

# Discovering protein functional sites with unsupervised techniques

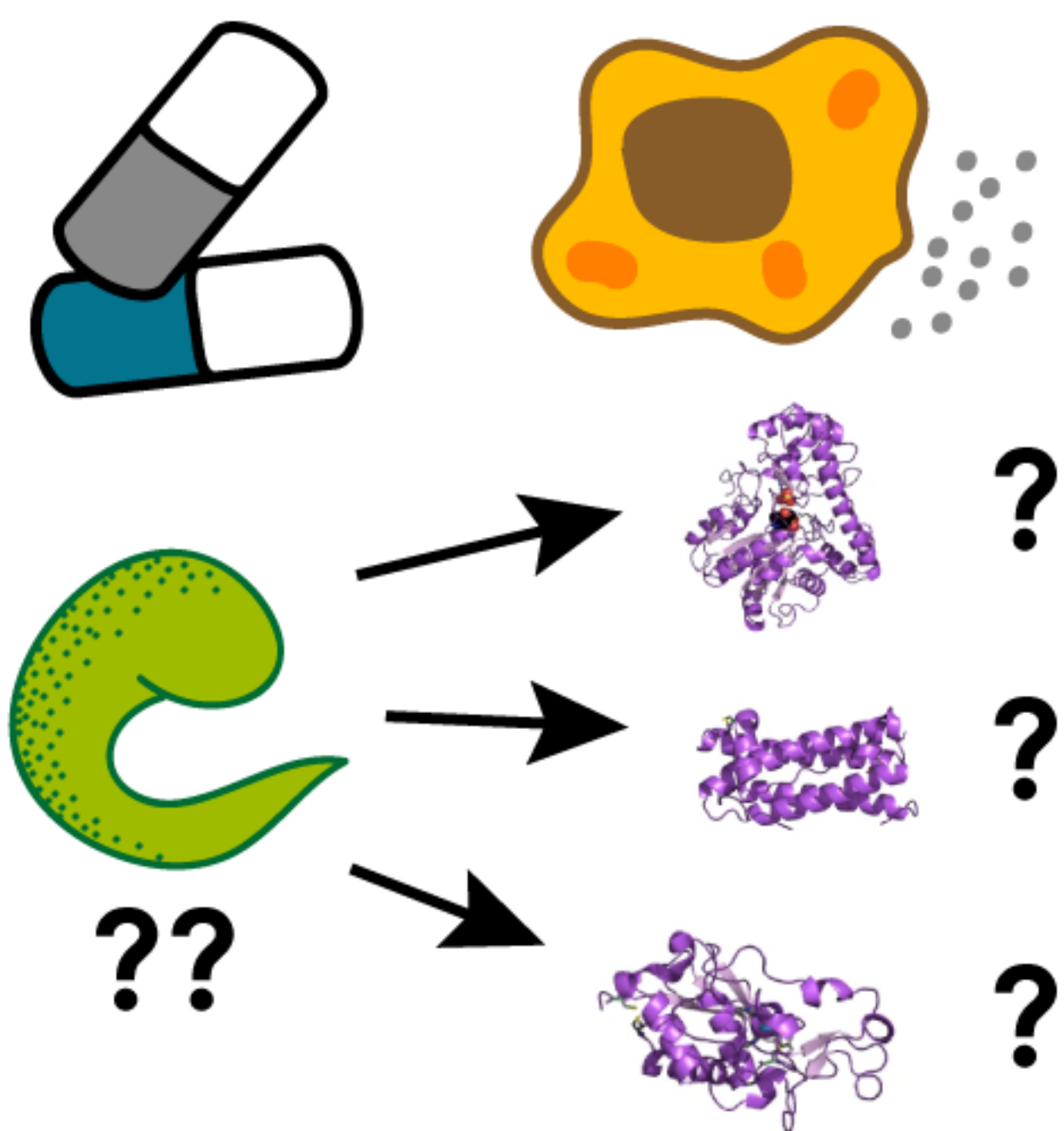
Shirley Wu<sup>1</sup>, Russ Altman<sup>2</sup>

1 Program in Biomedical Informatics, 2 Department of Bioengineering

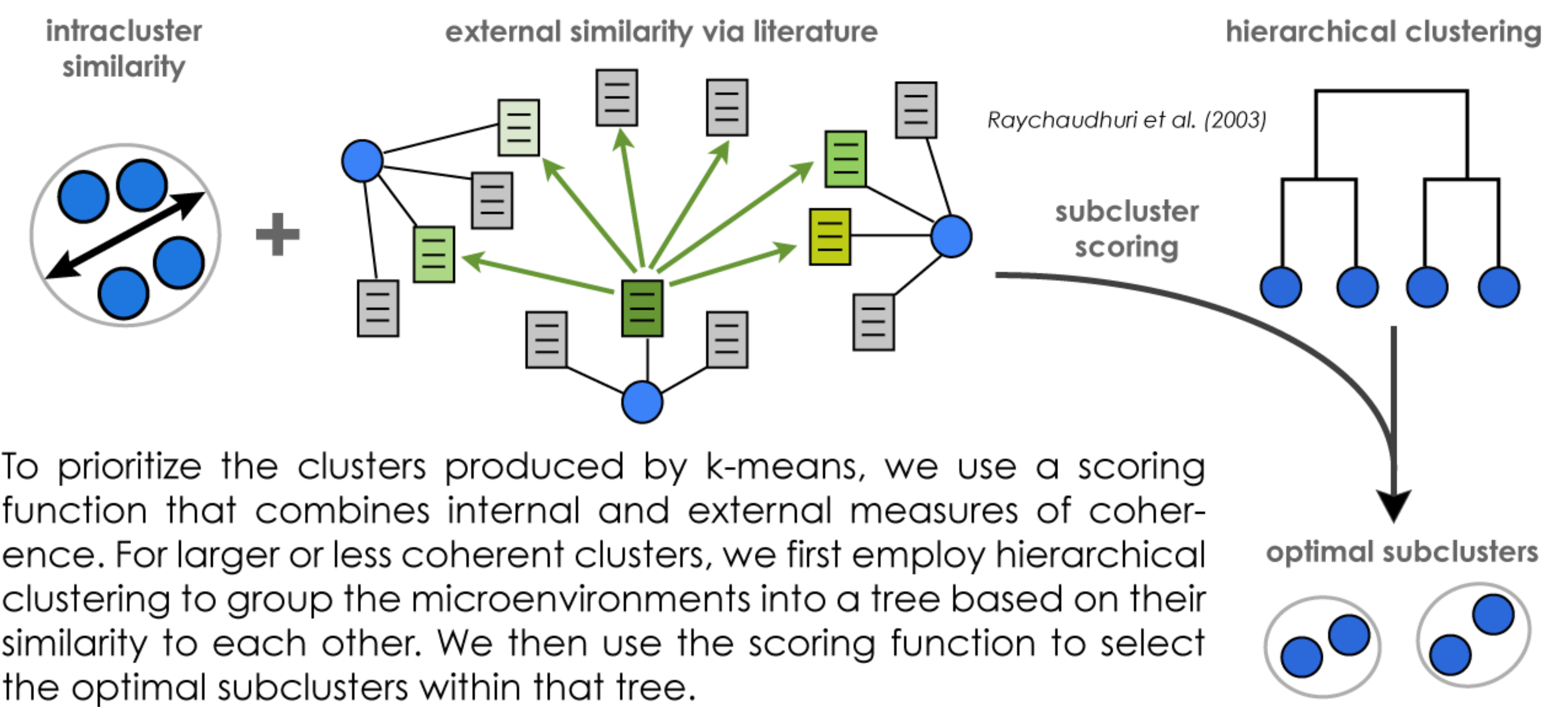
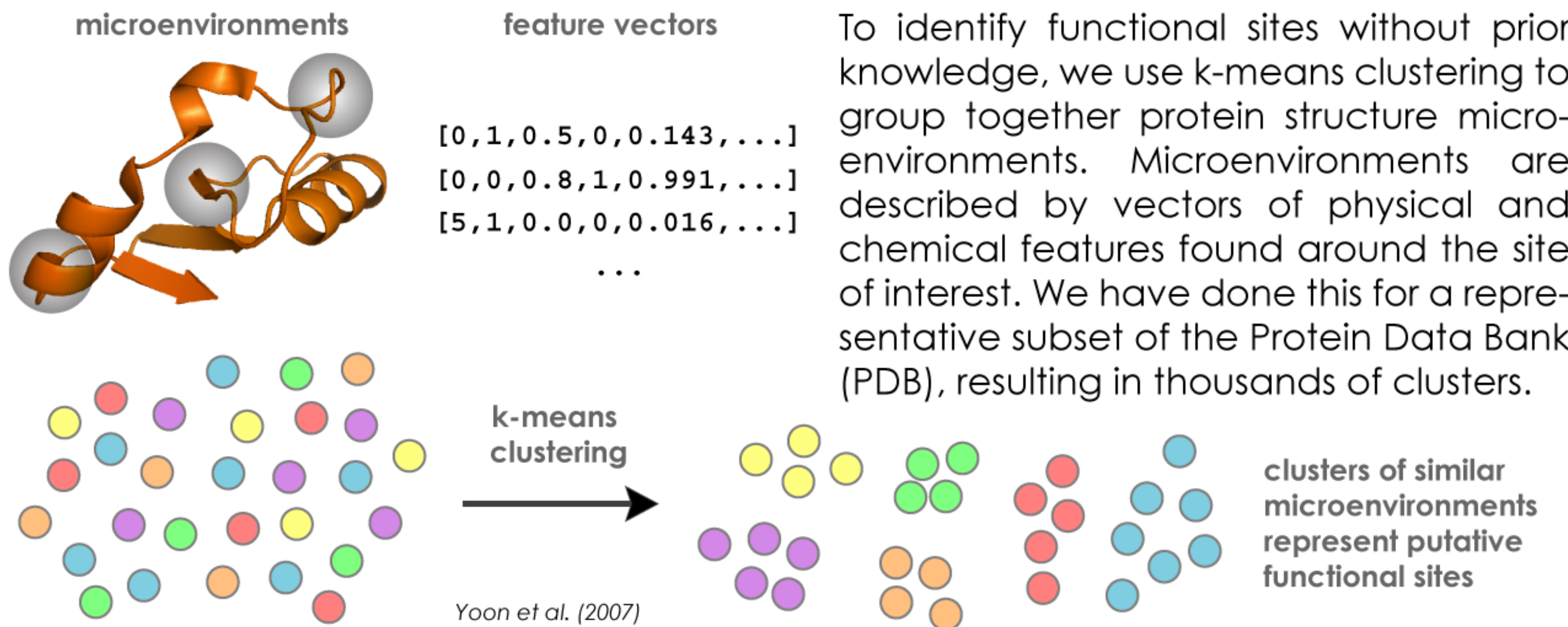
## Background

Characterizing protein function - what reactions they carry out, what molecules they bind, etc - is important for understanding biological processes. We can use this knowledge to engineer therapeutics and other beneficial biology.

Computational methods are fast and inexpensive, allowing high-throughput prediction of protein function. Most methods are supervised approaches, i.e. they use available data about known proteins and functions to make predictions. Large-scale genomics efforts, however, are increasing the rate at which we discover novel organisms and proteins, so we also need methods to identify new biological functions.



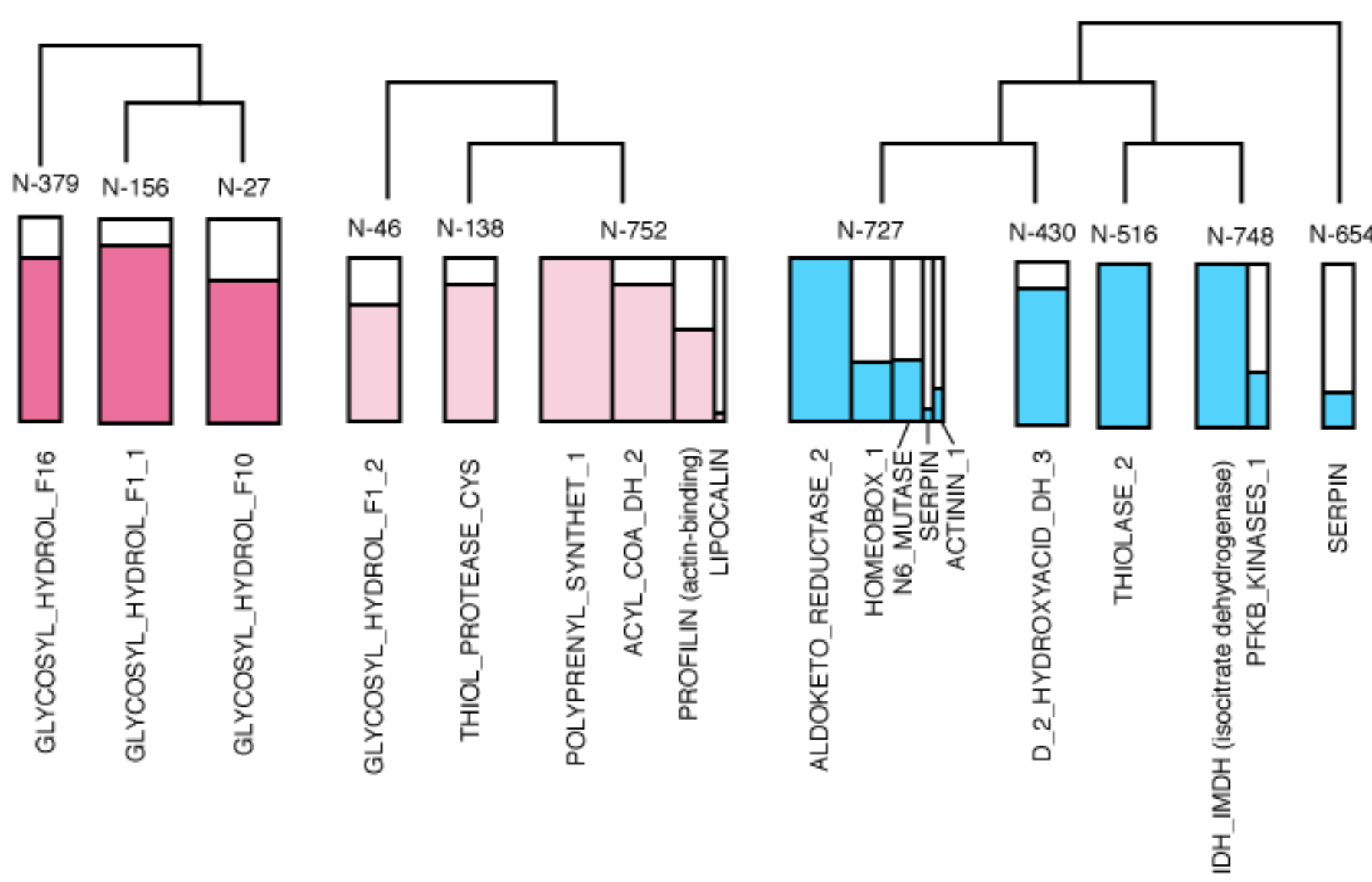
## Methods



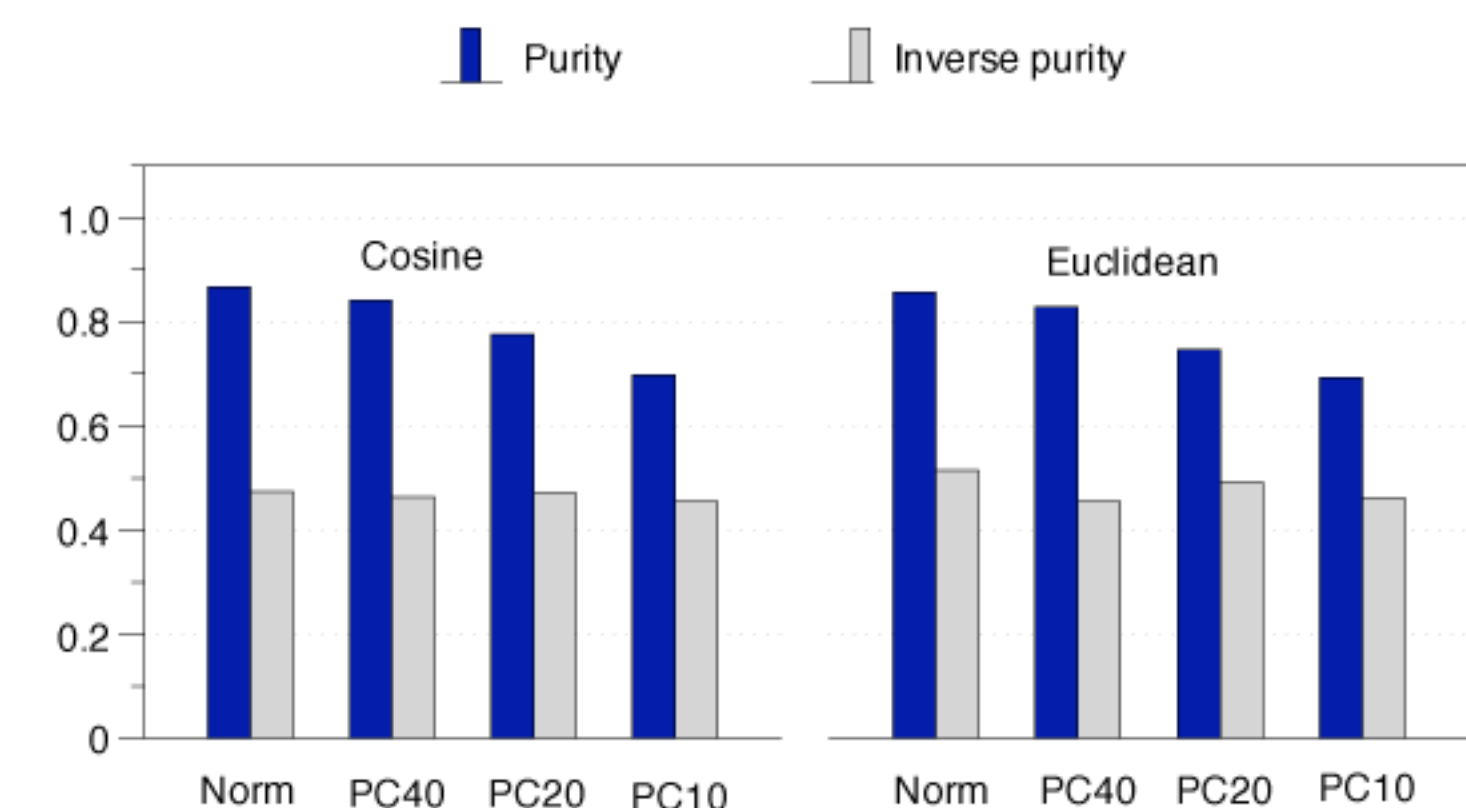
## Evaluation

To evaluate different distance metrics and normalizations of vectors for subcluster selection, we used a set of 1400 microenvironments with assignments to 168 known functions. We found that cosine similarity produced subclusters with slightly better purity with regards to known assignments (external coherence), and better silhouette widths, which balance intracluster tightness and intercluster separation (internal coherence). In addition, reducing the number of features using principal component analysis results in subclusters that are more internally coherent but less externally coherent.

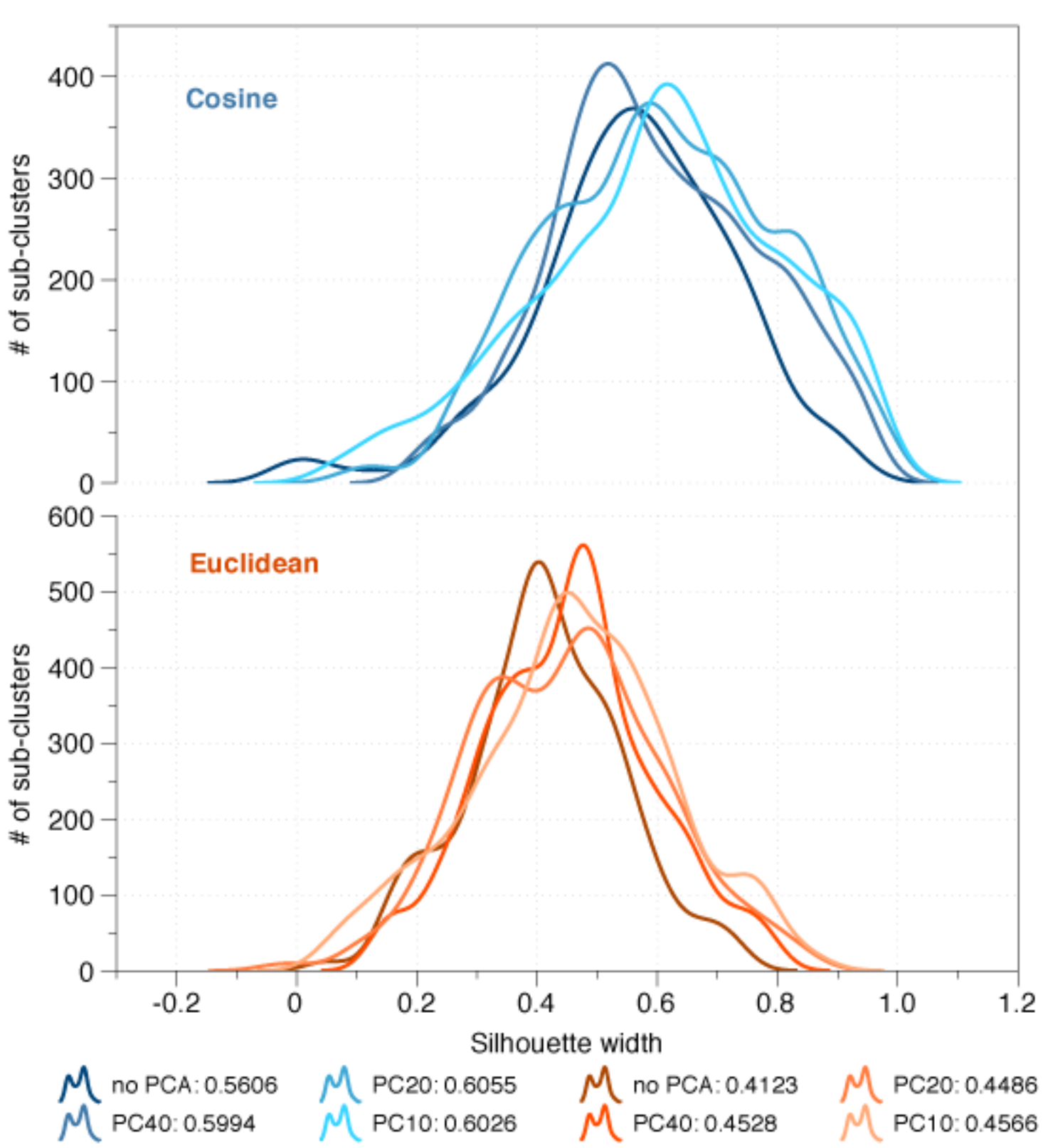
### Examples of subclusters with known functional assignments



### Purity and inverse purity of validation set



### Distribution of silhouette values for validation set



## Application

We are currently applying the subcluster selection strategy to the whole-PDB k-means clustering to determine compelling candidates for further analysis. We use a number of term enrichment methods to gain insight into the possible biological role of the microenvironment represented by each candidate subcluster.

### Cluster 257: 30 proteins



### MeSH terms

Insulin  
Hydrogen Bonding  
Ribosomal Proteins  
Peptides  
Pancreas  
Amino Acids  
Chymotrypsin  
Electrophoresis  
Protein Folding

### Raw text terms

structur monomer  
sequenc c-peptid  
monomer insulin  
conform insulin  
hexam crystal  
2zn insulin  
protein-protein  
coordin zn  
zn atom

Boyle et al. (2004)

### Gene Ontology terms

hormone activity  
glucose metabolic process  
receptor binding  
hexose metabolic process  
insulin receptor binding  
monosaccharide metabolic process  
negative regulation of catabolic process  
positive regulation of cytokine secretion  
insulin-like growth factor binding

## References

Yoon S, Ebert JC, Chung EY, De Micheli G, Altman RB. (2007) BMC Bioinformatics, 8:Suppl 4:S10.  
Raychaudhuri S, Chang J, Imam F, Altman RB. (2003) Nucleic Acids Res 31(15):4553-60.  
Boyle EI, Weng S, Gollub J, Jin H, Botstein D, Cherry JM, Sherlock G. (2004) Bioinformatics 20(18):3710-5.