# Using computational predictions to improve literature-based Gene Ontology annotations
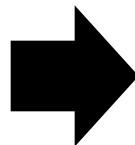
Julie Park, Ph.D.

*Saccharomyces* Genome Database • http://www.yeastgenome.org/
Department of Genetics • Stanford University School of Medicine

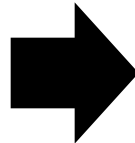# Attaining curation nirvana…

- Curation efficiency

- Annotation consistency

- Data accuracy

# …is not easy!

## Annotation errors

1. Mistakes in capturing the annotation

2. Outdated information

3. Missing annotations



*How can you find these errors?*

# Flavors of GO annotations

## 1. Literature-based – "Manual"
Individually assigned by biocurators based on the published literature

## 2. Computationally-predicted – "Computational"
Automatically generated by *in silico* methods such as protein signatures or computational algorithms

Sources of computational predictions in SGD
InterPro [1]
Swiss-Prot Keywords (SPKW) [2]
YeastFunc [3]
BioPixie [4]

1. Camon, et al (2003) *Genome Res.* 13:662-72
2. http://www.ebi.ac.uk/GOA/Swiss-ProtKeyword2GO.html
3. Tian, et al (2008) *Genome Biol.* 9 Suppl 1:S7
4. Huttenhower and Troyanskaya (2008) *Bioinformatics*. 24:i330-8

# Flavors of GO annotations

## 1. Literature-based – "Manual"
Individually assigned by biocurators based on the published literature

## 2. Computationally-predicted – "Computational"
Automatically generated by *in silico* methods such as protein signatures or computational algorithms
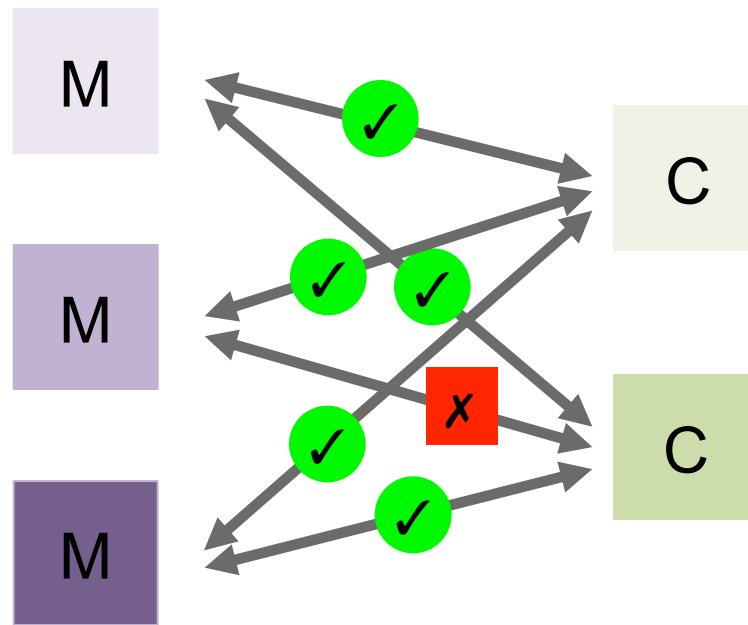
Sources of computational predictions in SGD
InterPro [1]
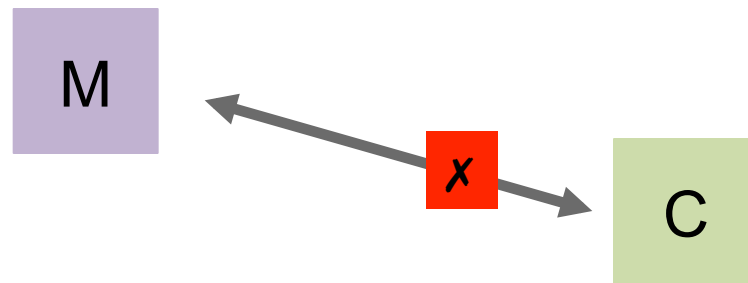Swiss-Prot Keywords (SPKW) [2]
YeastFunc [3]
BioPixie [4]

*Is it possible to take advantage of the strengths of computational predictions and leverage these annotations to improve manual ones?*

# CvManGO:
# Computational vs. Manual GO Annotations

# CvManGO:
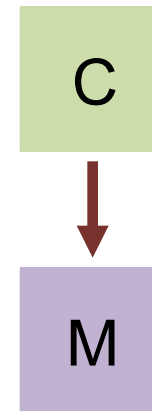# Computational vs. Manual GO Annotations



Do discrepancies between a literature-based annotation and a computational prediction indicate that the manual annotation needs to be updated?
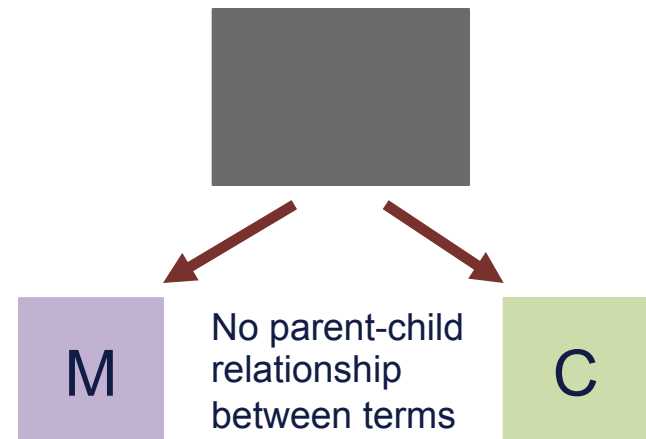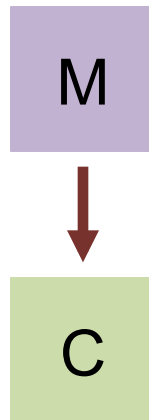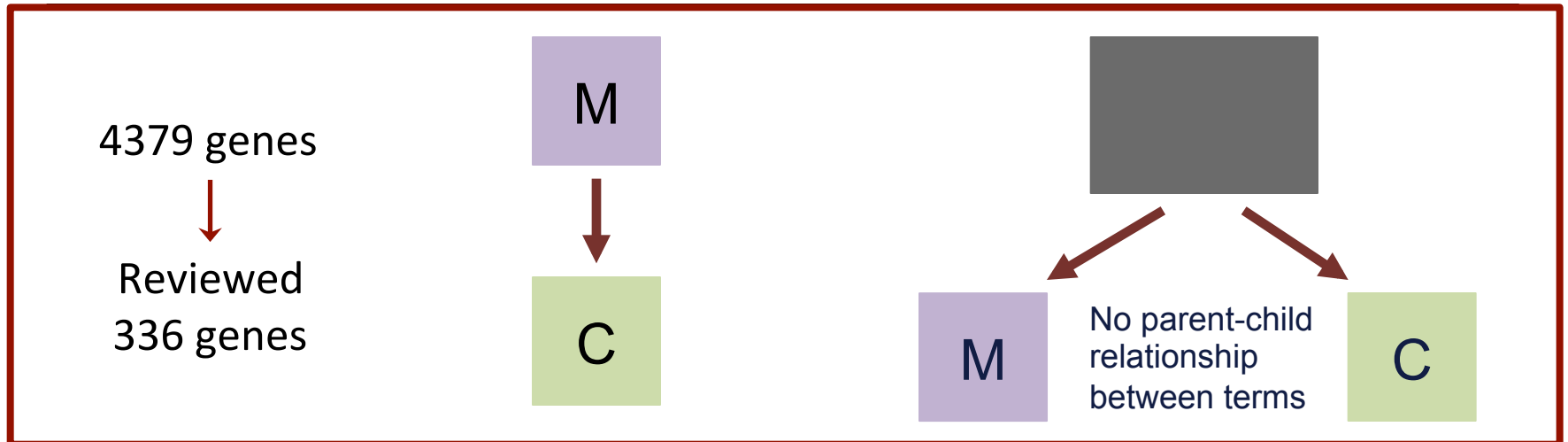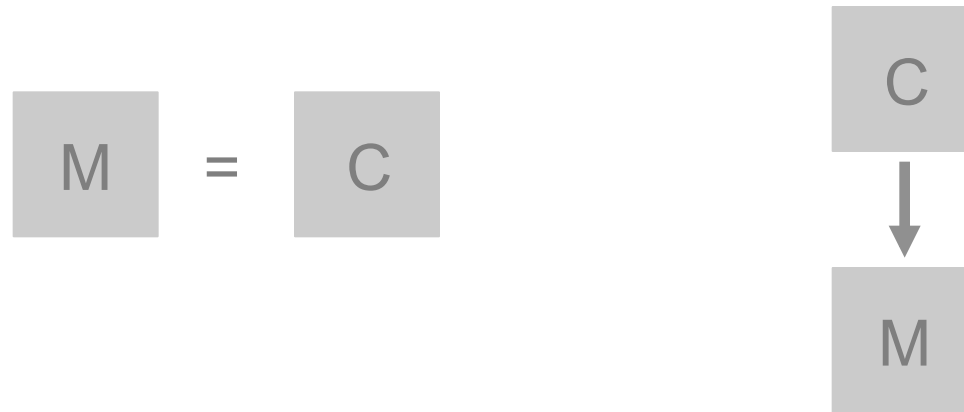
# CvManGO:
# Computational vs. Manual GO Annotations

OK

✓

M = C

C → M

---

Discrepancies

✗

M → C

No parent-child relationship between terms

M     C

# CvManGO:
# Computational vs. Manual GO Annotations



M = C

C
↓
M

4379 genes
↓
Reviewed
336 genes

M
↓
C

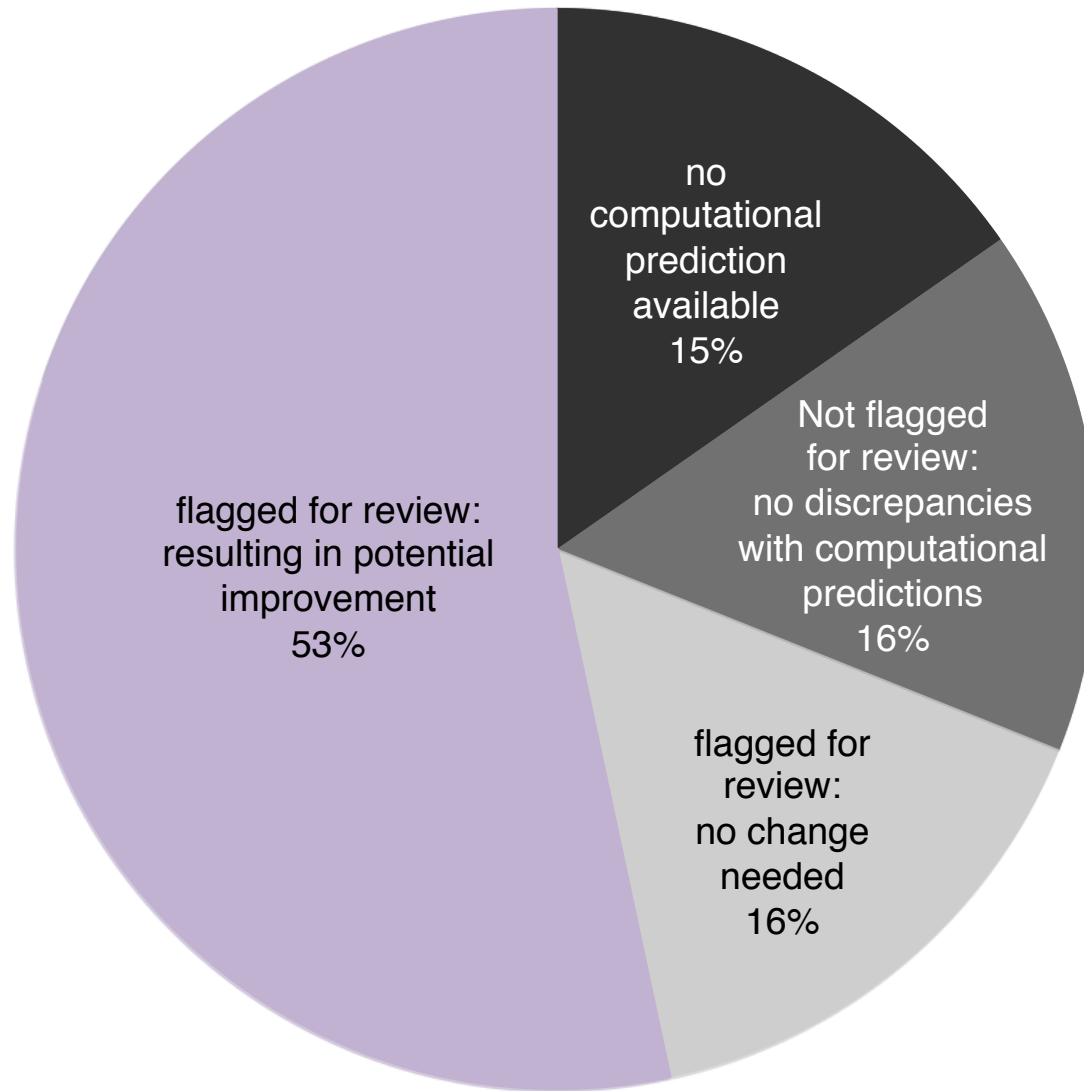M   No parent-child relationship between terms   C

*from October 2009 gene_association.sgd file*
*(6353 total genes)*

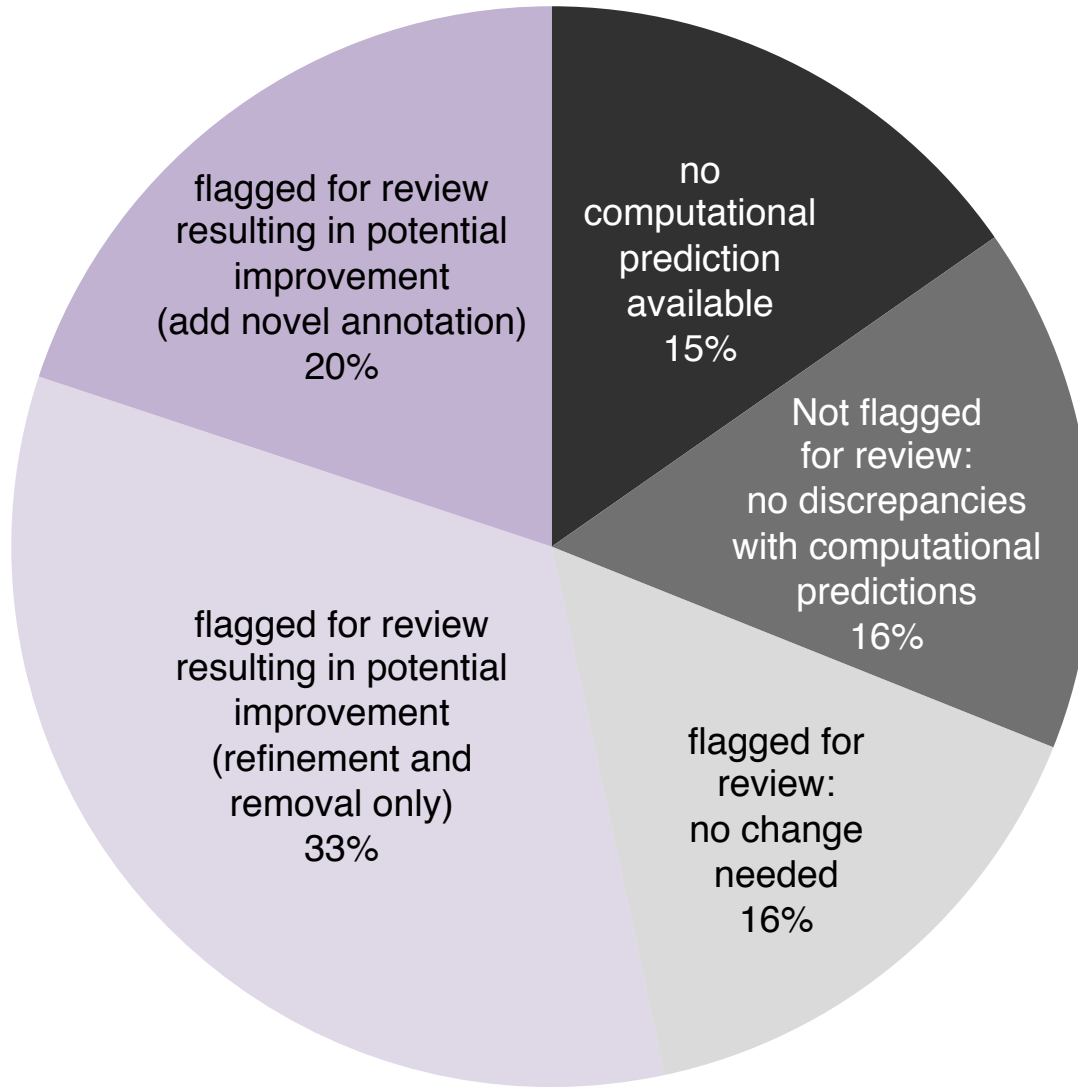# Discrepancies can identify genes that need updating
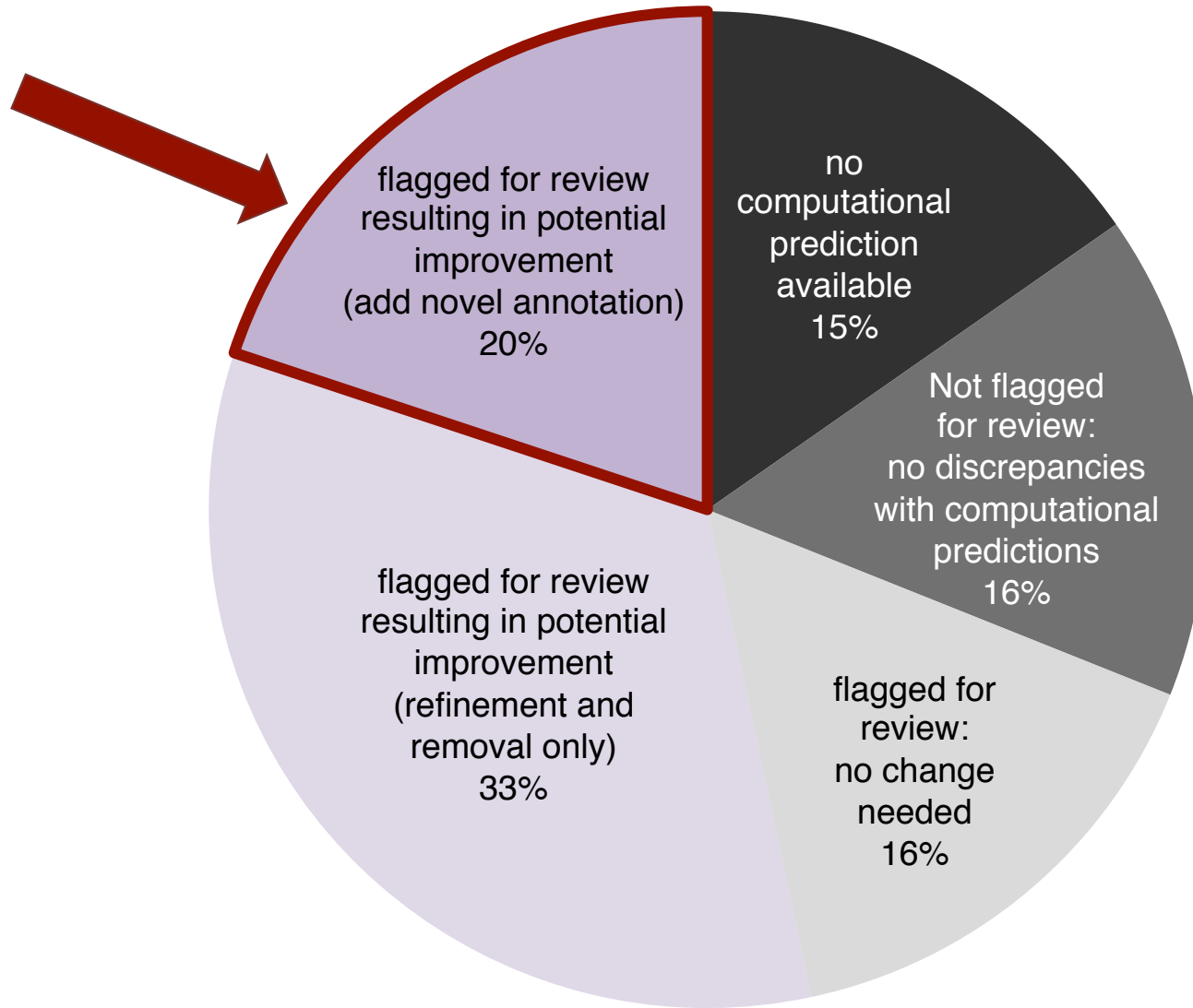
# Extrapolating to the entire genome



Still requires reviewing 4379/6353 genes—can we narrow this down further?

# Factoring in the type of update

# Factoring in the type of update



flagged for review resulting in potential improvement (add novel annotation) 20%

no computational prediction available 15%

Not flagged for review: no discrepancies with computational predictions 16%

flagged for review: no change needed 16%

flagged for review resulting in potential improvement (refinement and removal only) 33%

# Attributes of flagged genes

What are factors that enrich for genes missing annotations?

- Type of discrepancy
- GO aspect
- Amount of literature for a gene
- Source of the computational prediction
- Number of computational sources with discrepancies
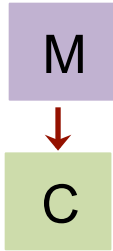
# Attributes of flagged genes

What are factors that enrich for genes missing annotations?
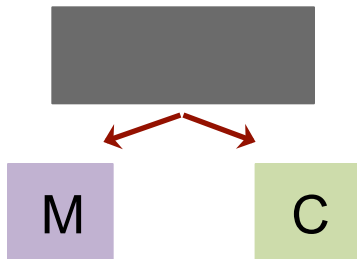
- **Type of discrepancy**
- GO aspect
- Amount of literature for a gene
- Source of the computational prediction
- Number of computational sources with discrepancies

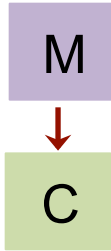# Analysis by Class of Discrepancies

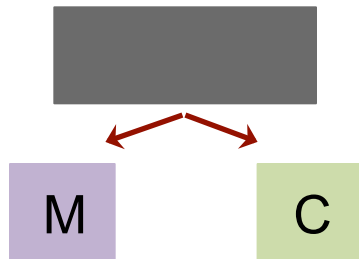M → C    Shallow class

Mismatch class

# Analysis by Class of Discrepancies

% updatable

M

C

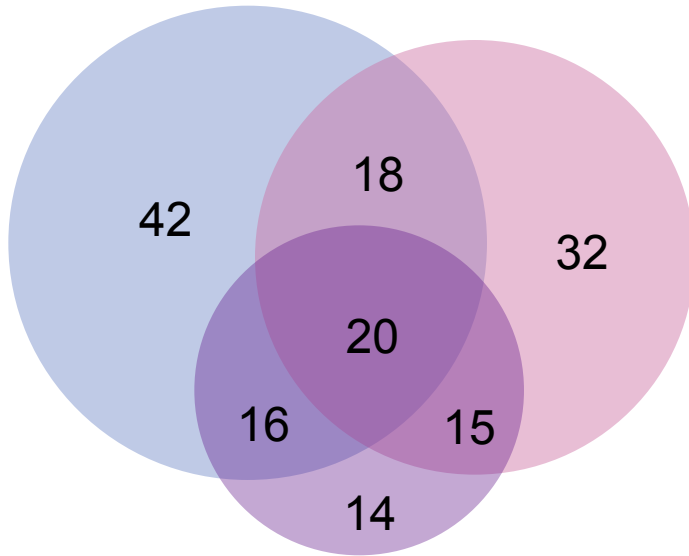Shallow class                    78.8%

M          C

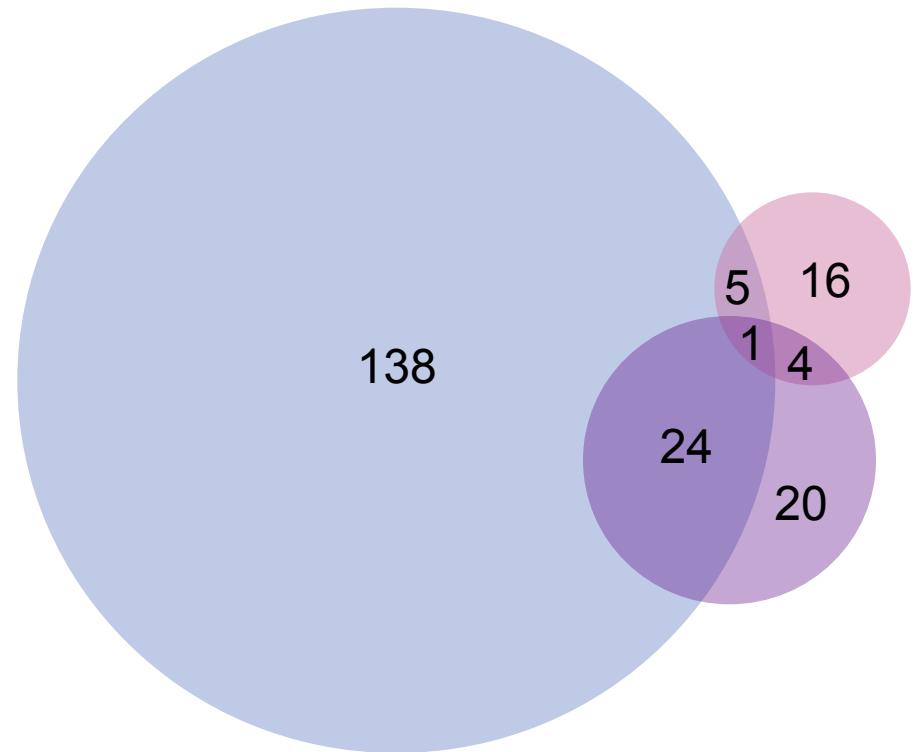Mismatch class                   59.2%

# Types of annotation updates by class



Mismatch

Shallow

Number of genes needing:
- annotation refinement
- annotation removal
- novel annotation addition

# Summary & Conclusions

- Majority of *S. cerevisiae* literature-based GO annotations are good

- Comparing manual vs. computational prediction can identify genes whose annotations need updating

- Additional work needs to be done to pinpoint these annotations and genes

# Summary & Conclusions

- Majority of *S. cerevisiae* literature-based GO annotations are good

- Comparing manual vs. computational prediction can identify genes whose annotations need updating

- Additional work needs to be done to pinpoint these annotations and genes

It works but
there is still work to do!

# Future plans

- Identify predictive features of genes that need updating
  - Are there specific GO terms used for manual curation more likely to be updated?
  - Do specific computational predictions indicate a GO term should be updated?
  - Examine node distance between GO terms used for computational and literature-based annotations
  - Examine contribution of annotation date and new publications
  - A combination of or all of the above?

- Evaluate the accuracy of computational predictions for *S. cerevisiae*

- Expand to evaluate annotations made based on orthology
  - Annotations from GOC PAINT project

- Develop a pipeline for curation prioritization at SGD

- Extend to other annotation projects

# GO Consortium



http://www.geneontology.org/

# *Saccharomyces* Genome Database staff



🐦 @yeastgenome          f http://on.fb.me/ksgskb

✉ sgd-helpdesk@lists.stanford.edu