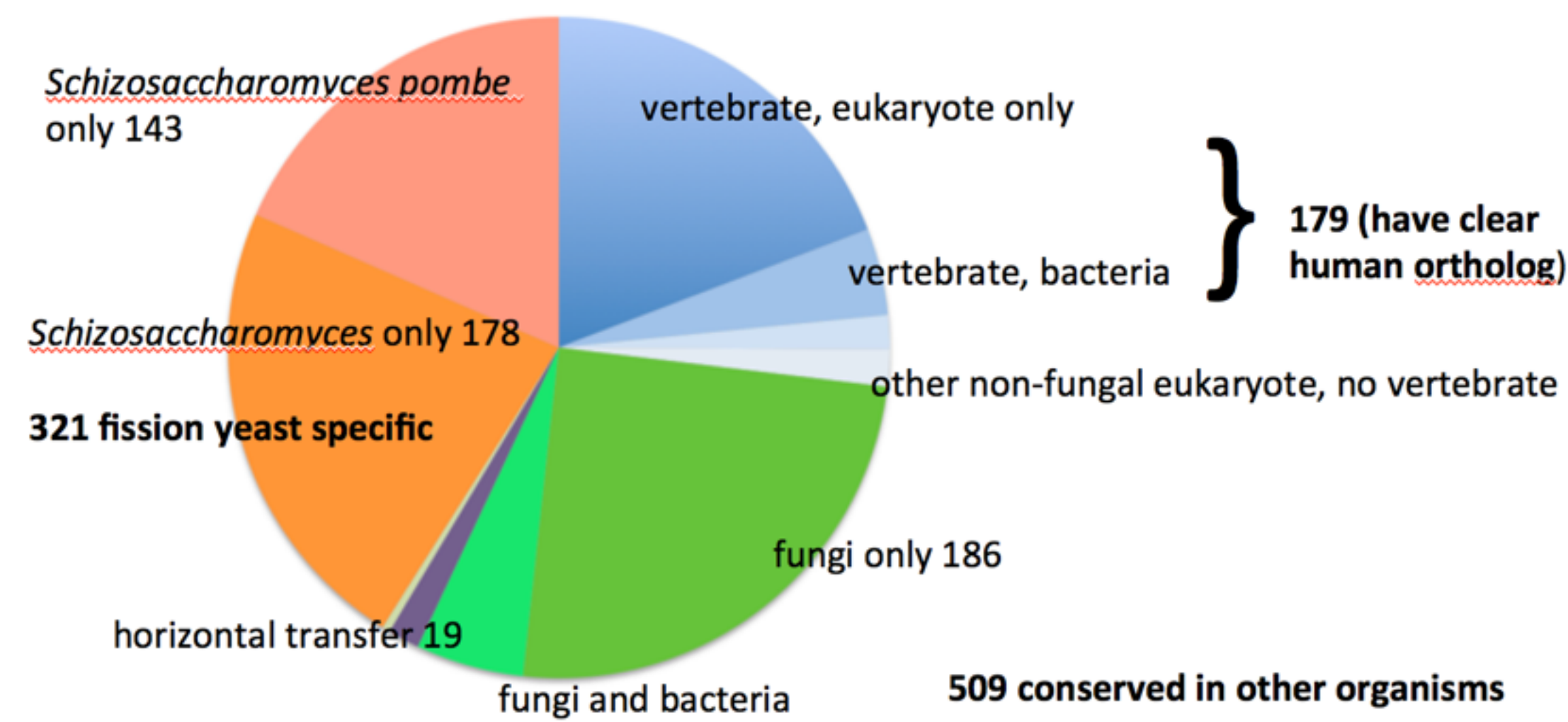
To identify priority targets for experimental investigation, PomBase provides an inventory of fission yeast proteins that are conserved among eukaryotes but whose broad biological roles remain unknown. A broad functional classification of the known proteome using a selection of Gene Ontology biological process categories ("GO Slim") has revealed correlations with features such as subcellular localization and morphological phenotype. Combining available data from genome-wide phenotype and localization experiments with insights from the functional classification of known proteins facilitates prediction of biological roles, and thereby guides specific experimental characterisation of unknown proteins.
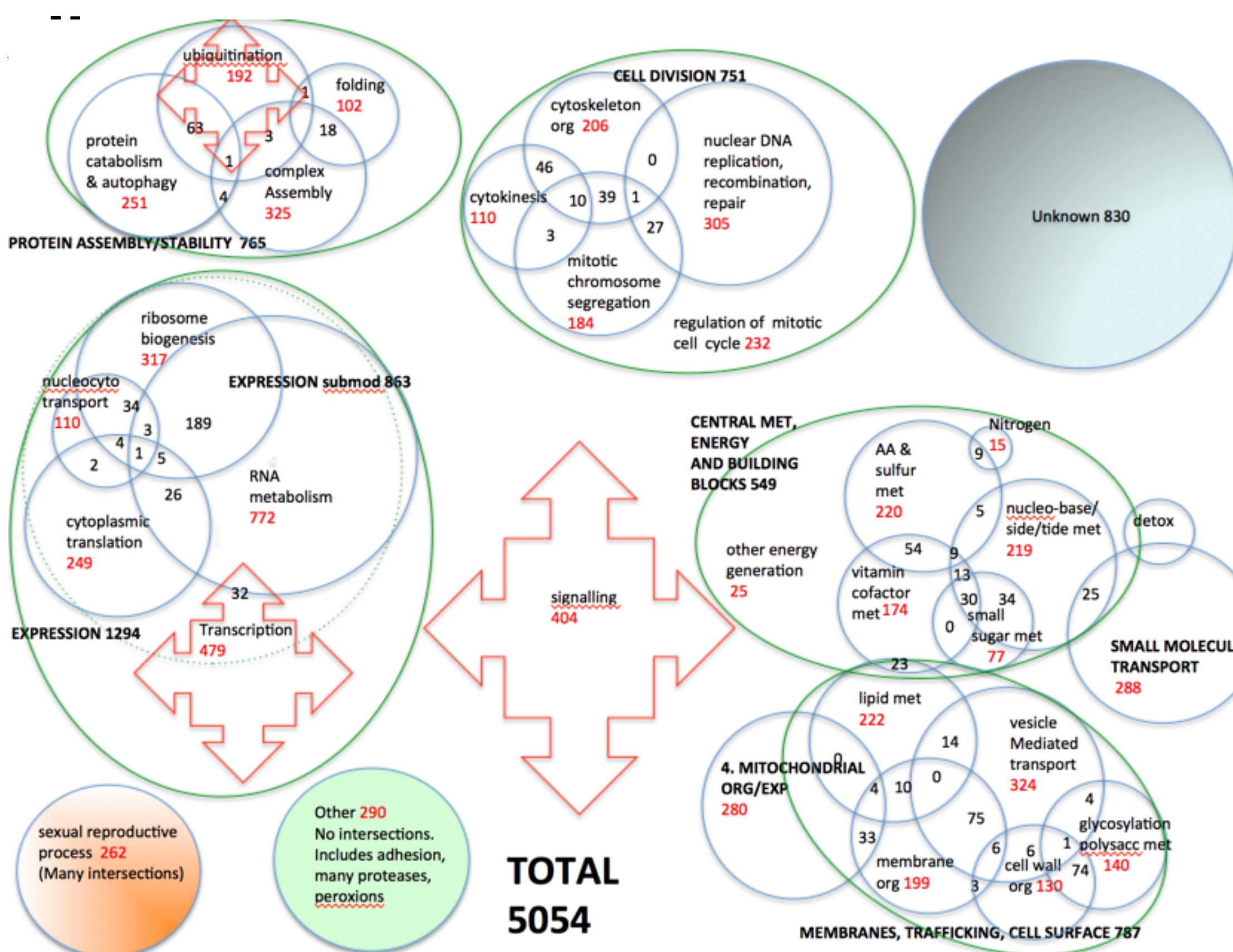
## We tend to study what we know



Legend: Published Role; Inferred Role; Conserved Unknown; Non-conserved, Unknown

lack of change in the conserved unknown inventory over the pcast decade

## Taxonomic distribution of Unknowns :



*Schizosaccharomyces pombe* only 143

vertebrate, eukaryote only

vertebrate, bacteria

} 179 (have clear human ortholog)

*Schizosaccharomyces* only 178

other non-fungal eukaryote, no vertebrate

321 fission yeast specific

fungi only 186

horizontal transfer 19

fungi and bacteria

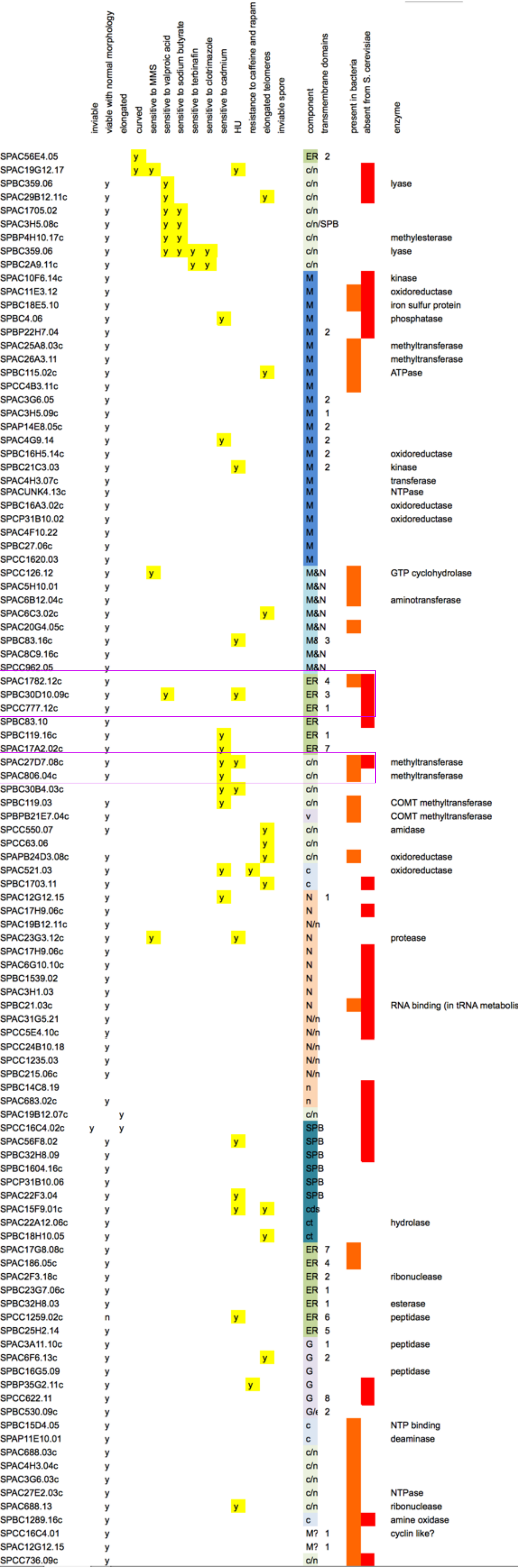509 conserved in other organisms

## Classifying Knowns: A Visual GO slim:



First step: Aim to assign unknowns to these broad classifiers

TOTAL 5054
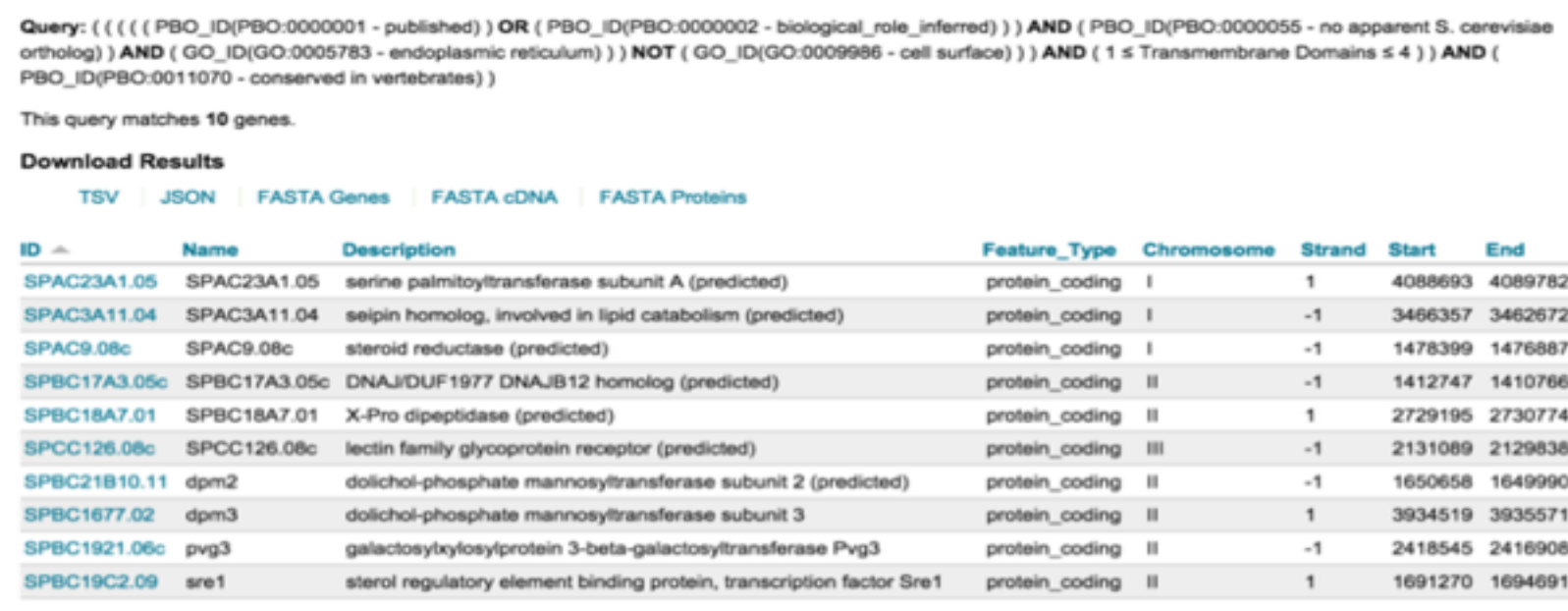
## Classifying Unknowns



Classification unknowns conserved to human

## Predicting processes E.g.1



1. ER localization
2. >1 <4 TM domain
3. Absent from *S. cerevisiae*
4. Conserved in vertebrates

"which known genes best match these profiles?"
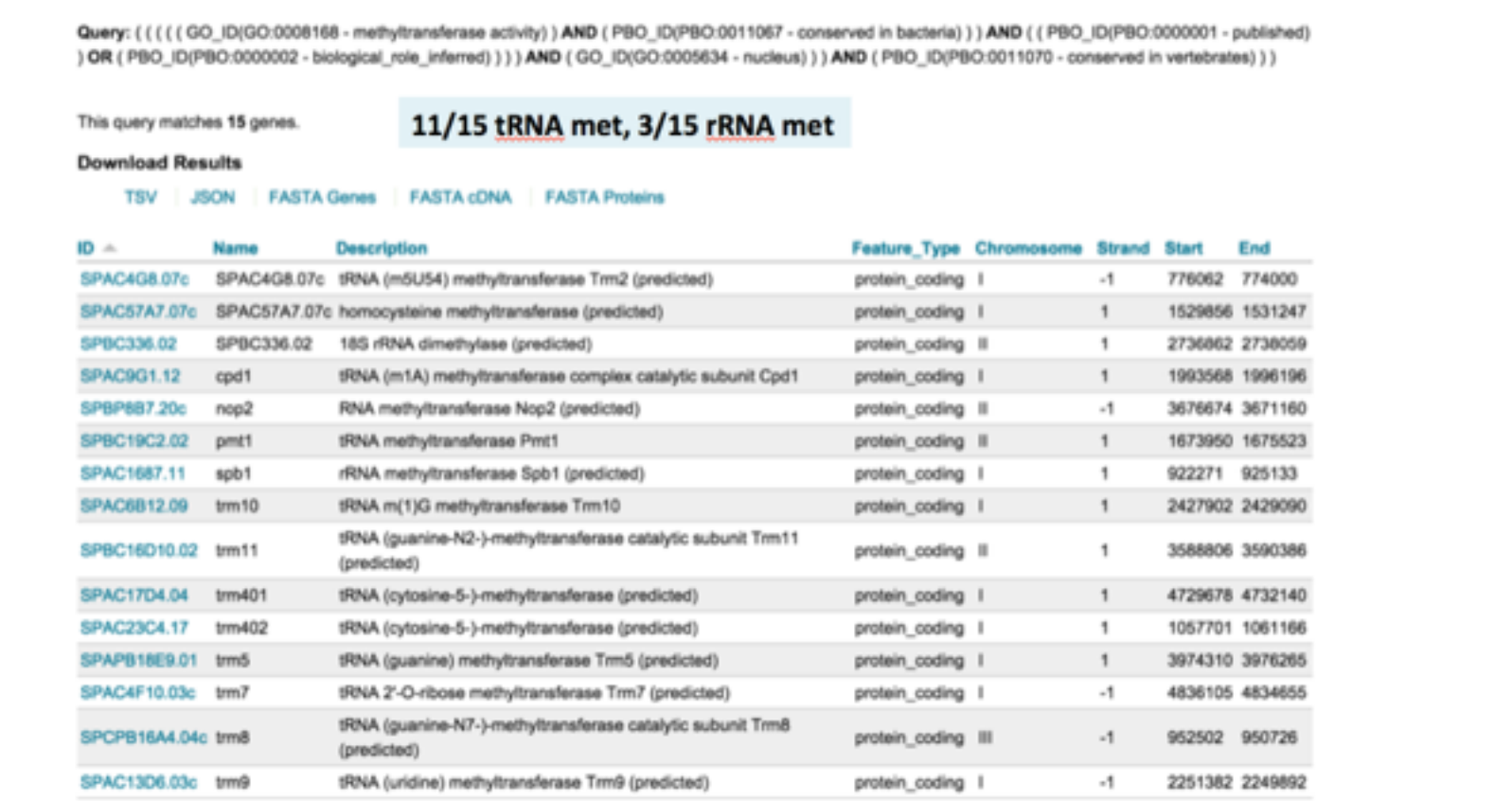
What are these genes enriched for?
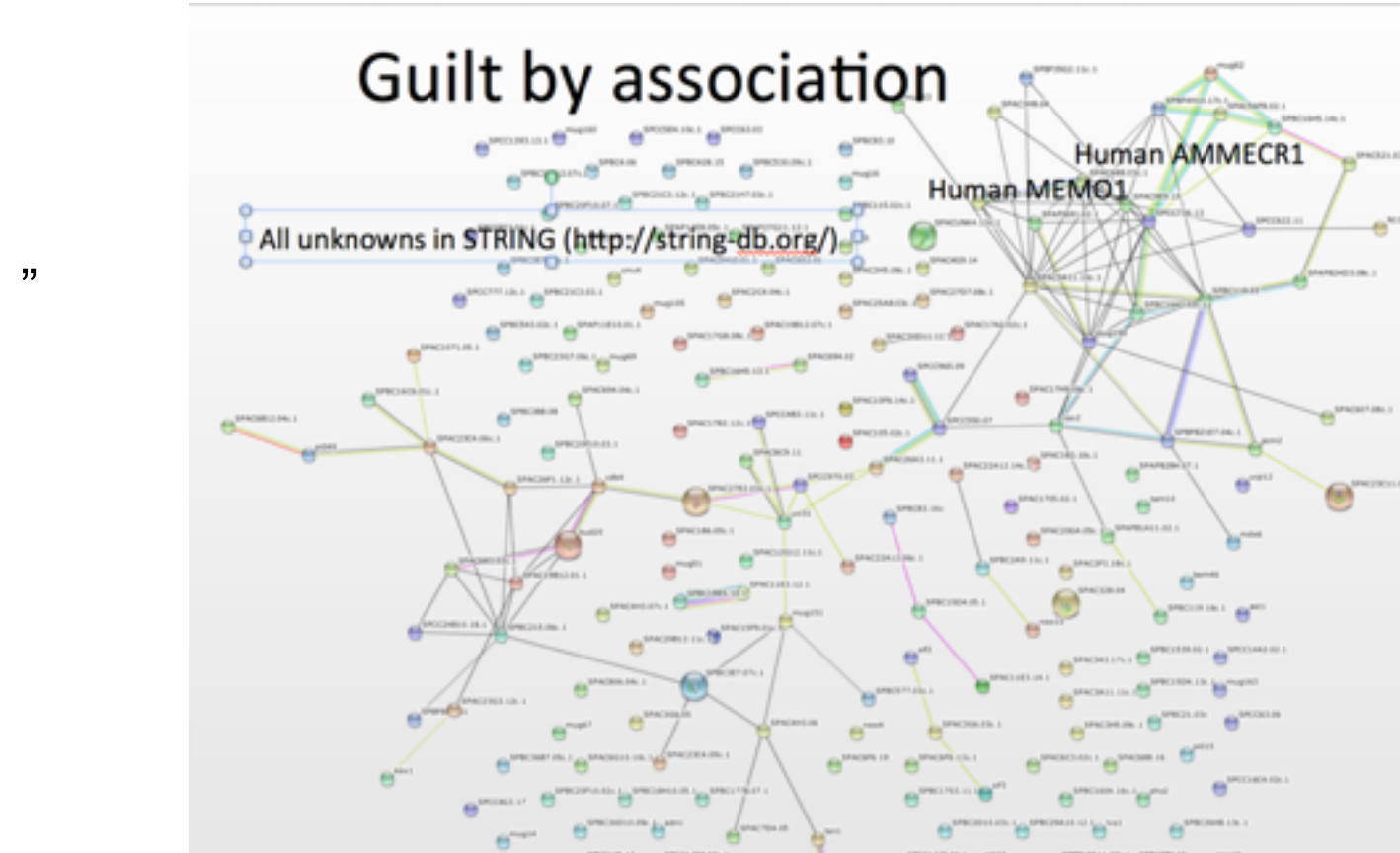
## Predicting processes E.g. 2



1. nuclear localization
2. methyltransferase domain
3. conserved in bacteria
4. Conserved in vertebrates

"which known genes best match these profiles?"

11/15 tRNA met, 3/15 rRNA met

## Predicting processes E.g. 3

### Guilt by association

All unknowns in STRING (http://string-db.org/)



Using Angeli http://bahlerweb.cs.ucl.ac.uk/cgi-bin/GLA/GLA_input
AMMECR subnetwork has connections to meiosis, possibly signalling

AUTHOR AFFILIATIONS
1. PomBase, Cambridge University, Tennis Court Rd, Cambridge, CB2 1QW, UK
2. PomBase, University College London, Darwin Building, Gower Street, London WC1E 6BT, UK