# Combining simulation and machine learning to recognize function in 4D
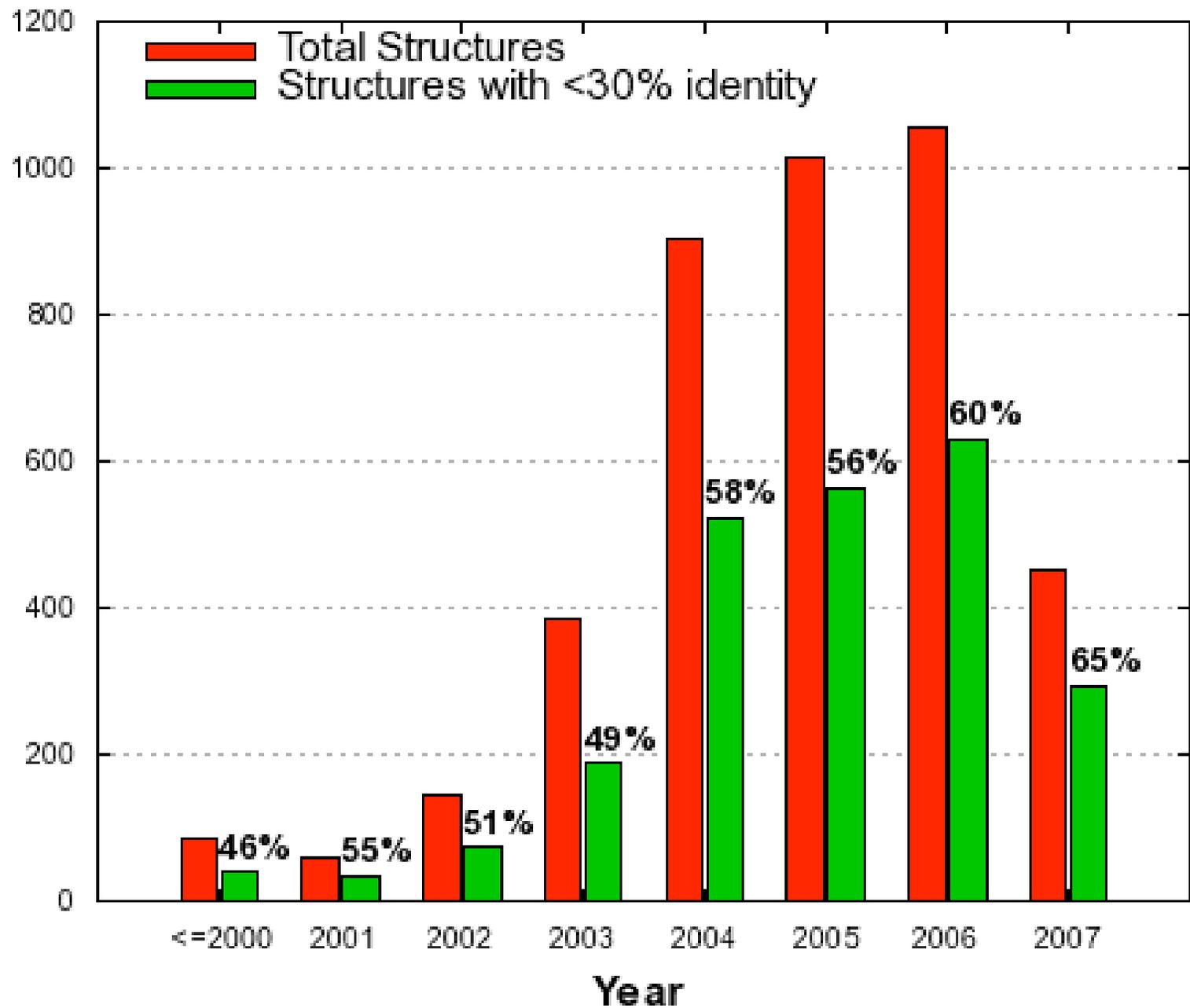
**Russ B. Altman, MD, PhD**

**Departments of Bioengineering, Genetics, Medicine, Computer Science Stanford University**

# Biological Motivations

- Structural genomics and prediction (structures without functions)
- Need to label putative molecular functions of new structures
- Need to discover new molecular functions
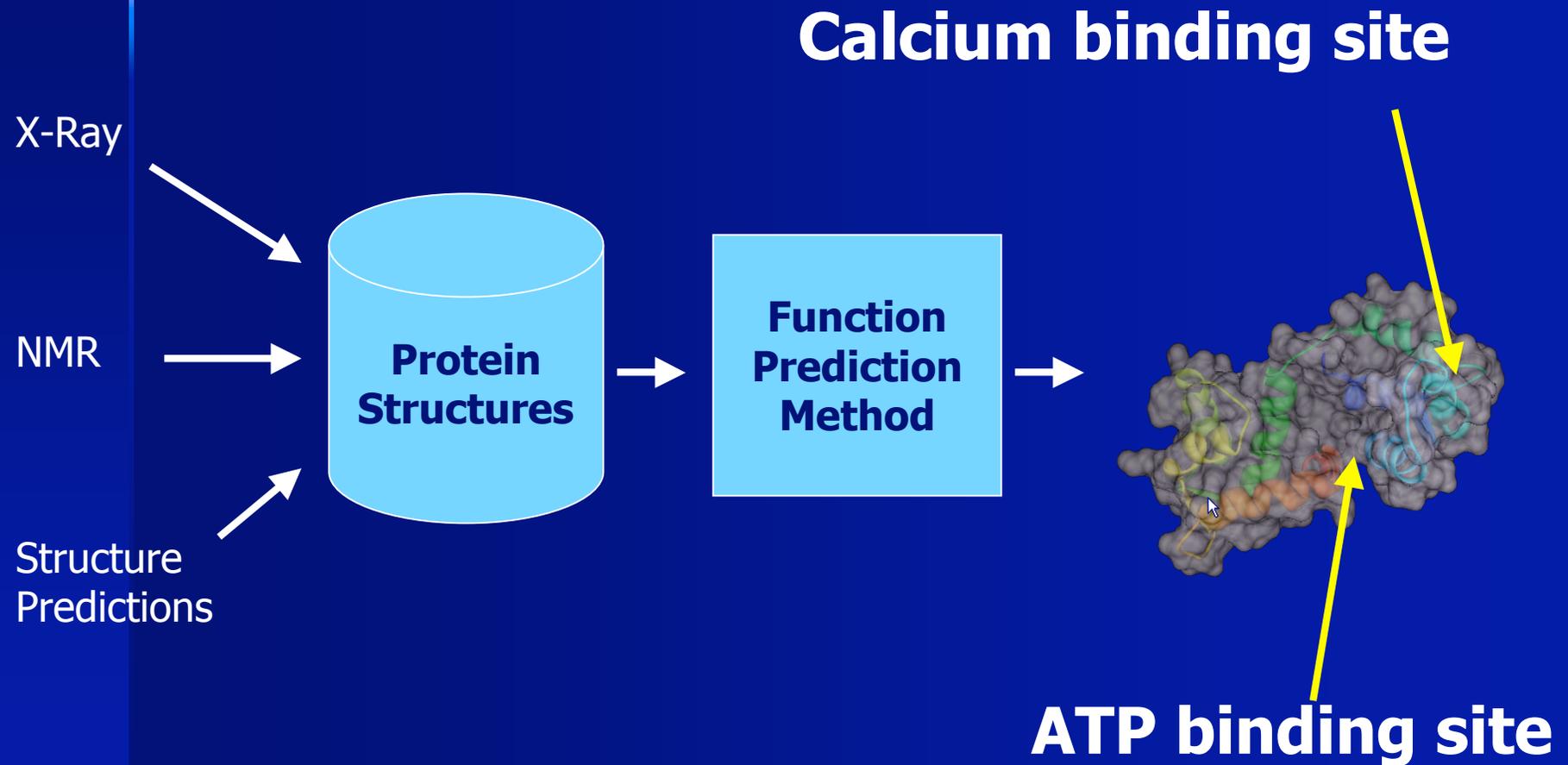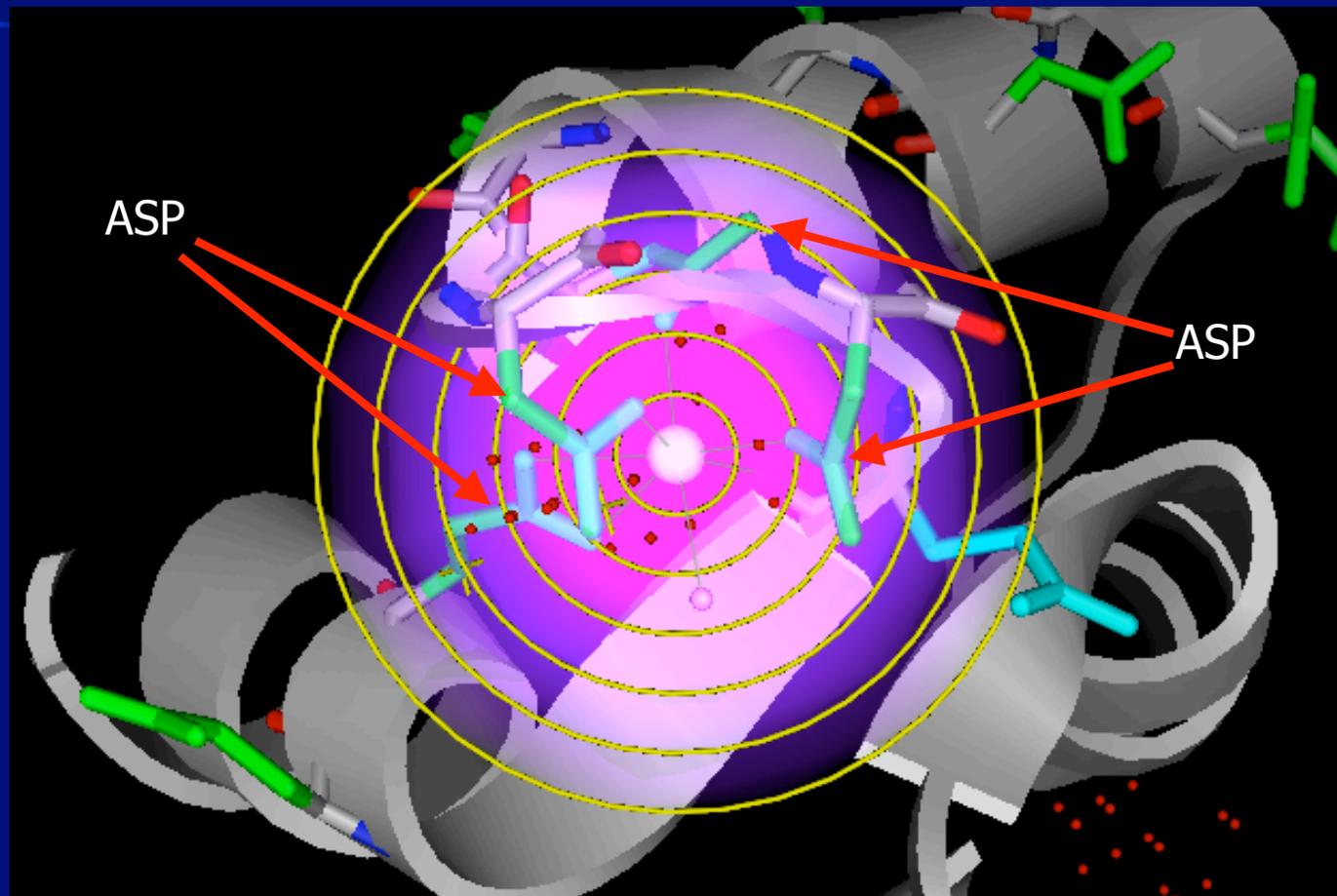- Increasing examples of polyfunctional proteins.

# Outline

- Introduction to FEATURE & WebFEATURE

- Using molecular simulation to improve FEATURE

# Protein Function Prediction & Annotation

**Calcium binding site**

X-Ray

NMR

Protein Structures

Function Prediction Method

Structure Predictions
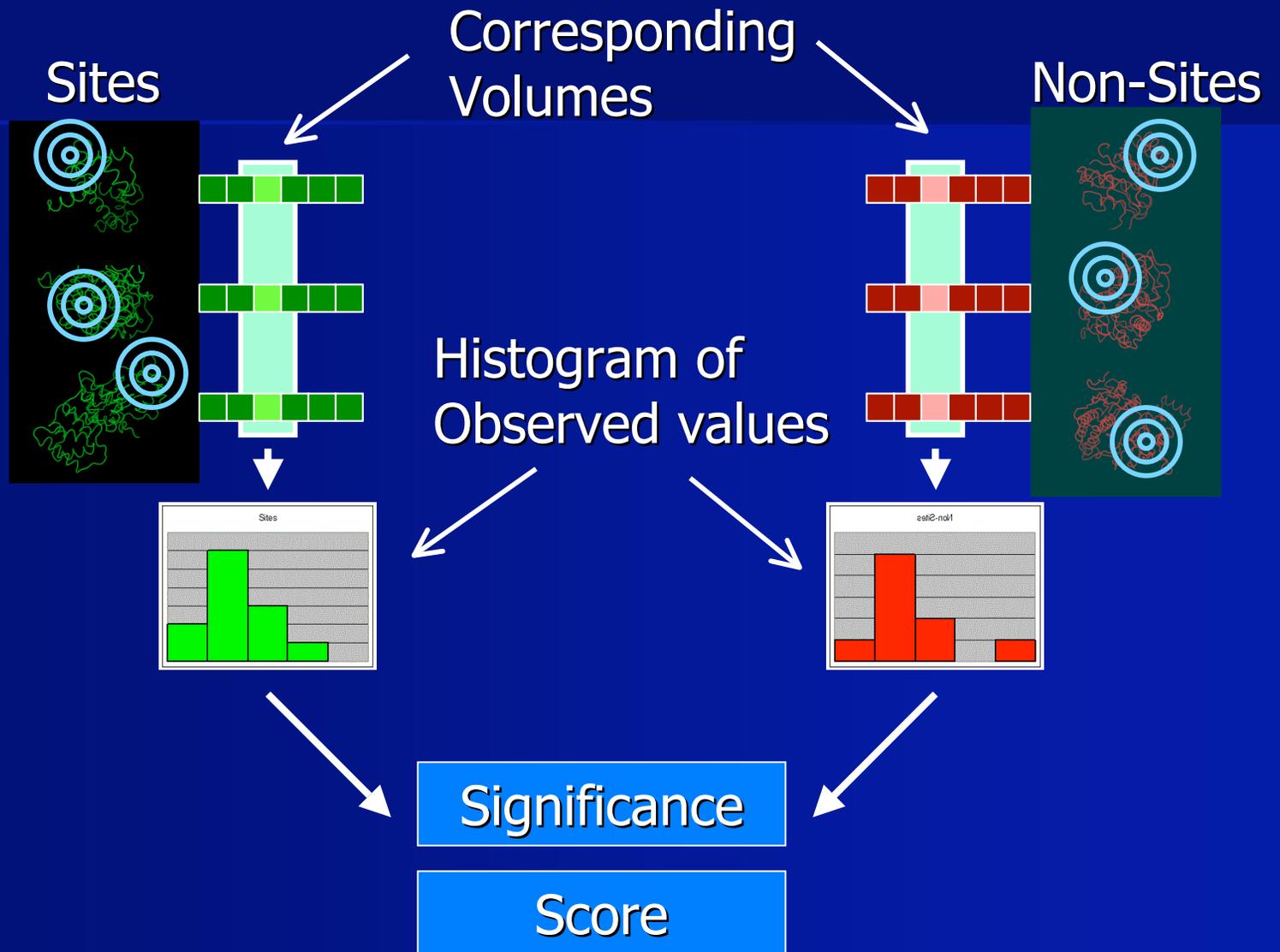
**ATP binding site**

# Radial Microenvironment

## Table 2.1 List of SeqFEATURE physicochemical properties.

*SeqFEATURE uses properties from the atom level up to the secondary structure level to characterize the physicochemical environment around a functional site.*

| AtomName | C | N | O | S |
|---|---|---|---|---|
| | ANY | OTHER | | |
| ChemicalGroup | Hydroxyl | Amide | Amine | Carbonyl |
| | RingSystem | Peptide | | |
| AtomProperties | VDWVolume | Charge | NegCharge | PosCharge |
| | ChargeWithHis | Hydrophobicity | Mobility | SolventAccessibility |
| ResidueName | ALA | ARG | ASN | ASP |
| | CYS | GLN | GLU | GLY |
| | HIS | ILE | LEU | LYS |
| | MET | PHE | PRO | SER |
| | THR | TRP | TYR | VAL |
| | HOH | OTHER | | |
| ResidueProperties | Hydrophobic | Charged | Polar | NonPolar |
| | Basic | Acidic | | |
| SecondaryStructure | 3Helix | 4Helix | 5Helix | Bridge |
| | Strand | Turn | Bend | Coil |
| | Het | Unknown | | |

# FEATURE Training

Sites

Corresponding Volumes

Non-Sites

Histogram of Observed values

Significance

Score
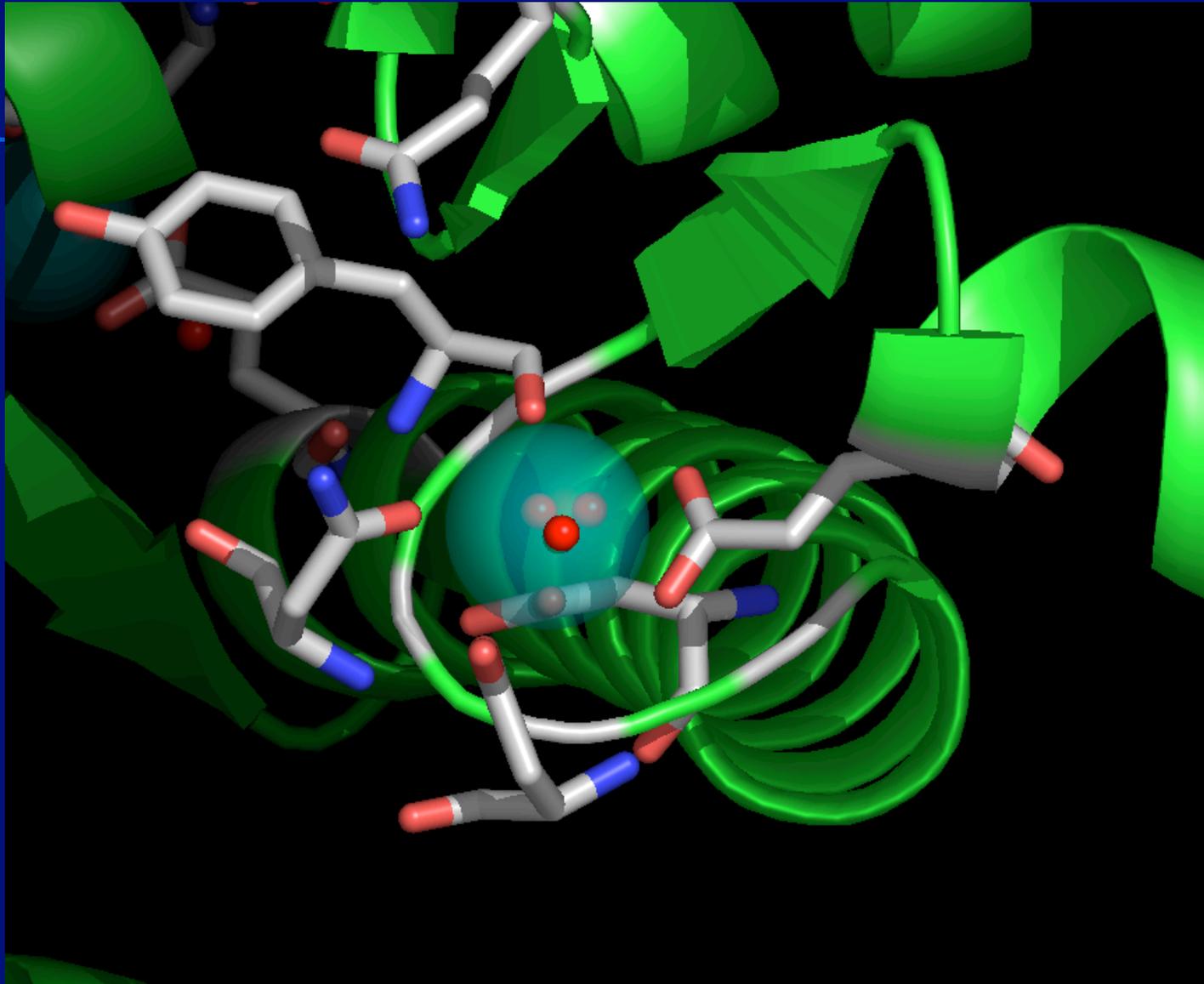
# FEATURE vectors

- 44 vectors of physical/chemical features
- Counted in 6 shells
- 44 x 6 = 264 feature-shell combinations that summarize the abundance of a feature in a shell (e.g. # of oxygens in shell that is 2-3 Angstroms from center)
- Vector = 264 continuous or discrete values describing environment around site center.
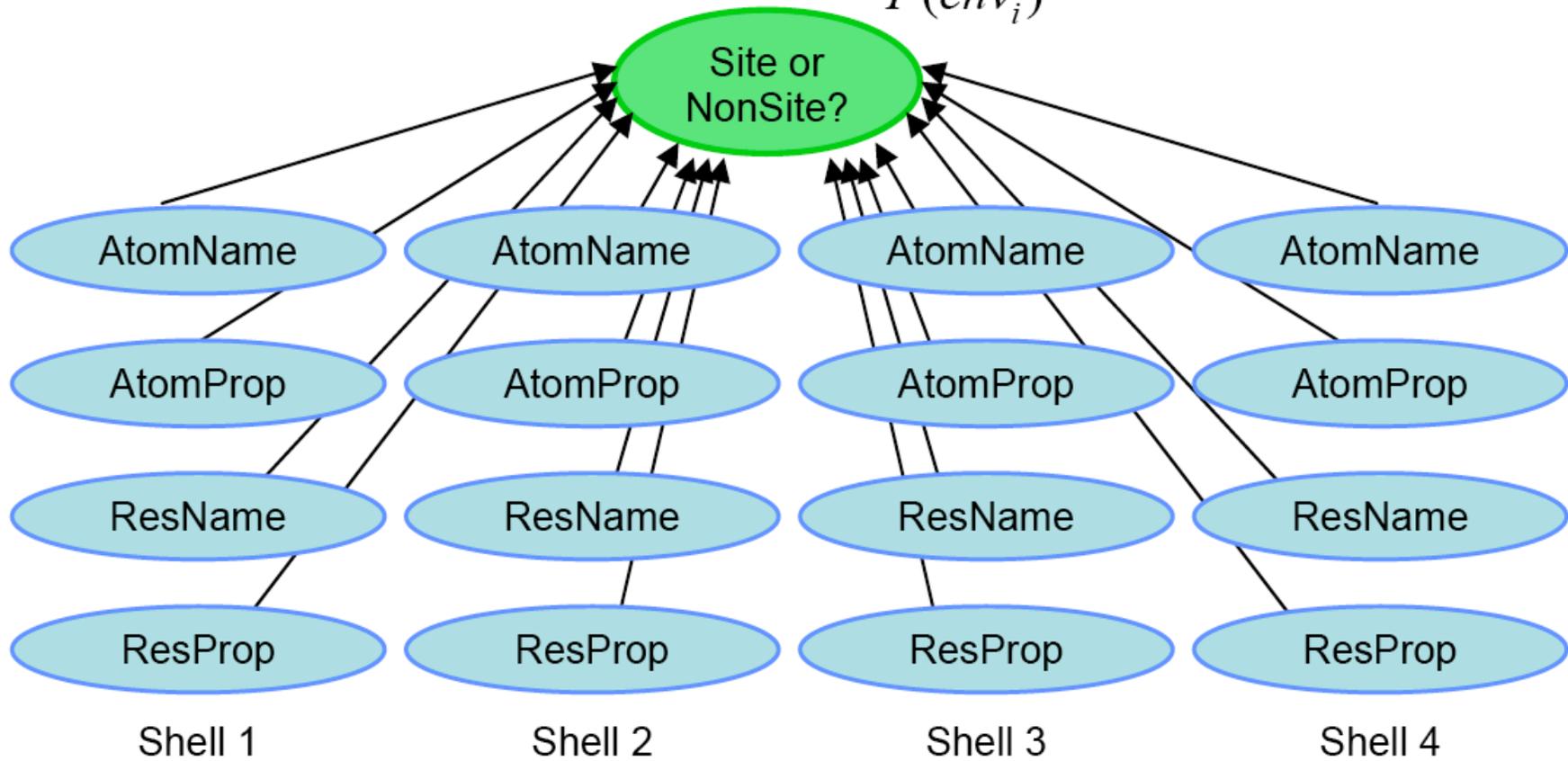
# Example: Calcium binding site

$$P(Site \mid env) = \prod \frac{P(env_i \mid site)}{P(env_i)} P(site)$$

Site or NonSite?

| AtomName | AtomName | AtomName | AtomName |
| AtomProp | AtomProp | AtomProp | AtomProp |
| ResName | ResName | ResName | ResName |
| ResProp | ResProp | ResProp | ResProp |
| Shell 1 | Shell 2 | Shell 3 | Shell 4 |

# WebFEATURE

## Automated function prediction in protein structures

- **WebFEATURE**
- FEATURE
- About
- Methods
- Usage
- Data
- Publications
- Projects
- Metals
- Helix Group

### Scan a structure for function
*(For assistance, please see our detailed instructions)*

**Step 1: Choose a structure**

Structure: PDB ID [____] or Upload (PDB format) [_____] [Browse...]

**Step 2: Choose a type of site to scan for**

Model: [ATP binding site ▼] [info]

**Step 3: Choose run mode**

Run mode: ⦿ Interactive ○ Email: [_____]

**Step 4: Submit and view results**

[Scan it!] [Clear Form]

WebFEATURE uses Jmol which requires Java 1.5.0 for visualization of results. These pages require Java and Javascript for proper interactivity. Visualization and interactivity can also be performed off-line using software packages such as RasMol, PyMol, and Chimera.
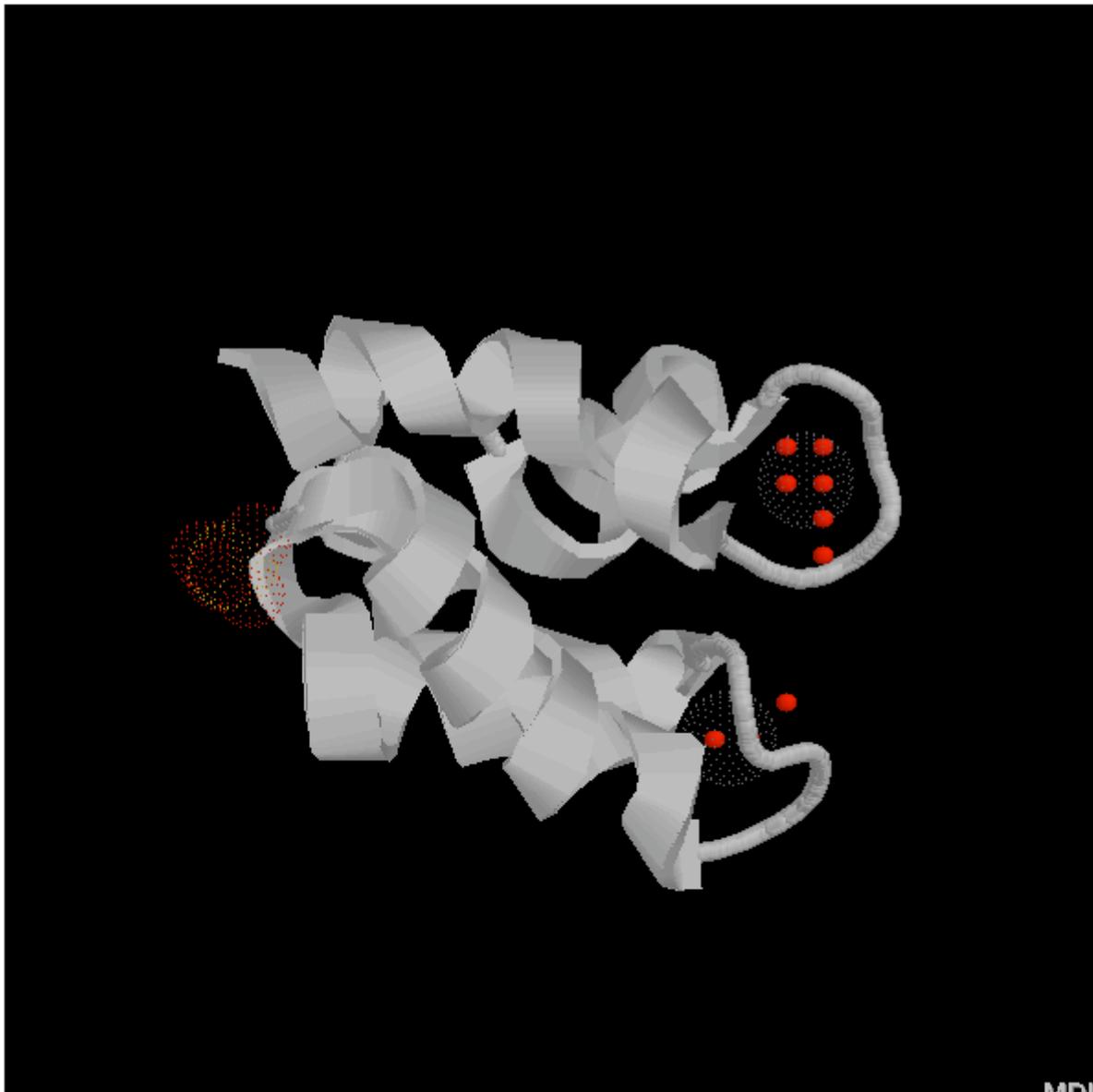
---

**EXAMPLES:**

| Model | PDBID | Description | Runtime |
|---|---|---|---|
| ○ Calcium binding site | 3icb | Intestinal Calcium Binding Protein | 5 secs |
| ○ ATP binding site | 1csn | Casein Kinase-1 | 45 secs |
| ○ Chloride binding site | 1pml | Tissue Plasminogen Activator Kringle 2 | 25 secs |
| ○ Diffuse bound Mg for RNA | 1ajf | P5B Stem-Loop from Group I Intron | 5 secs |
| ○ Site bound Mg for RNA | 1gid | P4-P6 RNA Ribozyme Domain from Group I Intron | 60 secs |
| ○ Trypsin model (SeqFeature) | 1bqy | Plasminogen Activator (Tsv-Pa) From Snake Venom | 5 secs |

---

Principal investigator: Russ Altman

**Shirley Wu**
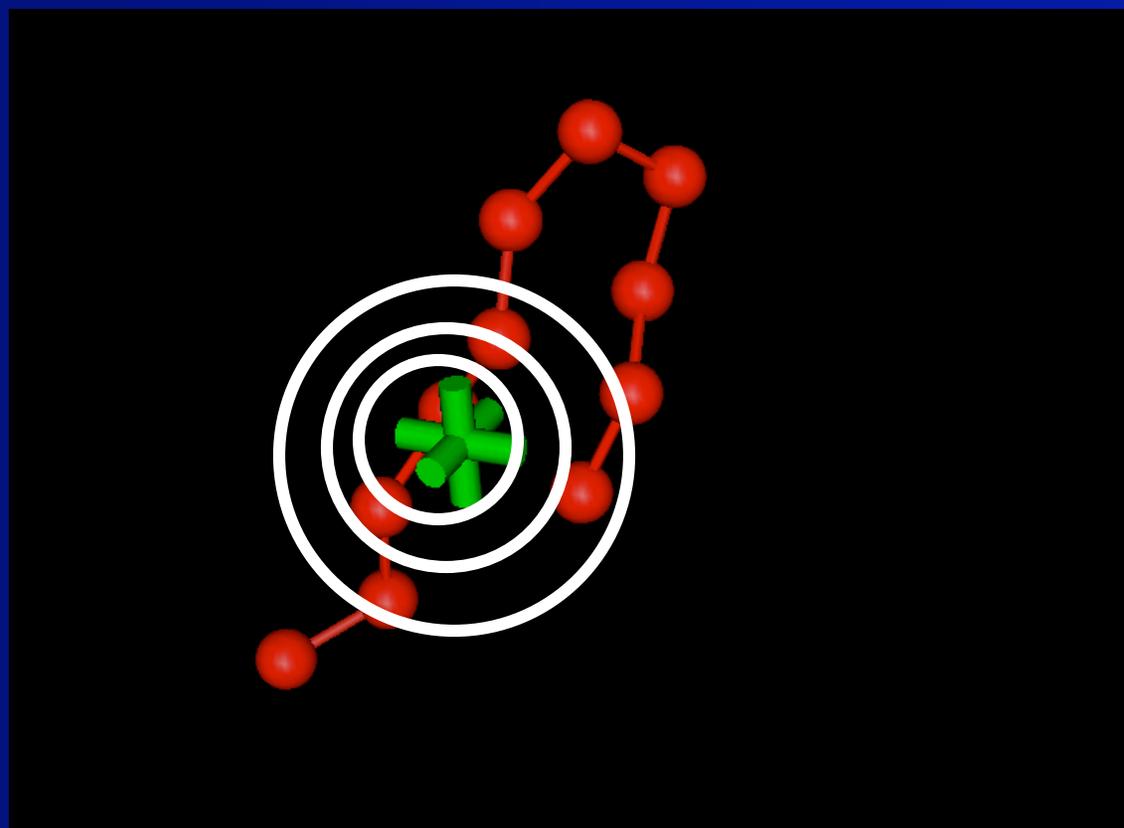
## Challenge:
Create a library of models available to FEATURE

## Solution:
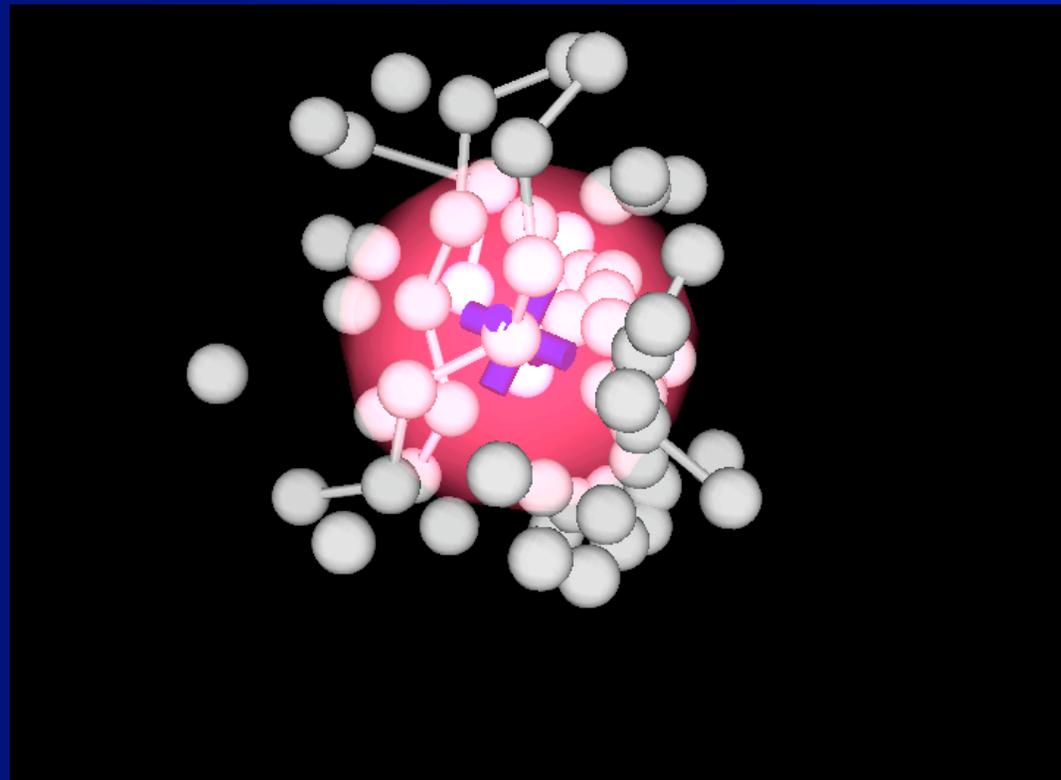Use 1D sequence motifs as "seeds" for 3D motifs

# SeqFEATURE

- Build from Sequence Motif Databases

- Automatically creates 3D motifs from 1D sequence motifs

- Hypothesis: 3D motifs perform better than 1D motifs in identifying functional sites
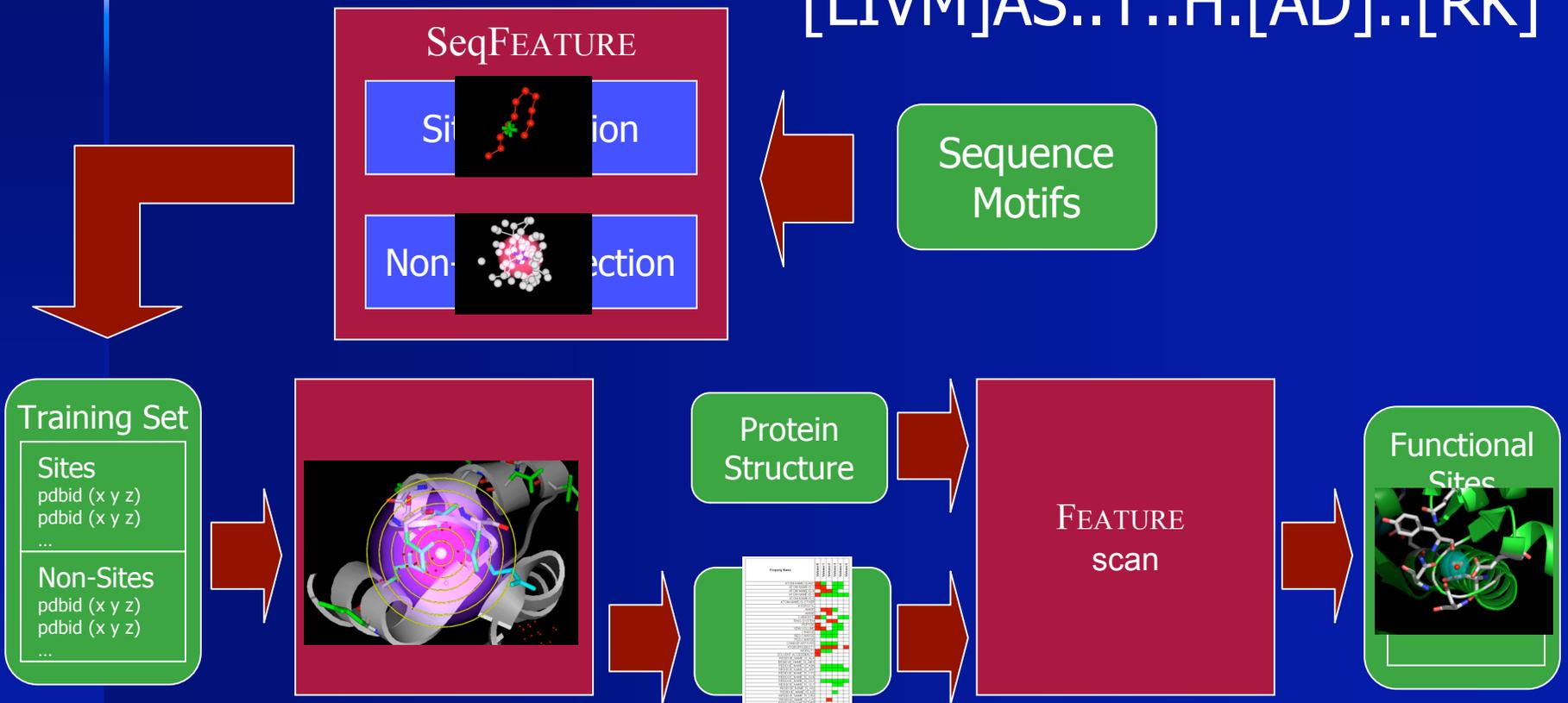
# Extracting 3D motif examples from 1D motif

# Extracting non-sites for background statistics

- Random amino acid with similar atom density and outside the site area

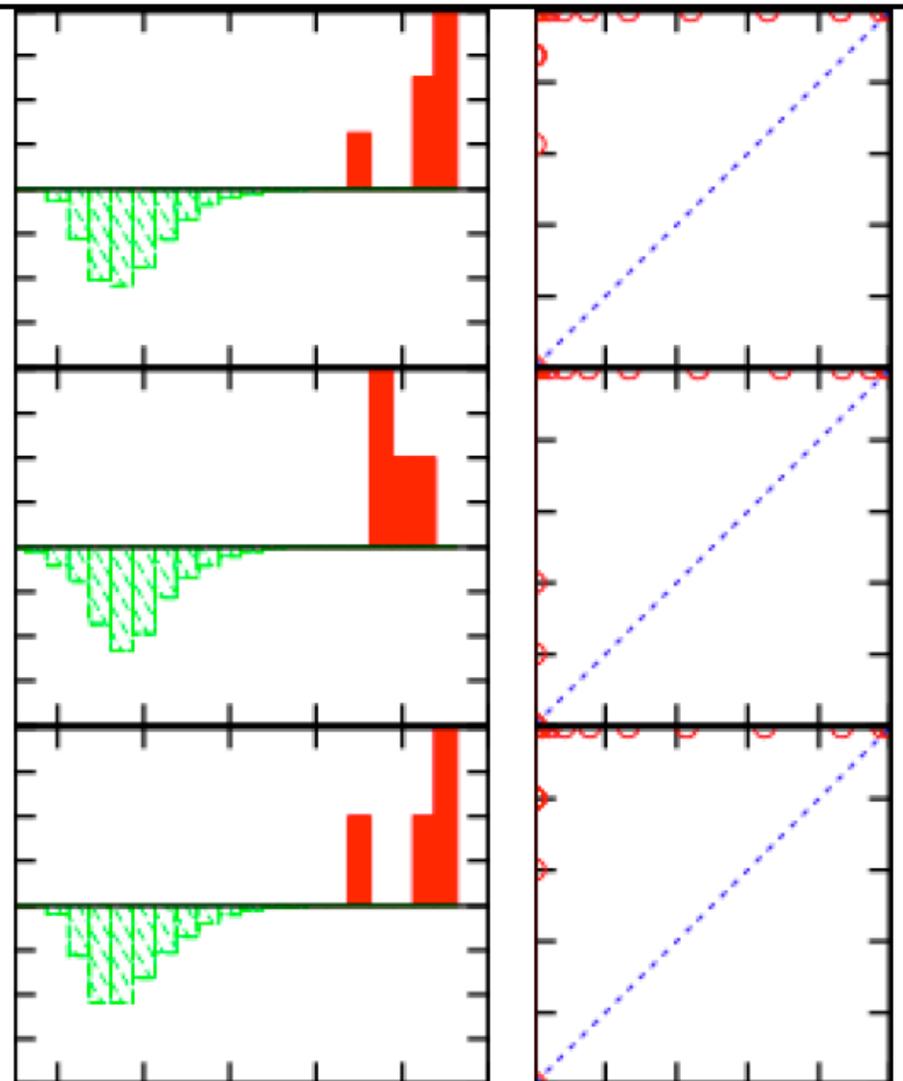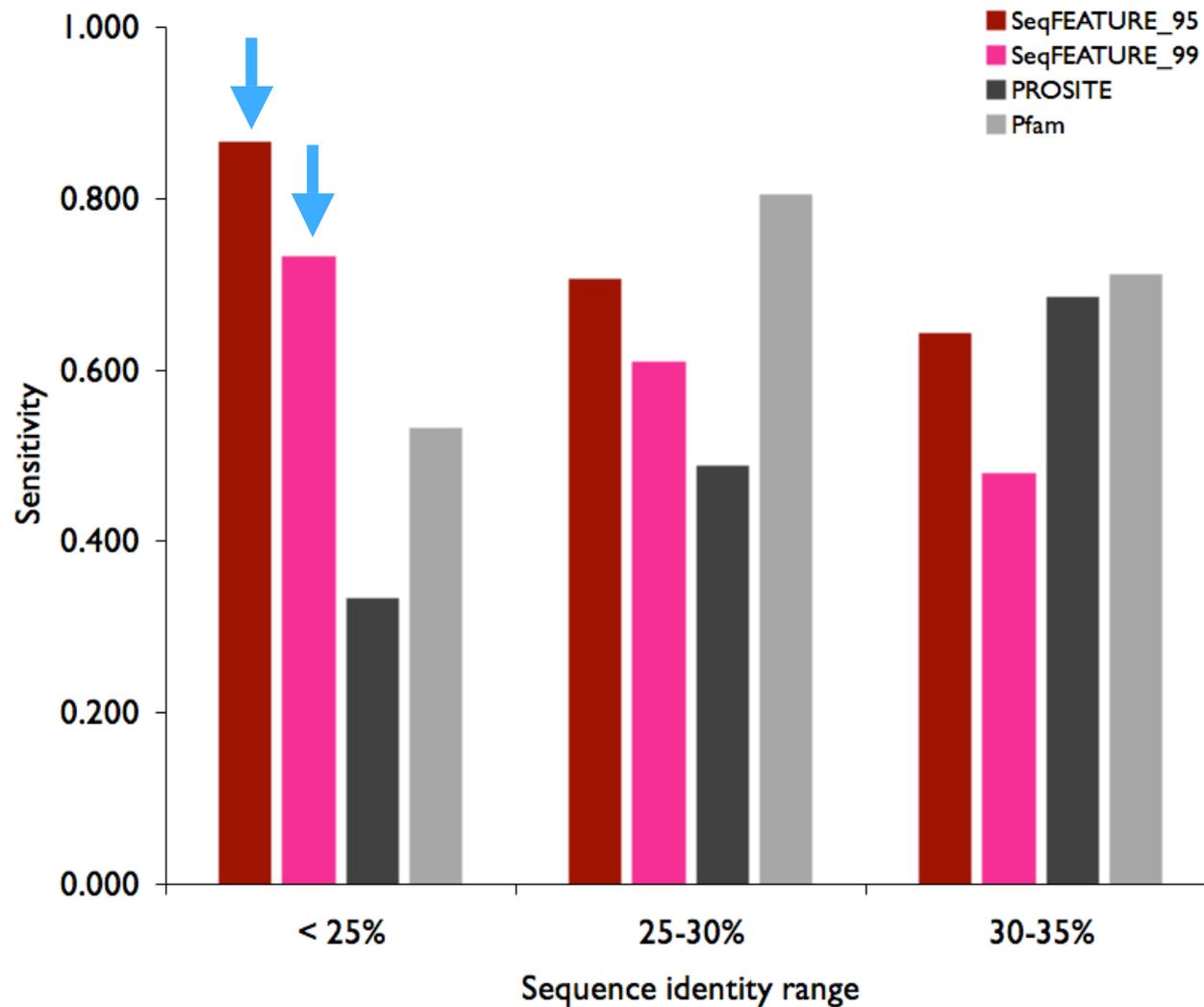| Model Performance | z-score Distribution | z-score ROC |
|---|---|---|
| **Model:** TYR_PHOSPHATASE_1.3.CYS.SG<br>**AUC:** 1.0<br>**zAUC:** 1.0<br>**psAUC:** None<br>**psPPV:** 0.9824<br>**psSens:** 0.9076 | | |
| **Model:** RNASE_T2_1.4.HIS.NE2<br>**AUC:** 1.0<br>**zAUC:** 1.0<br>**psAUC:** None<br>**psPPV:** 1.0<br>**psSens:** 1.0 | | |
| **Model:** HIPIP.7.CYS.SG<br>**AUC:** 0.9999<br>**zAUC:** 1.0<br>**psAUC:** None<br>**psPPV:** 1.0<br>**psSens:** 1.0 | | |

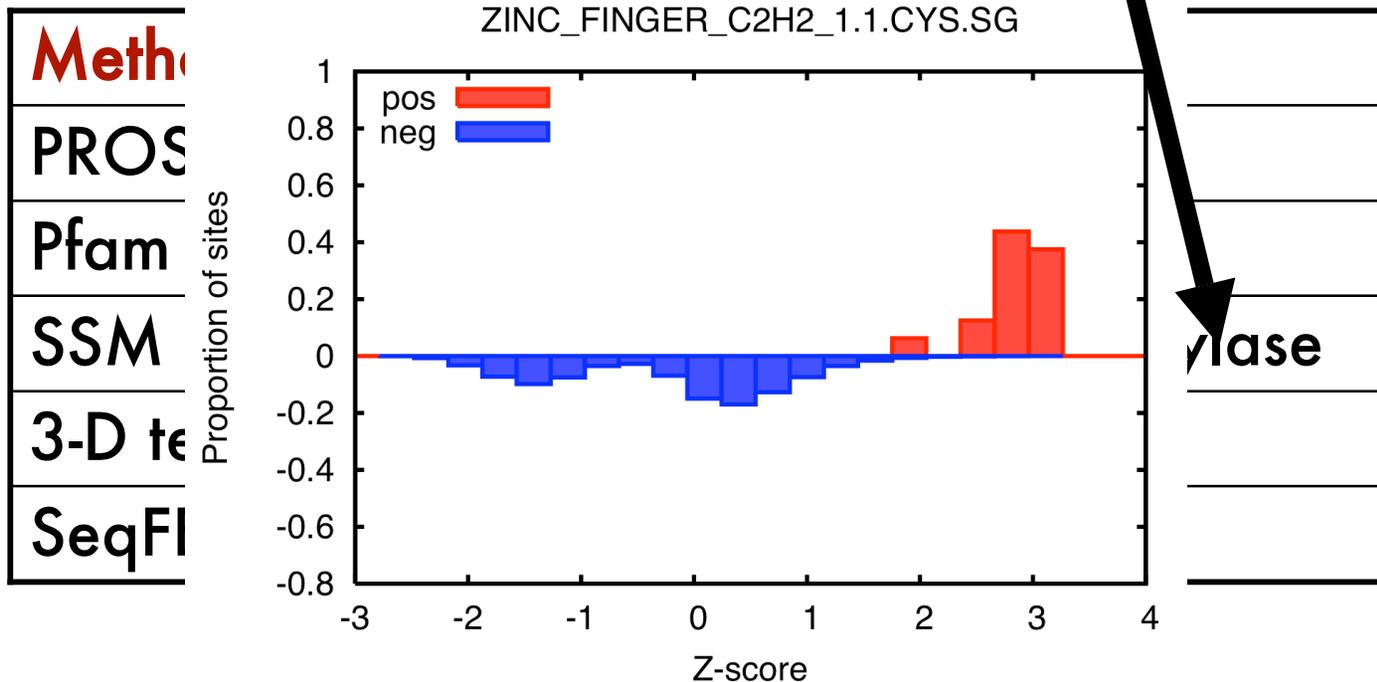# (seq)FEATURE performs better when sequence and structural similarity are low

# 1Z84: Galt-like protein

| SeqFEATURE Model | Site | z-score | Cutoff |
|---|---|---|---|
| ZINC_FINGER_C2H2_1.1.CYS.SG | CYS 63 | 4.713 | 3.177 |
| ZINC_FINGER_C2H2_1.3.CYS.SG | CYS 66 | 2.795 | 3.071 |

| Metho | | ylase |
|---|---|---|
| PROS | | |
| Pfam | | |
| SSM | | |
| 3-D te | | |
| SeqF | | |



ZINC_FINGER_C2H2_1.1.CYS.SG

# 1Z84: Galt-like protein

**2GLI**
**(zinc-finger protein)**

**1Z84**
**(prediction)**

# Beta-lactamase prediction (left)

(c)

(d)

# Go to other slides