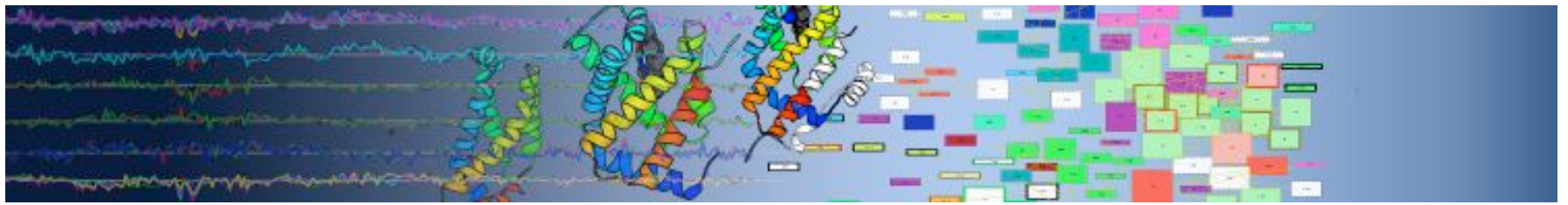# Protein Function Prediction Using Structural Homology

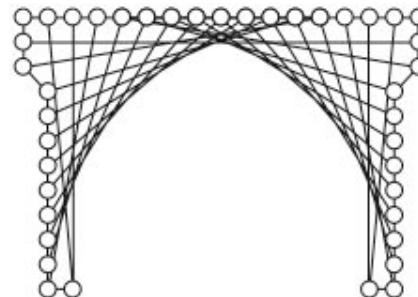## Kevin Drew,

Lars Malmstroem,
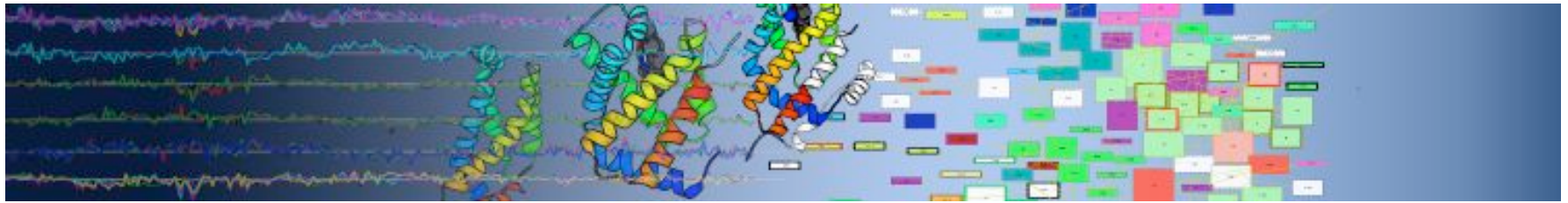
Glenn Butterfoss,

Richard Bonneau

CENTER FOR GENOMICS
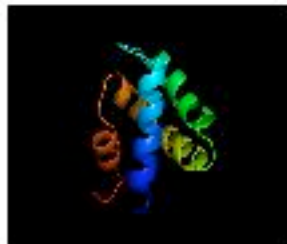AND SYSTEMS BIOLOGY
NEW YORK UNIVERSITY

COURANT INSTITUTE

# Genome Annotation using Homology
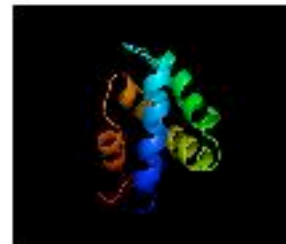


Sequence Homology  40%-60% Coverage

Score = 107 bits (264), Expect = 6e-23
Identities = 63/160 (39%), Positives = 97/160 (60%), Gaps = 7/160 (4%)

```
Query:1  MSVMYKKILYPTDFSETAEIALKHVKTLKAEEVILLDEREIKKRDIFSLLLGVA  60
         M M++K+L+PTDFSE A  A++  + ++  EVILLDE   +++     L+ G +
Sbjct:1  MIFMFRKVLFPTDFSEGAYRAVEVFEKMEVGEVILLDEGTLEE-----LMDGYS  55
...
```
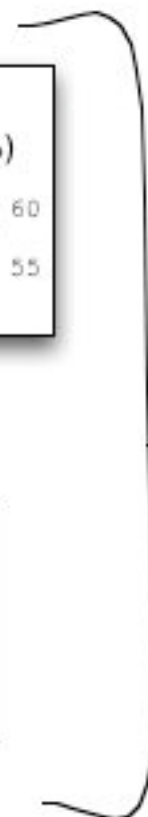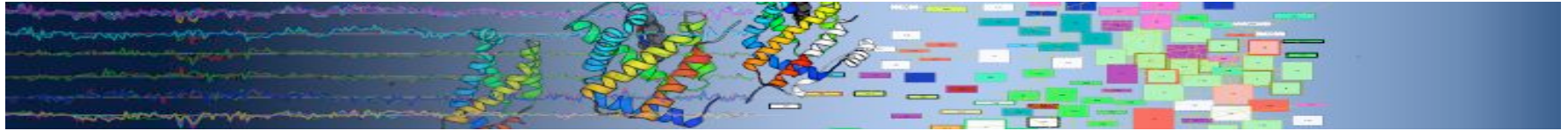
Structural Homology:

=

Function Annotations

# Structural Homology: Example

Bacteriocin AS-48, Casp 4

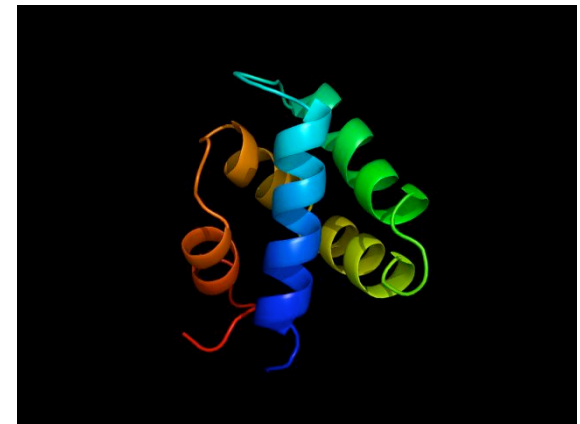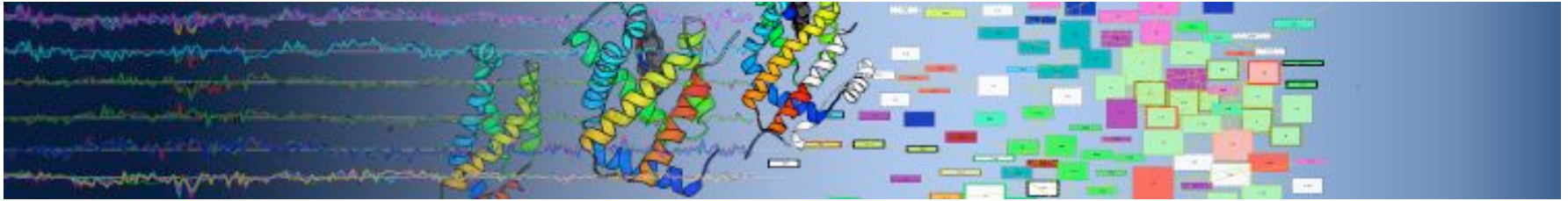|  | 1E68 |  | 1NKL |
|---|---|---|---|
| Sequence: | MAKEFGIPAAVAGTVLNVVEAGGW VTTIVSILTAVGSGGLSLLAAAGRES IKAYLKKEIKKKGKRAVIAW | 4%= | GYFCESCRKIIQKLEDMVGPQPNEDTVTQAAS QVCDKLKILRGLCKKIMRSFLRRISWDILTGKKP QAICVDIKICKE |

Structure:         =

Function:    Cyclic Bacterial Lysin     =       NK Lysin    3

Bonneau, R., Tsai, J., Ruczinski, I., Baker, D. Functional Inferences from Blind ab Initio Protein Structure Predictions. J. Structural Biology. (2001)

# Gene Ontology (GO)



Molecular Function GO:0003674

Binding GO:0005488

Protein Binding GO:0005515

Clathrin Binding GO:0030276

specificity
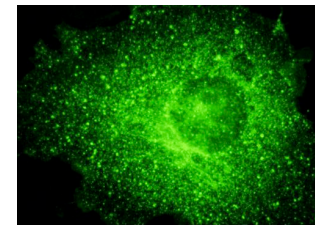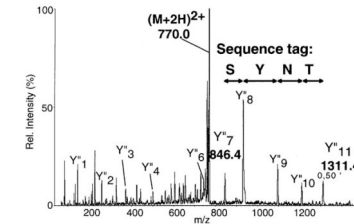
Molecular Function GO Graph

4

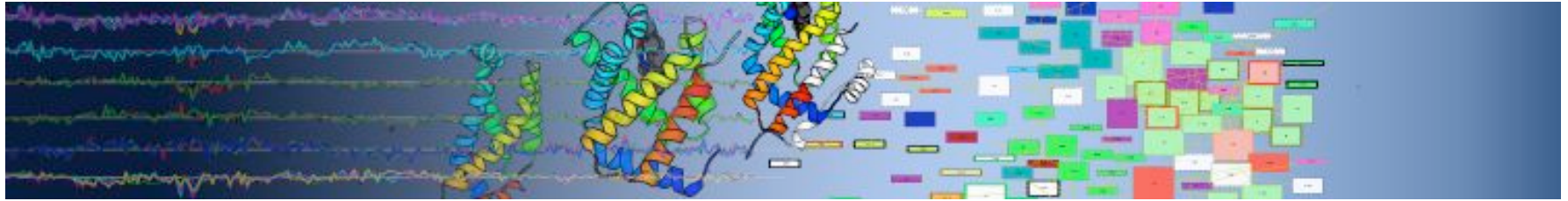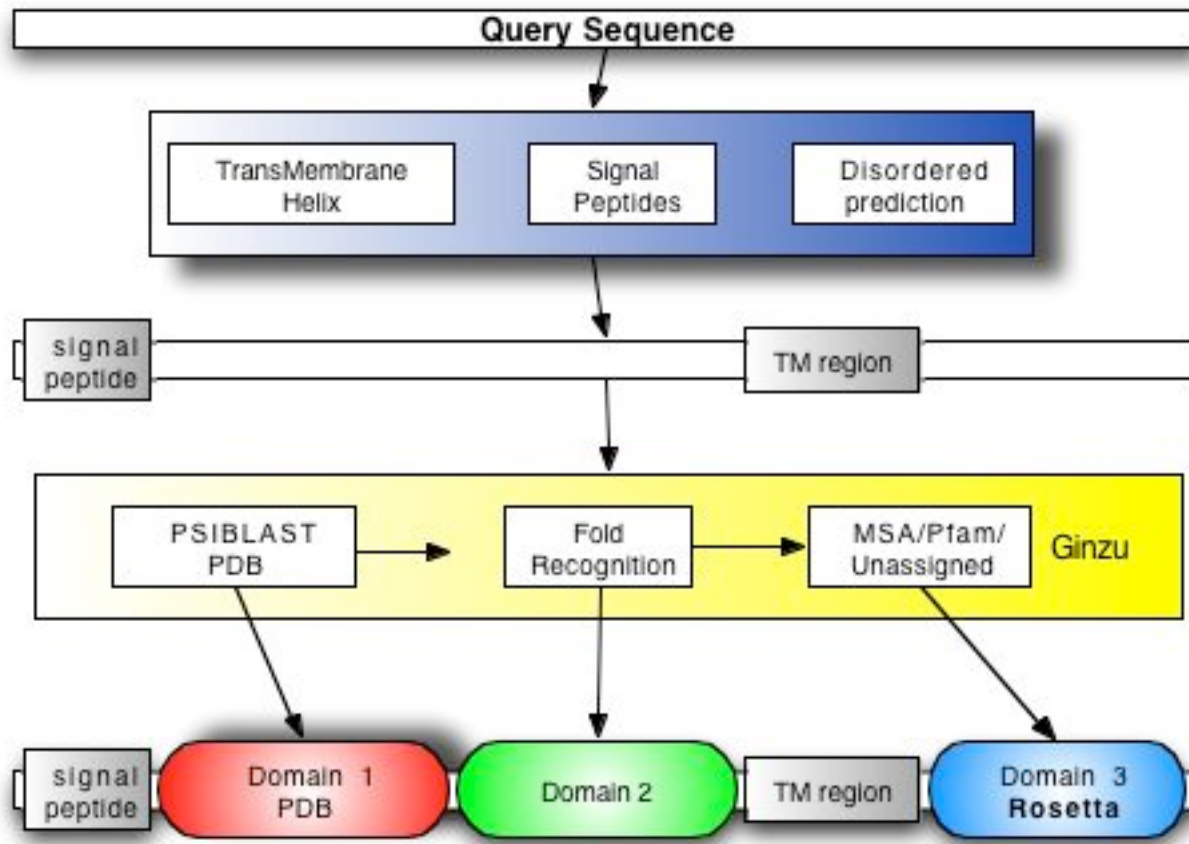# Additional Evidence of Function for Integration with Structure

- GO Biological Process

- GO Cellular Component

- Experimental Data

    - Mass Spec Pull Down

    - Fluorescent Localization

- Generally boosts confidence of predictions

# Protein Domain Prediction

# Function Prediction Overview

# Function Prediction Overview

# Function Prediction Overview

# Function Prediction Overview

# Matching Predicted Structures to Known Structures



$$\log\left(\frac{P_{MCM}}{1 - P_{MCM}}\right) = a \cdot zscore + b \cdot CO + c \cdot converg + d \cdot \left|\log\left(\frac{L_{Astral}}{L_{predicted}}\right)\right| + C$$

$$P(Structure) = 1 - \prod_{k=1}^{n}(1 - p_k)$$

11

# Training Data Derived from GO and Known Structures



Predicted Structures

Known Structures

Gene Ontology Terms

P(Structure)

P(Function|Structure)

P(Function)

Training Set Reduction

GO: 1.6 million sequences

GO + Astral
Blast-hits: 643,173

Cluster Centers
280,511

Removal
of Benchmark
< 280,511

# Naïve Bayes

- In words: what is the probability of a variable, y, is true given features, **x**, over the probability y is false given the features **x**.
  - Take the log and if its >0 its more likely to be true than false.

- y = molecular function and **x** = {sf, bp, cc}

$$LL_X = log\left(\frac{P(y = TRUE)}{P(y = FALSE)}\right) + \sum_{j=1}^{d} log\left(\frac{P(x_j|y = TRUE)}{P(x_j|y = FALSE)}\right)$$

# Full Function Prediction Formula

$$LL_X = log\left(\frac{P(y = TRUE)}{P(y = FALSE)}\right) + \sum_{j=1}^{d} log\left(\frac{P(x_j|y = TRUE)}{P(x_j|y = FALSE)}\right)$$

Naive Bayes

$$LL_{PLS} = log\left(\frac{P(Function)}{P(\overline{Function})}\right) + \sum_{i=1}^{N}\left[P(Structure_i) * log\left(\frac{P(Structure_i|Function)}{P(Structure_i|\overline{Function})}\right)\right] + \sum_{j=P,L} log\left(\frac{P(x_j|Function)}{P(x_j|\overline{Function})}\right)$$

# Full Function Prediction Formula

$$LL_X = log\left(\frac{P(y = TRUE)}{P(y = FALSE)}\right) + \sum_{j=1}^{d} log\left(\frac{P(x_j|y = TRUE)}{P(x_j|y = FALSE)}\right)$$

Naive Bayes

Structure Contribution

$$LL_{PLS} = log\left(\frac{P(Function)}{P(\overline{Function})}\right) + \sum_{i=1}^{N}\left[P(Structure_i) * log\left(\frac{P(Structure_i|Function)}{P(Structure_i|\overline{Function})}\right)\right] + \sum_{j=P,L} log\left(\frac{P(x_j|Function)}{P(x_j|\overline{Function})}\right)$$



Predicted Structures    Known Structures    Gene Ontology Terms

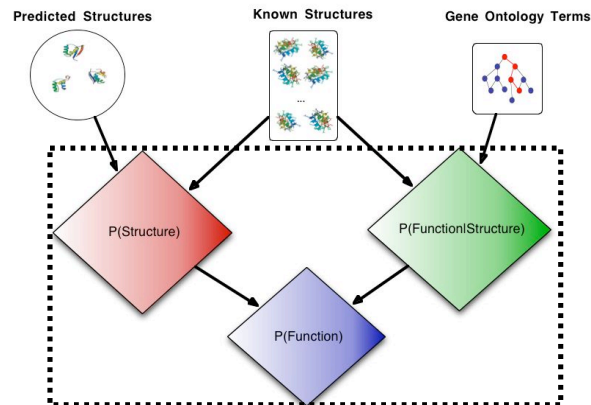P(Structure)    P(Function|Structure)    P(Function)

# Full Function Prediction Formula

$$LL_X = log\left(\frac{P(y = TRUE)}{P(y = FALSE)}\right) + \sum_{j=1}^{d} log\left(\frac{P(x_j|y = TRUE)}{P(x_j|y = FALSE)}\right)$$

Naive Bayes

$$LL_{PLS} = log\left(\frac{P(Function)}{P(\overline{Function})}\right) + \sum_{i=1}^{N}\left[P(Structure_i) * log\left(\frac{P(Structure_i|Function)}{P(Structure_i|\overline{Function})}\right)\right] + \sum_{j=P,L} log\left(\frac{P(x_j|Function)}{P(x_j|\overline{Function})}\right)$$

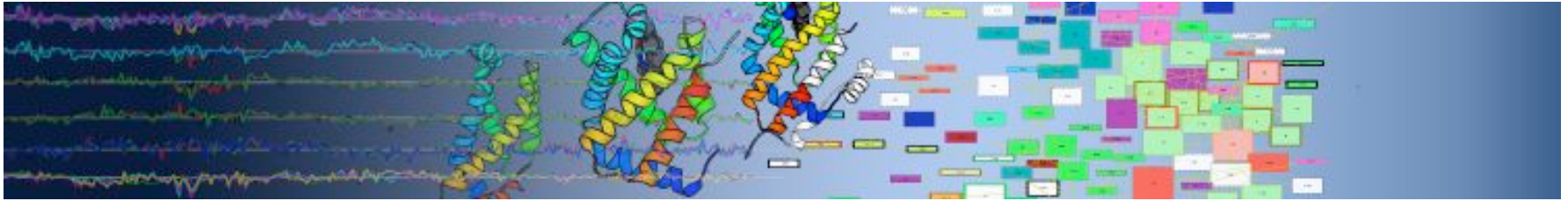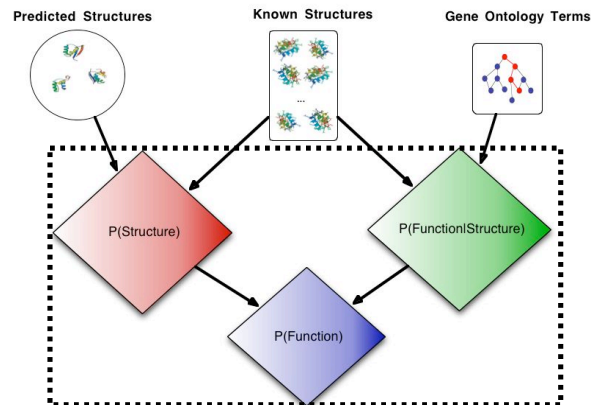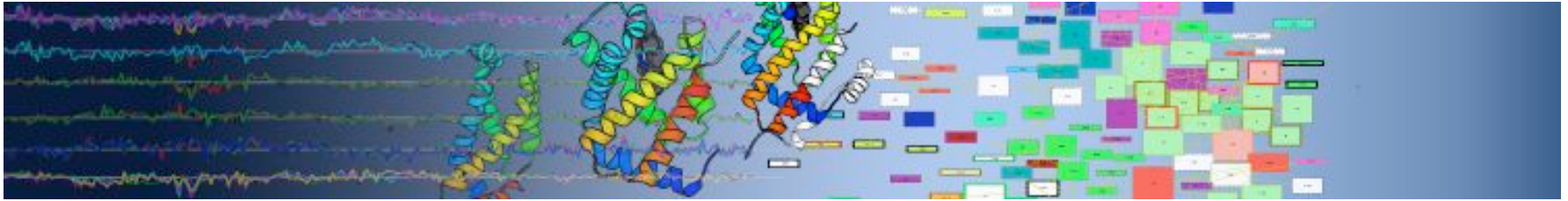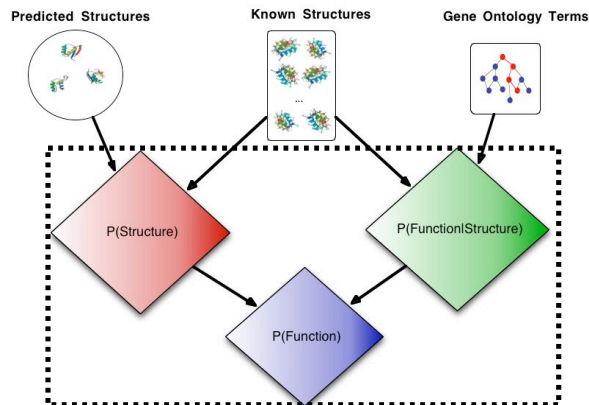# Full Function Prediction Formula

$$LL_X = log\left(\frac{P(y=TRUE)}{P(y=FALSE)}\right) + \sum_{j=1}^{d} log\left(\frac{P(x_j|y=TRUE)}{P(x_j|y=FALSE)}\right)$$

Naive Bayes

Prior

Additional Evidence

$$LL_{PLS} = log\left(\frac{P(Function)}{P(\overline{Function})}\right) + \sum_{i=1}^{N}\left[P(Structure_i)*log\left(\frac{P(Structure_i|Function)}{P(Structure_i|\overline{Function})}\right)\right] + \sum_{j=P,L} log\left(\frac{P(x_j|Function)}{P(x_j|\overline{Function})}\right)$$
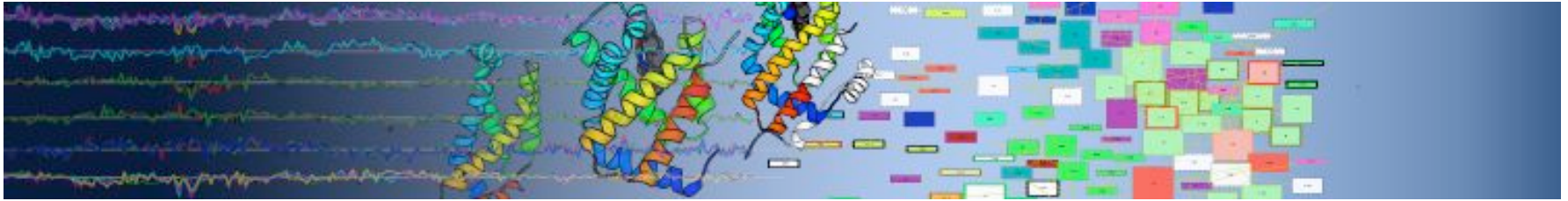
Predicted Structures    Known Structures    Gene Ontology Terms

P(Structure)    P(Function|Structure)

P(Function)

17

# Results: Solved Structures

How accurate are we when we predict SCOP Superfamily for PDB Structures?



histogram of scop_benchmark : 565 true / 988 total

grey = all mcm scores, cadetblue = correct based on since solved, KS-test: D= 0.5 p-value= 0e+00

# How accurate are we when we predict Structure for Swissprot Proteins?

Histogram of swissprot_benchmark : 3709 true / 6143 total

Grey Bar = total number of structure predictions
Green Bar = number of correct structure predictions
% above bar = percent correct

**Low Confidence**

**High Confidence**

Grey = all mcm scores, seagreen = correct based on since solved, KS-test: D=0.67 p-value= 0e+00

MCM Score

Count

How do we measure preformance of function predictions?

Accuracy - # of Correct / # of Total

Coverage - # of Proteins with High Confidence Predictions

Specificity and Uniqueness - # of proteins with annotations
(background probability)

# How accurate are our function predictions using structure only?

**Histogram of Function Prediction for swissprot_benchmark : s predictors**

3083 Domains

Grey Bar = total number of function predictions
Green Bar = number of correct functions predictions
% above bar = percent correct

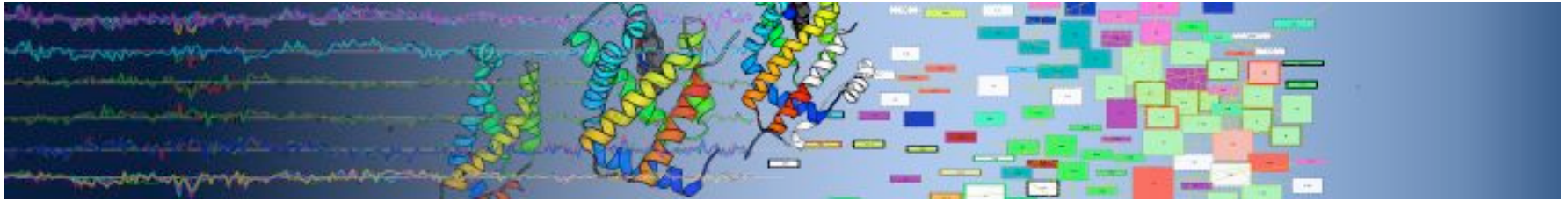0.24  0.26  0.35  0.31  0.45  0.61  0.44  0.57  0.82  0.73  0.52  0.80  0.38  0.25  0.80  0.97  0.89  0.71

Frequency
1500
1000
500
0

0    2    4    6    8

**Low Confidence**          Log Likelihood Ratio          **High Confidence**

# How accurate are our function predictions using GO process & structure?



Histogram of Function Prediction for swissprot_benchmark : ps predictors

# What does structure provide over GO process alone?
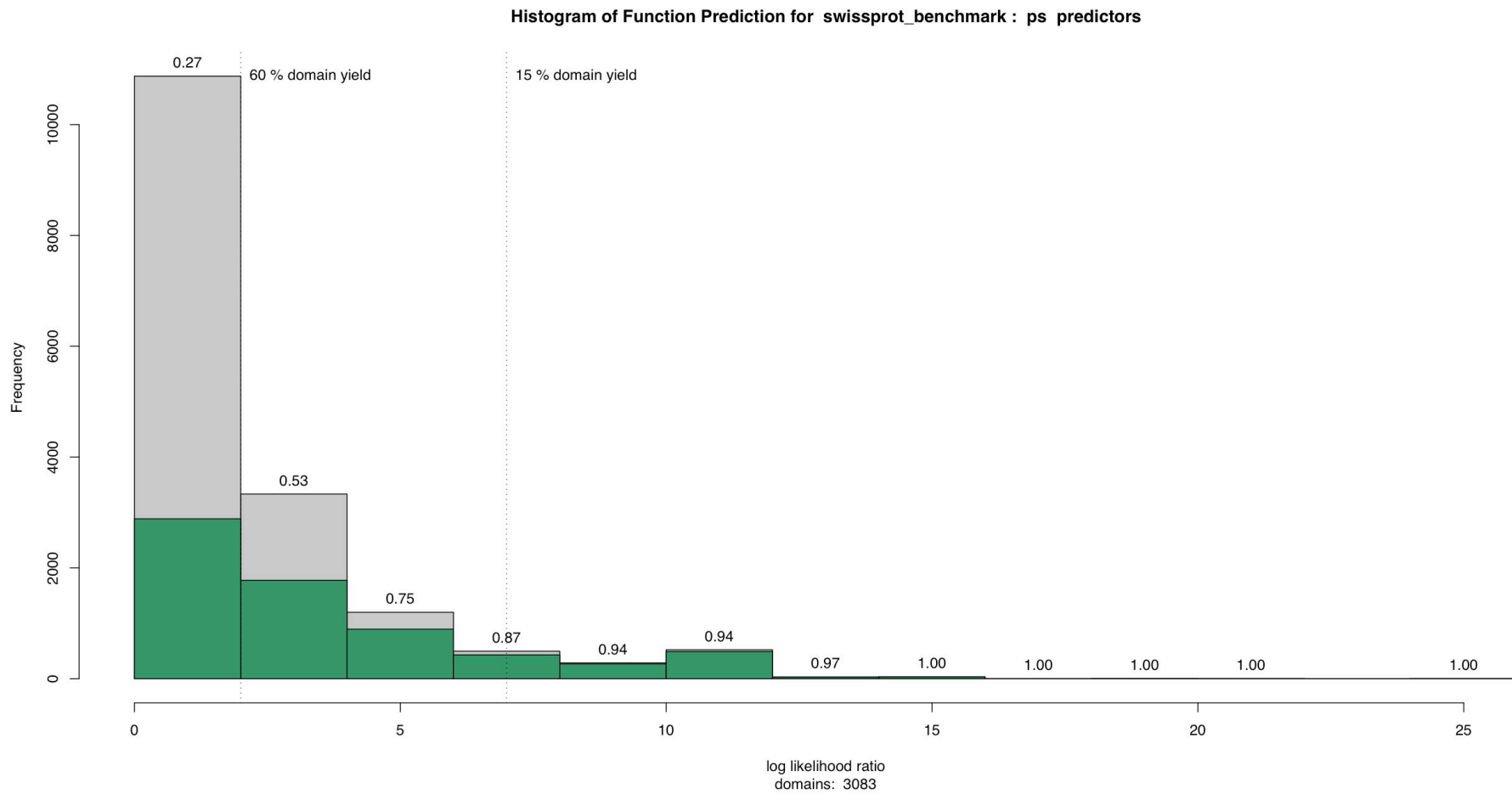


Histogram of Function Prediction for swissprot_benchmark : p predictors

# What does structure provide over GO process alone?

Protein Coverage: PS vs P

Process (orange)
Process & Structure (green)
# = number of domains

# Uniqueness and Specificity of GO Functions

## Unique Functions by Evidence



Localization: 27
PL: 0
LS: 5
PLS: 7
Process: 64
PS: 32
Structure: 90

Swissprot LLR >= 2

## Sampling of Predicted Terms

| GO ID | GO Name | Percent of Genes with Terms |
|---|---|---|
| GO:0005198 | structural molecule activity | 0.03 |
| GO:0003735 | structural constituent of ribosome | 0.02 |
| GO:0003676 | nucleic acid binding | 0.17 |
| GO:0003723 | RNA binding | 0.04 |
| GO:0016491 | oxidoreductase activity | 0.16 |
| GO:0046872 | metal ion binding | 0.11 |
| GO:0016787 | hydrolase activity | 0.24 |
| GO:0043167 | ion binding | 0.12 |
| GO:0043169 | cation binding | 0.11 |
| GO:0005509 | calcium ion binding | 0.01 |

. . .

| GO ID | GO Name | Percent of Genes with Terms |
|---|---|---|
| GO:0004550 | nucleoside diphosphate kinase activity | 0.0009 |
| GO:0005496 | steroid binding | 0.001 |
| GO:0042379 | chemokine receptor binding | 0.0006 |
| GO:0030234 | enzyme regulator activity | 0.01 |
| GO:0016788 | hydrolase activity, acting on ester bonds | 0.04 |
| GO:0008289 | lipid binding | 0.005 |
| GO:0004812 | aminoacyl-tRNA ligase activity | 0.01 |
| GO:0005506 | iron ion binding | 0.03 |
| GO:0005216 | ion channel activity | 0.003 |

# Conclusions

- Method for predicting function using Rosetta protein predictions

- Accurately match protein predictions to known structures

- Accurately prediction functions

- Integrate multiple pieces of evidence to increase coverage

- Predict specific and unique functions

# Acknowledgements

### Bonneau Lab

Glenn Butterfoss
Thadeous Kacmarczyk
Peter Waltman
Aviv Madar
Kevin Belasco
Alex Pine
Richard Bonneau

### University of Washington

Lars Malmstroem
David Baker
Trisha N. Davis
Michael Riffle
Yeast Resource Center

### NYU

Sasha Levy
Peter McKenney
Jane Carlton
Dennis Shasha

Kris Gunsalus

Fabio Piano

Patrick Eichenberger
Biology Department
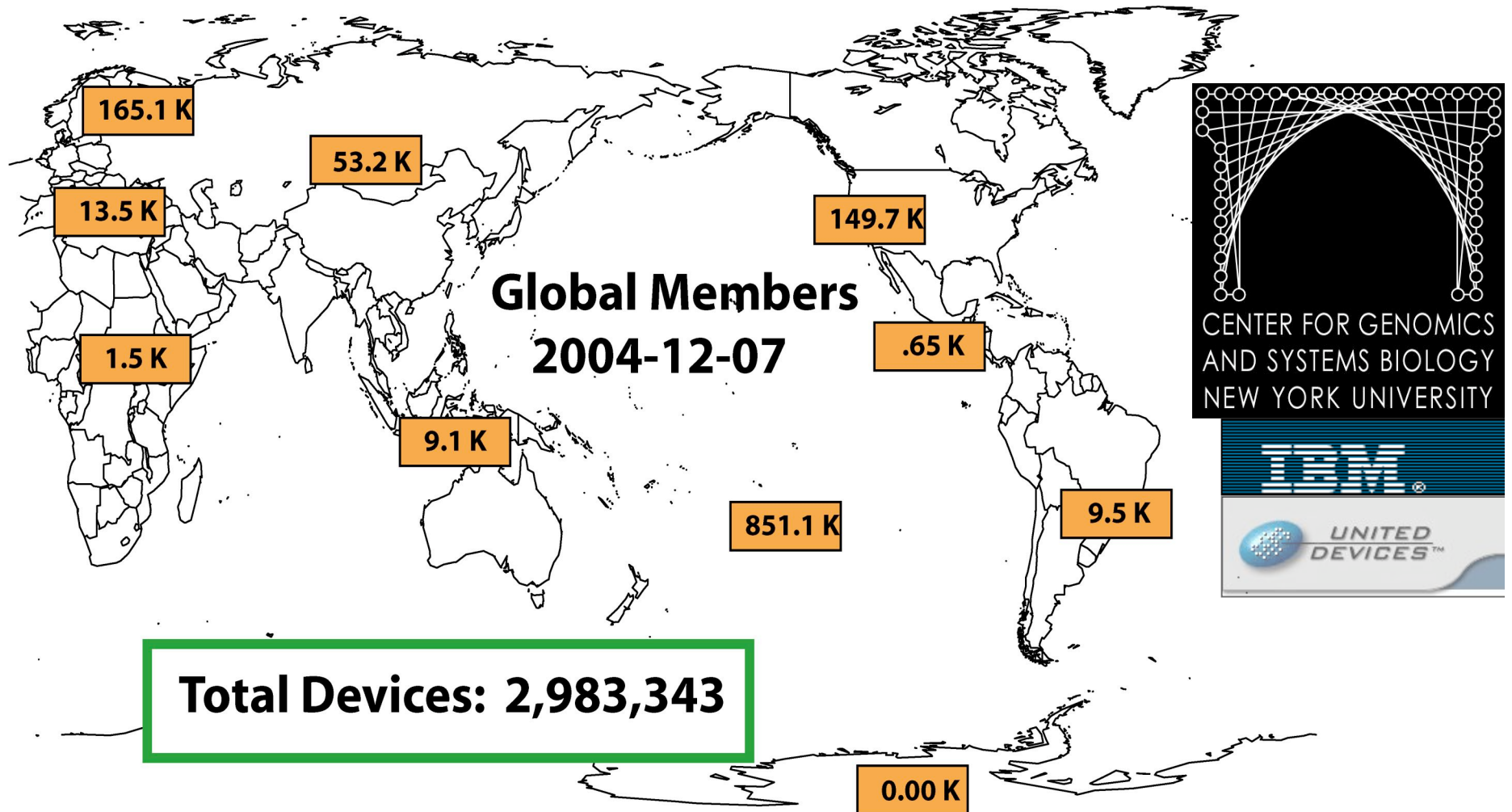
### IBM

Viktors Berstis
Keith J Uplinger
Bill Boverman

### Funding

DOD
DOE
NSF

### Rosetta-Commons