

Predicting Metalloprotein Binding Sites: A Machine Learning Approach

Supervisor: Dr. Kaveh Kavousi

Presenter : Fereshteh Noroozi

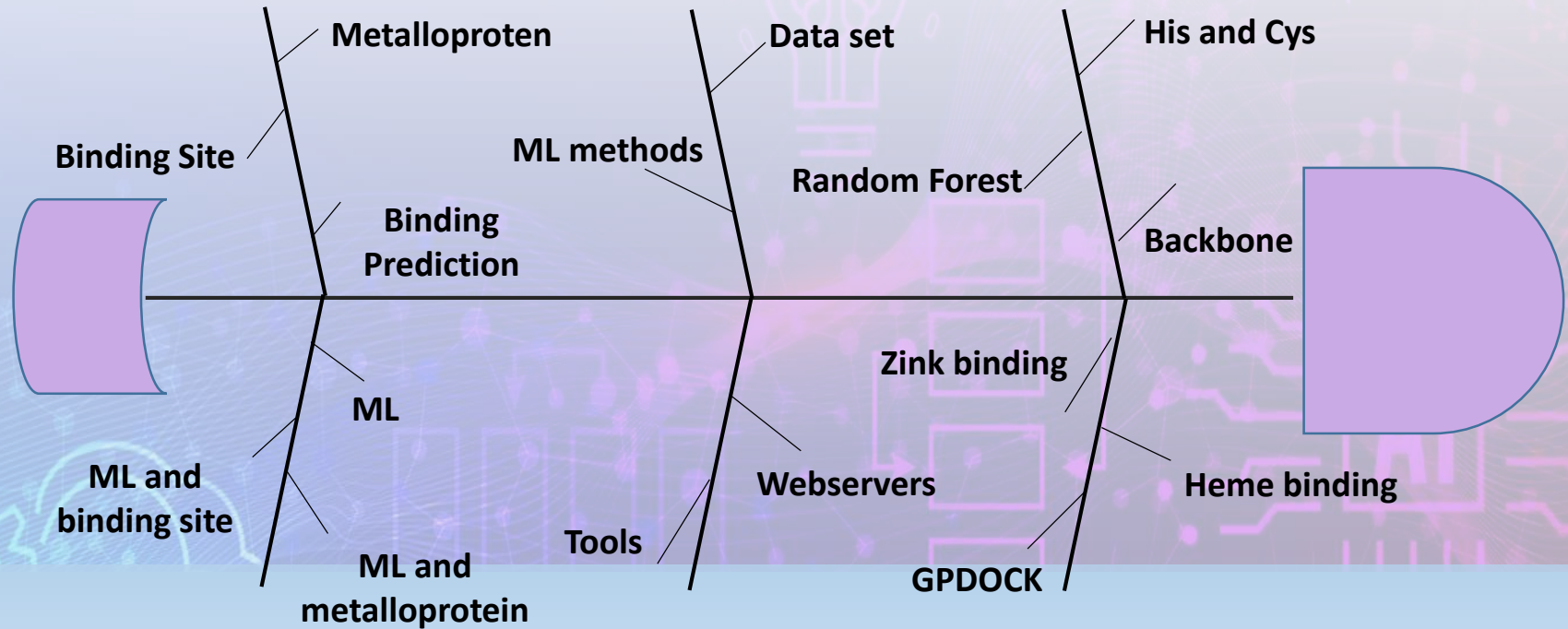
2024

Table of content

Introduction

Method

Results



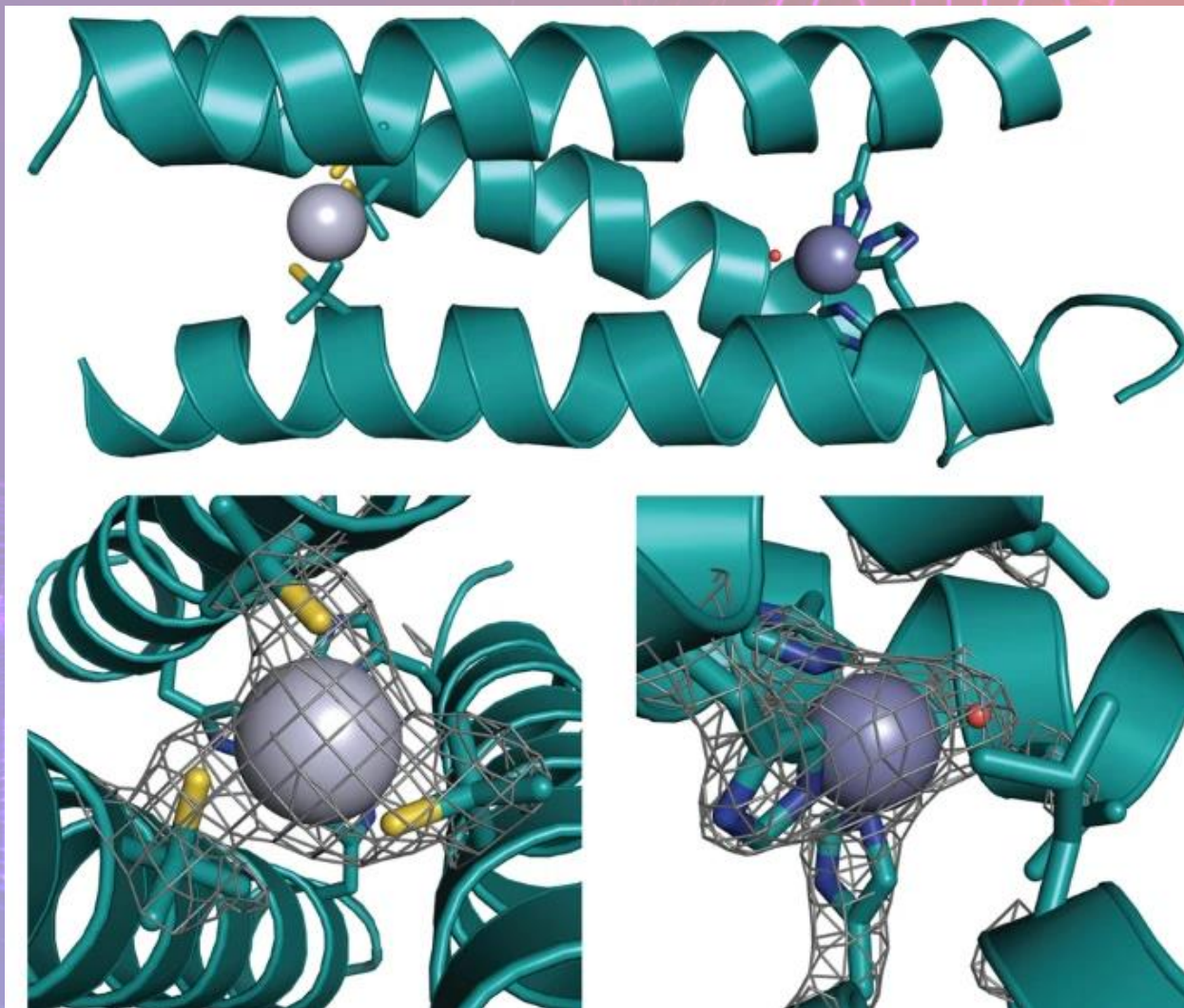
Introduction





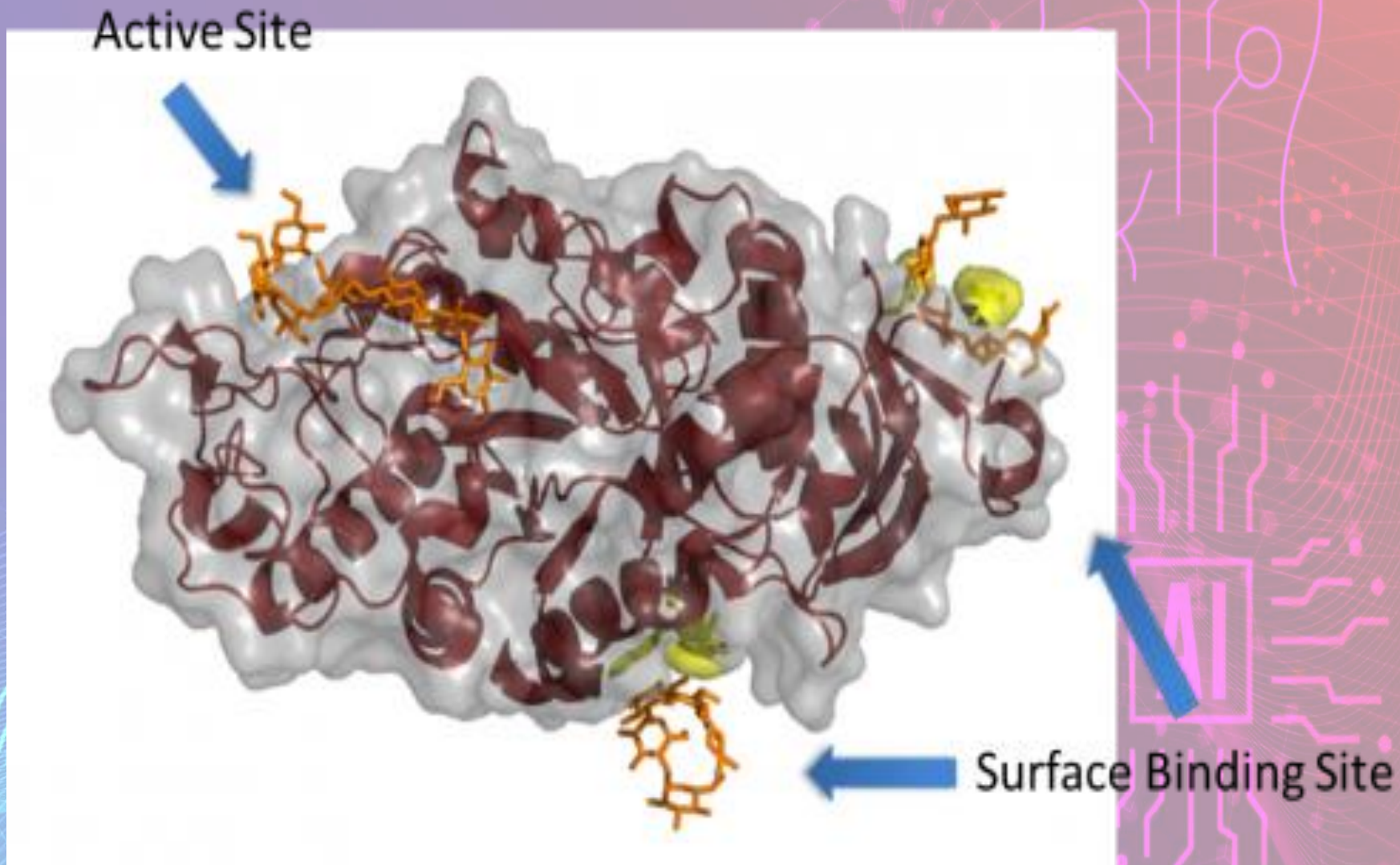
Metalloproteins

- Metalloproteins are a class of proteins that contain metal ions as essential components for their structure and function





Crucial Binding Site Significance



PDB ID 1rp8-alpha-amylase isozyme 1



Diverse Methods for Binding Sites Prediction

1

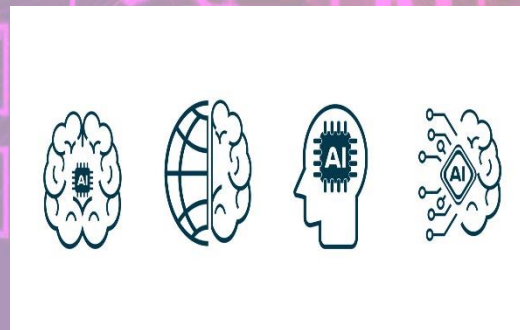
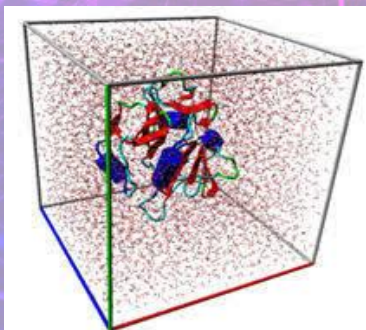
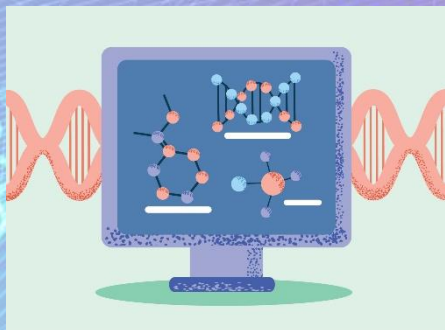
- **Sequence-Based Methods**

2

- **Structure-Based Methods**

3

- **AI Approaches**





Discovering with Machine Learning

1

- Autonomous extraction of influential patterns.

2

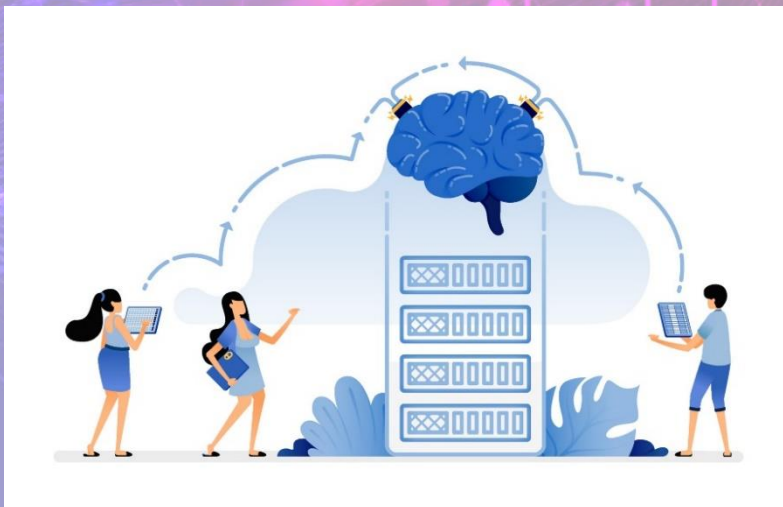
- Significance in Data Science for precise estimations.

3

- Empowering data scientists through data-driven predictions.

4

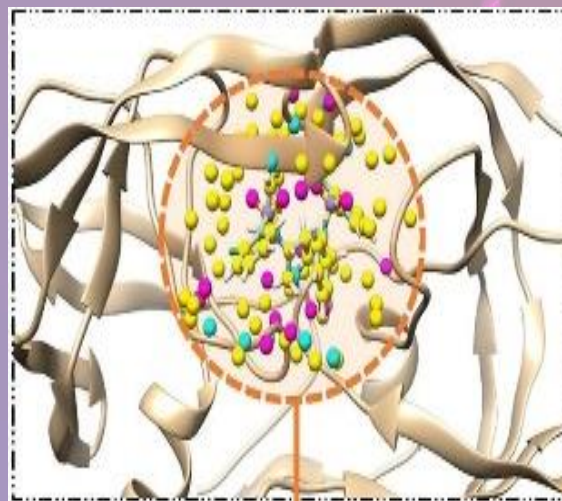
- Emerging Technologies





Utilizing Machine Learning for Prediction of Metalloprotein Binding Sites

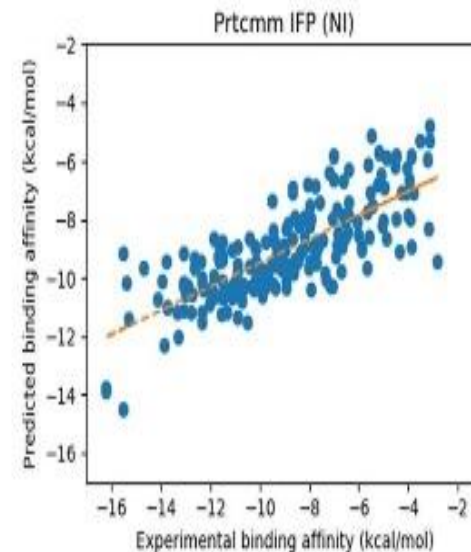
- Naive Bayesian
- Logistic Regression
- K-Nearest Neighbors
- Support Vector Machine
- Metalloprotein Prediction



IFP

0	0	1	1	0	1	0	0	1	1	0	1
---	---	---	---	---	---	---	---	---	---	---	---

Machine Learning





Purpose





Topic Statistics

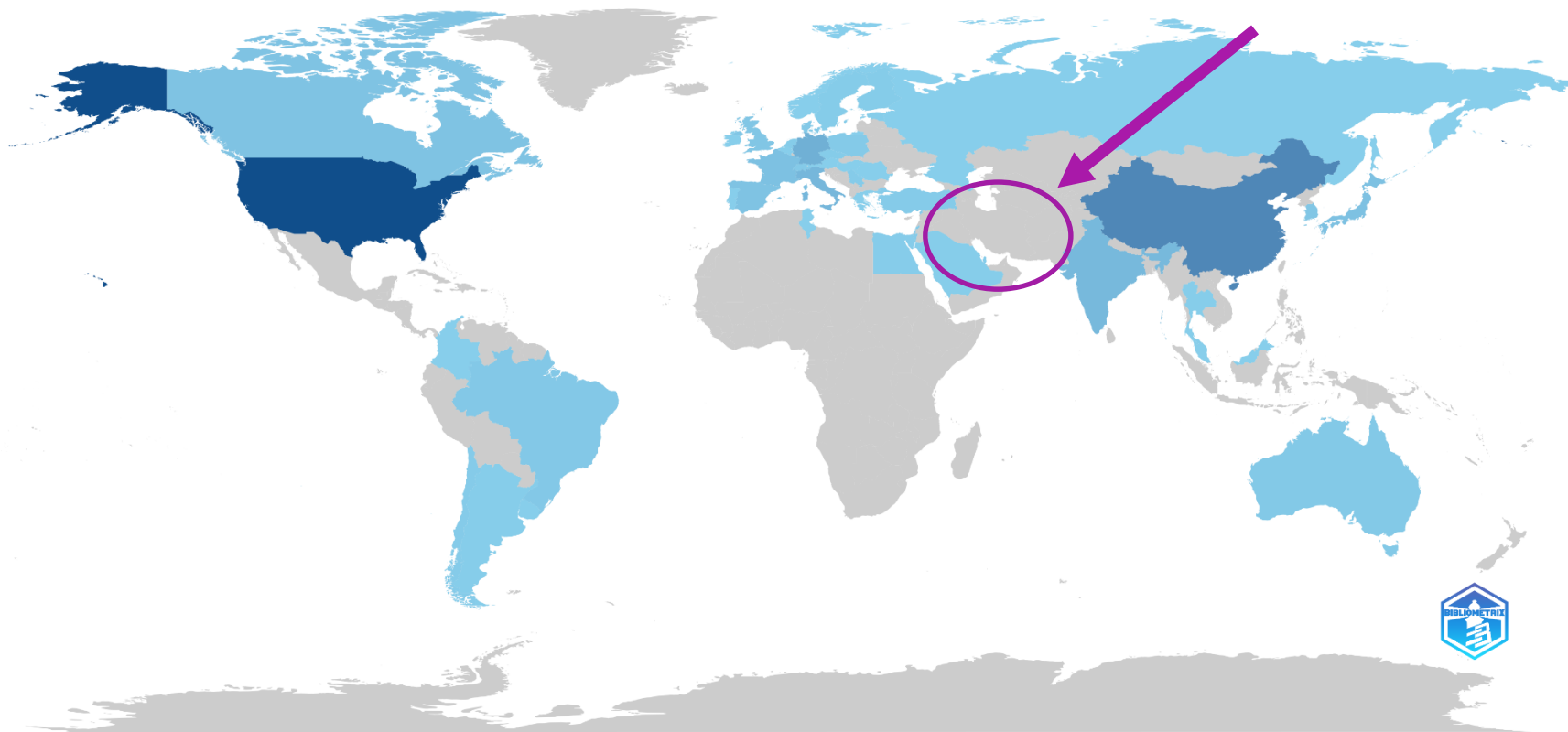
Annual Scientific Production





Topic Statistics

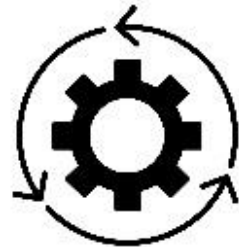
Country Scientific Production



Research background

Paper Title	Year	Model
Identifying Cysteines and Histidines in Transition-Metal-Binding Sites Using Support Vector Machines and Neural Networks	2006	SVM
Predicting zinc binding at the proteome level	2007	SVM
SCMHBP: prediction and analysis of heme binding proteins using propensity scores of dipeptides	2014	SVM
Prediction of Metal Ion Binding Sites in Proteins from Amino Acid Sequences by Using Simplified Amino Acid Alphabets and Random Forest Model	2017	Random Forest Model
Identifying metal binding amino acids based on backbone geometries as a tool for metalloprotein engineering	2021	Random Forest Model
GPDOCK: highly accurate docking strategy for metalloproteins based on geometric probability	2023	logistical regression model

Methods





Data Collection

Initial Data Collection

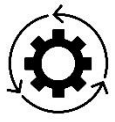
PDB
Uniprot
MetalPDB
PROSITE

Data Filtering

2.5 Å Res
HSSP value
cd-hit
UniRef 50
S 90%

Final Data

Clustering
Annotation
PyMOL
MBS
Balancing the Dataset
Selection of Ligands



Feature Selection

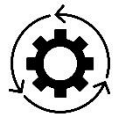
- Position-specific evolutionary profiles
- Global descriptor
- Conservation Features
- Conformational Similarity
- BLOSUM 50 Substitution Matrix
- Hydrophobicity
- Amino Acid Composition (AAC)
- Dipeptide Composition (DPC)

Sequence Length Relative to Average = $\frac{L - \bar{L}}{\bar{L}}$

Where:

- L represents the sequence length of the protein chain being considered.
- \bar{L} represents the average sequence length of all protein chains in the training set.

PSI-BLAST	Conservation of CYS	Conservation of HIS
4	00100	00010
3	01000	10000
0	00010	00001
2	10000	00100
1	00001	01000



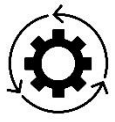
Feature Selection

(a). Atomic sets

	α : Heavy atoms of Ligand		β : Heavy atoms of Protein
	γ : O/N/S in Ligand		δ : O/N/S in Protein
	ϵ : Not O/N/S in Ligand		ζ : Not O/N/S in Protein
	η : Metal ions in Ligand		θ : Metal ion

(b). 36 features

$\sum_i \sum_j^{\alpha, \beta} (d_{ij} < 3.5)$	$\sum_i \sum_j^{\alpha, \beta} (d_{ij} < 4.5)$	$\sum_i \sum_j^{\alpha, \beta} (d_{ij} < 5.5)$	$\sum_i \sum_j^{\gamma, \delta} (d_{ij} < 3.5)$	$\sum_i \sum_j^{\gamma, \delta} (d_{ij} < 4.0)$	$\sum_i \sum_j^{\gamma, \delta} (d_{ij} < 4.5)$
$\sum_i \sum_j^{\gamma, \delta} (d_{ij} < 5.0)$	$\sum_i \sum_j^{\gamma, \delta} (d_{ij} < 5.5)$	$\sum_i \sum_j^{\gamma, \zeta} (d_{ij} < 3.5)$	$\sum_i \sum_j^{\gamma, \zeta} (d_{ij} < 4.0)$	$\sum_i \sum_j^{\gamma, \zeta} (d_{ij} < 4.5)$	$\sum_i \sum_j^{\gamma, \zeta} (d_{ij} < 5.0)$
$\sum_i \sum_j^{\gamma, \zeta} (d_{ij} < 5.5)$	$\sum_i \sum_j^{\epsilon, \delta} (d_{ij} < 3.5)$	$\sum_i \sum_j^{\epsilon, \delta} (d_{ij} < 4.0)$	$\sum_i \sum_j^{\epsilon, \delta} (d_{ij} < 4.5)$	$\sum_i \sum_j^{\epsilon, \delta} (d_{ij} < 5.0)$	$\sum_i \sum_j^{\epsilon, \delta} (d_{ij} < 5.5)$
$\sum_i \sum_j^{\epsilon, \zeta} (d_{ij} < 3.5)$	$\sum_i \sum_j^{\epsilon, \zeta} (d_{ij} < 4.0)$	$\sum_i \sum_j^{\epsilon, \zeta} (d_{ij} < 4.5)$	$\sum_i \sum_j^{\epsilon, \zeta} (d_{ij} < 5.0)$	$\sum_i \sum_j^{\epsilon, \zeta} (d_{ij} < 5.5)$	$\sum_i \sum_j^{\gamma, \theta} (d_{ij} < 3.0)$
$\sum_i \sum_j^{\gamma, \theta} (d_{ij} < 3.5)$	$\sum_i \sum_j^{\gamma, \theta} (d_{ij} < 4.0)$	$\sum_i \sum_j^{\gamma, \theta} (d_{ij} < 4.5)$	$\sum_i \sum_j^{\epsilon, \theta} (d_{ij} < 4.0)$	$\sum_i \sum_j^{\epsilon, \theta} (d_{ij} < 4.5)$	$\sum_i \sum_j^{\epsilon, \theta} (d_{ij} < 5.0)$
$\sum_i \sum_j^{\epsilon, \theta} (d_{ij} < 5.5)$	$\sum_i \sum_j^{\eta, \delta} (d_{ij} < 5.5)$	$\sum_i \sum_j^{\gamma, \delta} (min_{ij})$	$\sum_i \sum_j^{\gamma, \zeta} (min_{ij})$	$\sum_i \sum_j^{\epsilon, \delta} (min_{ij})$	$\sum_i \sum_j^{\epsilon, \zeta} (min_{ij})$



Model Selection-Support Vector Machine(SVM)

- **Gaussian kernels**(radial basis function (RBF) kernel)

$$K(x_i, x_j) = \exp(-\gamma \|x_i - x_j\|^2)$$

Width of the Gaussian distribution Euclidean distance

- **C regularization parameter**
- **K-fold cross validation**
- **LIBSVM package**

Model	γ	C
Binary SVM	0.05	0.1
Multiclass SVM	0.05	5



Model Selection-Random forest consists of multiple decision trees

- **Weka package**
- **Nested cross-validation**(double cross-validation)
- **K-fold**
- **GridSearchCV**(scikit-learn in Python)

Weka Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Classifier

Choose RandomForest -P 100 -I 100 -num-slots 1 -K 0 -M 1.0 -V 0.001 -S 1

Test options

☐ Use training set

☐ Supplied test set Set...

☒ Cross-validation Folds 10

☐ Percentage split % 66

More options...

(Nom) class

Start Stop

Result list (right-click for options)

10:40:49 - trees.RandomForest

Classifier output

Correctly Classified Instances 326 92.8775 %

Incorrectly Classified Instances 25 7.1225 %

Kappa statistic 0.8428

Mean absolute error 0.1287

Root mean squared error 0.2255

Relative absolute error 27.951 %

Root relative squared error 47.0057 %

Total Number of Instances 351

=== Detailed Accuracy By Class ===

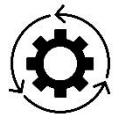
	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC
	0.865	0.036	0.932	0.865	0.897	0.844
	0.964	0.135	0.927	0.964	0.946	0.844
Weighted Avg.	0.929	0.099	0.929	0.929	0.928	0.844

=== Confusion Matrix ===

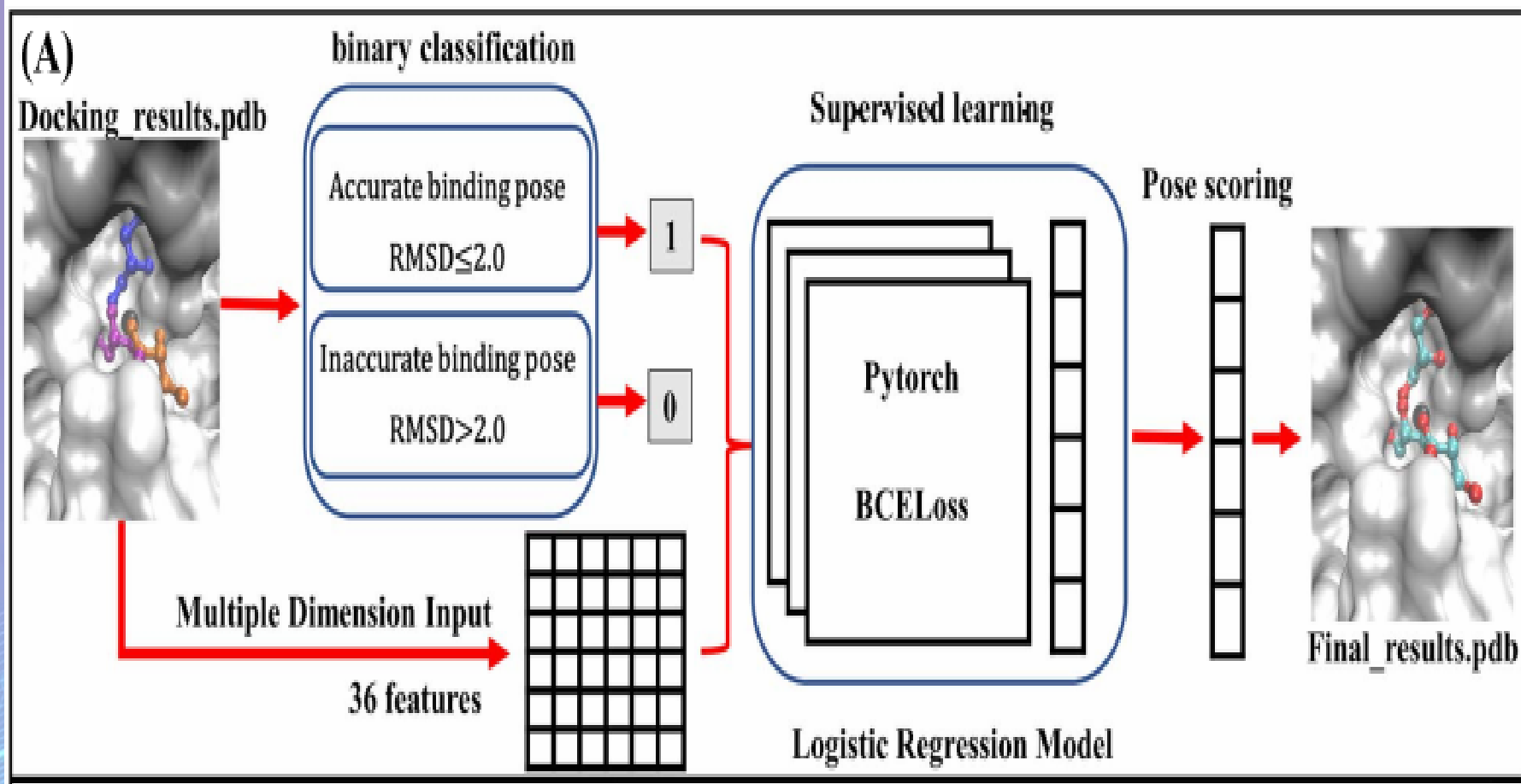
a b <-- classified as

109 17 | a = b

8 217 | b = g



Model Selection-Logistic regression model



$$P(y = 1|X) = \frac{1}{1+e^{-z}}$$

$$L(y, \hat{y}) = -\frac{1}{N} \sum_{i=1}^N [y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)]$$



Evaluation Procedure

Overall Accuracy



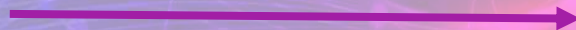
$$Accuracy = \frac{TP+TN}{TP+TN+FP+FN}$$

Precision (Positive Predictive Value)



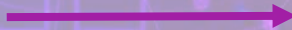
$$Precision = \frac{TP}{TP+FP}$$

Recall (Sensitivity)

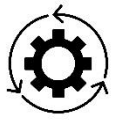


$$Recall (Sensitivity) = \frac{TP}{TP+FN}$$

Specificity (True Negative Rate)



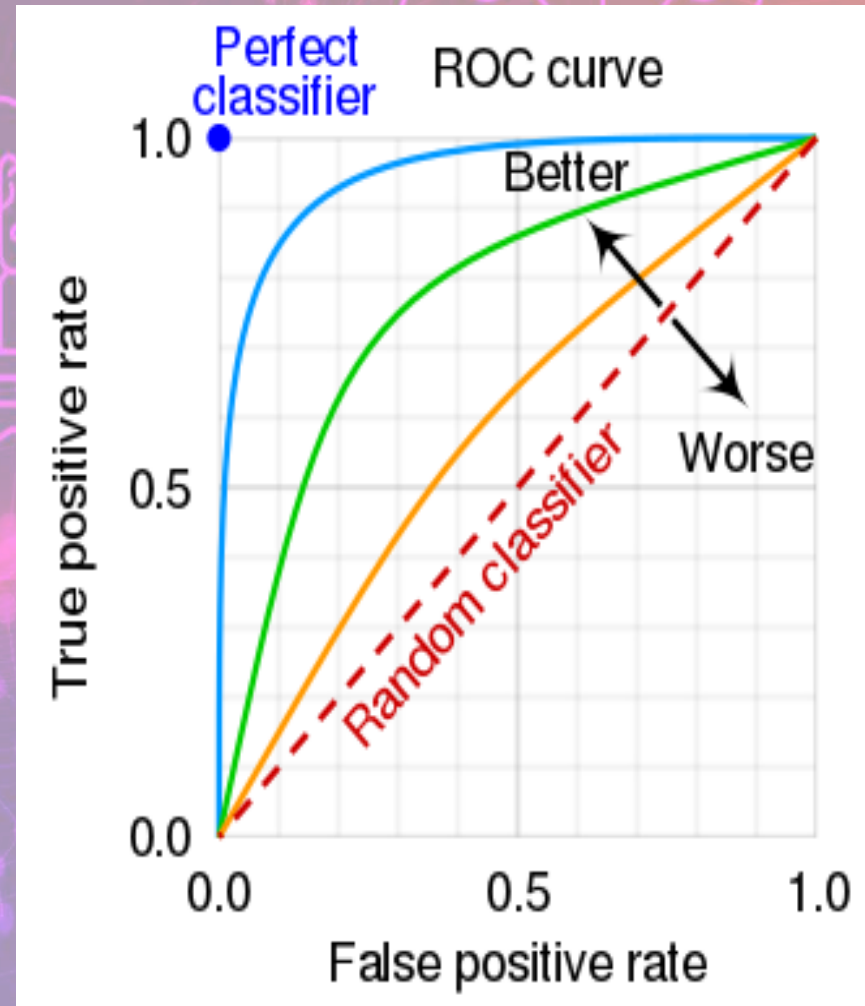
$$Specificity = \frac{TN}{TN+FP}$$

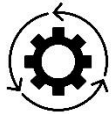


Evaluation Procedure

- Receiver Operating Characteristic (ROC) Curve
- Area Under the ROC Curve (AUC)
- Recall-Precision Curve (AURPC)
- Matthews Correlation Coefficient (MCC)

$$MCC = \frac{TP \times TN - FP \times FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}}$$





Webserver-MAHOMES Web Server

MAHOMES II

Metal Activity Heuristic of Metalloprotein and Enzymatic Sites (MAHOMES) II - Predicts if a protein bound metal ion is enzymatic or non-enzymatic

Overview

The ability to distinguish enzyme from non-enzyme sites remains an unsolved, challenging task. We've developed MAHOMES, a machine learning based tool which classifies metals bound to proteins as enzymatic or non-enzymatic. We intend to build on the previous work to make MAHOMES II, a more stable and robust version with a web server.

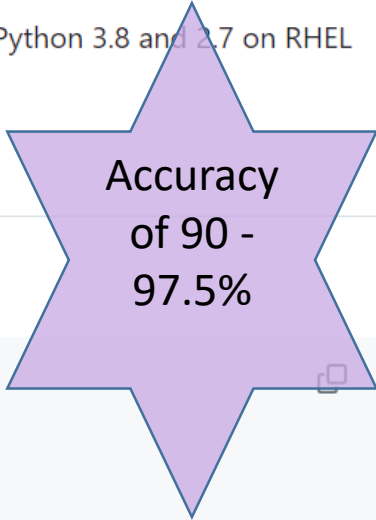
System requirements

Feature calculations also require using Rosetta, FindGeo, and blues which we run using Python 3.8 and 2.7 on RHEL 8 operating system.

Installation guide

set up virtual environment:

```
$ virtualenv --version # check for virtual enviroment
$ pip install virtualenv # download using pip if no version is found
$ virtualenv -p /usr/bin/python3 venv # create new virtual environment
$ source venv/bin/activate # switch to new env
$ pip install -r requirements.txt # add packages to environment
```



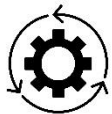
Accuracy
of 90 -
97.5%



Webserver-MAHOMES Web Server

Feature group
Local protein density
solvation
Pocket void
pKa
Rosetta
Electrostatics
Pocket hydrophobicity
BLUUES
SolvEnergy
Metal coordination geometry

Ponent	Description
Backend	Implemented in Python 3
Framework	Utilizes the Flask framework
Templates	HTML client-side interfaces created using Jinja
Job Metadata Storage	Stored in JSON files
Database	Metadata and scheduling information managed in an SQLite3 database
Job Execution Program	Monitors the SQLite database for new submissions, executes jobs, and sends result links via email
Web Server	Hosted on the Slusky Lab web server, operating as a virtual machine in the University of Kansas's enterprise data center



Webserver-ZincBinder Web Server



[Home](#) [Prediction](#) [Algorithm](#) [Developers](#) [Help](#) [Contact](#)

For detailed methodology and further information, see the [Help file](#).

Usage: Paste your sequence in the textarea provided or upload the file containing the sequence in [Fasta format](#) into the sequence field below and press the [Run Prediction](#) button.

Any other line numbers or whitespaces will be removed. [Tips](#)

Query title (optional)	1A8T <i>Query ID</i>
Input sequence format	Fasta Format <i>User can paste query sequences</i>
Query sequence	>1A8T:A PDBID:CHAIN:SEQUENCE AQKSVKISDDISITQLSDKVYTYVSLAEKGCWGMVPSNCGMIVNNHQAALLDTPINDAQTEMLVNWVTDLSLHA KVTTTIP NHWICDCICGLGYLQRKGVQSYANQMTIDLAKKGLPVPEHGFDTSLTVSLDGMPLQCYLLCGGHATDNV VWLPTENIL FGCCMLKDNQTTSIGNISDADVTAWPKTLDKVKAKFPSARYVVPCHGNYGGTELDITKQIVNQYESTSKP <i>Or Upload sequence file</i>
OR Upload Sequence file	Choose File <i>No file chosen</i> <i>Here user can select threshold</i>
SVM Threshold :	0.1 <i>1</i>
E-mail address for job completion alert (optional):	sri.abhishikha@gmail.com <i>Enter your mail Id</i>
<i>Run prediction</i> <input type="button" value="Run Prediction"/> <input type="button" value="Clear"/>	

Zincbinder Prediction Result

Query Search Detail

JOB-ID	zinc_7495
Number of Query Sequences	1
Predicted on	2:00:14 pm

Prediction Result

Protein-ID	zinc binding residue	position	znbinding score
1A8TAPDBIDCHAINSEQUENCE	H	82	1.5
1A8TAPDBIDCHAINSEQUENCE	H	84	2.0
1A8TAPDBIDCHAINSEQUENCE	D	86	0.7
1A8TAPDBIDCHAINSEQUENCE	C	87	0.8
1A8TAPDBIDCHAINSEQUENCE	H	145	1.5
1A8TAPDBIDCHAINSEQUENCE	C	164	1.3
1A8TAPDBIDCHAINSEQUENCE	H	206	2.3

85.37%
sensitivity
with 86.20%
specificity



Webserver-ZincBinder Web Server

Dataset Description	Count
Total zinc-bound protein structures	1996
Total zinc-binding site IDs	3896
Protein chains obtained from PDB structures	5169
Protein chains interacting with zinc (HETEROATOM)	3924
Resolution range of protein structures (PISCES parameter)	0-2.5Å



Tools

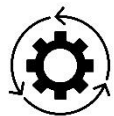
Metal Binding Site Prediction: 1D and 3D Approaches

Metal 1D Approach

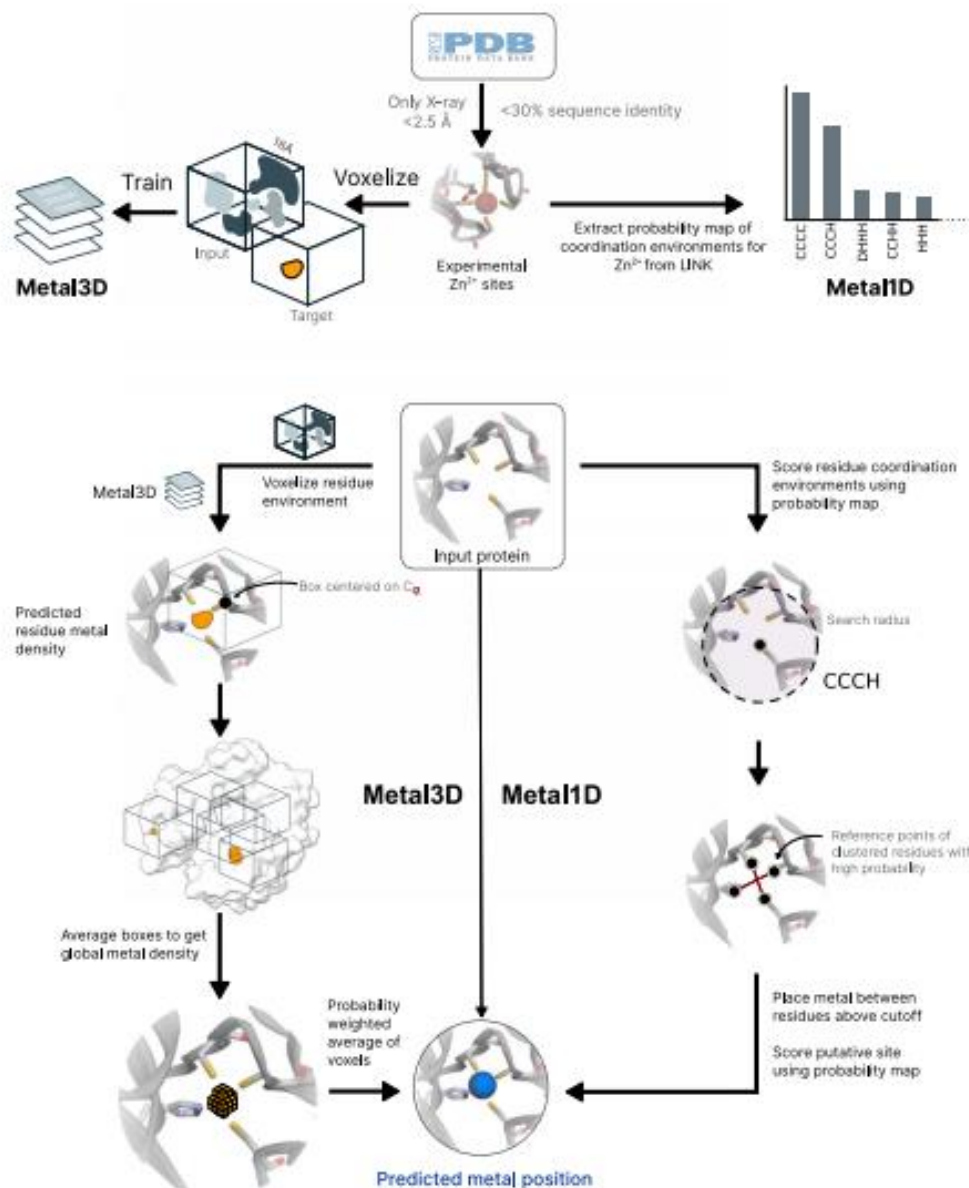
Probability Map
Generation(ProbMapGenerator.ipynb),
BioPandas python library

Metal 3D Approach

Voxelization Process(moleculekit
Python library),



Tools



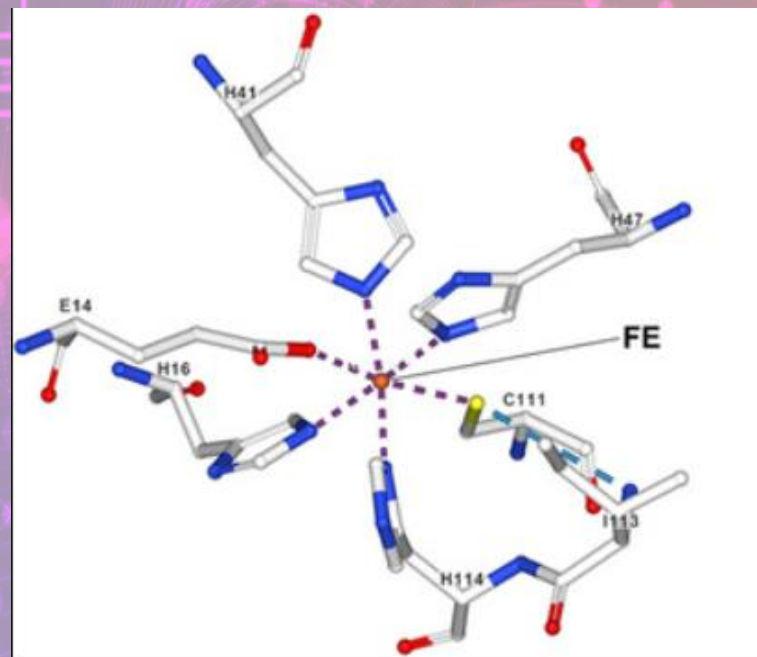
Results



Identifying Cys and His in Transition-Metal-Binding Sites Using-SVM

- Prediction of histidine in two states
- Prediction of cysteine in three states
- SVM trained to locally classify the binding state of single HIS and CYS
- Referring to metal-binding amino acids as "ligands."

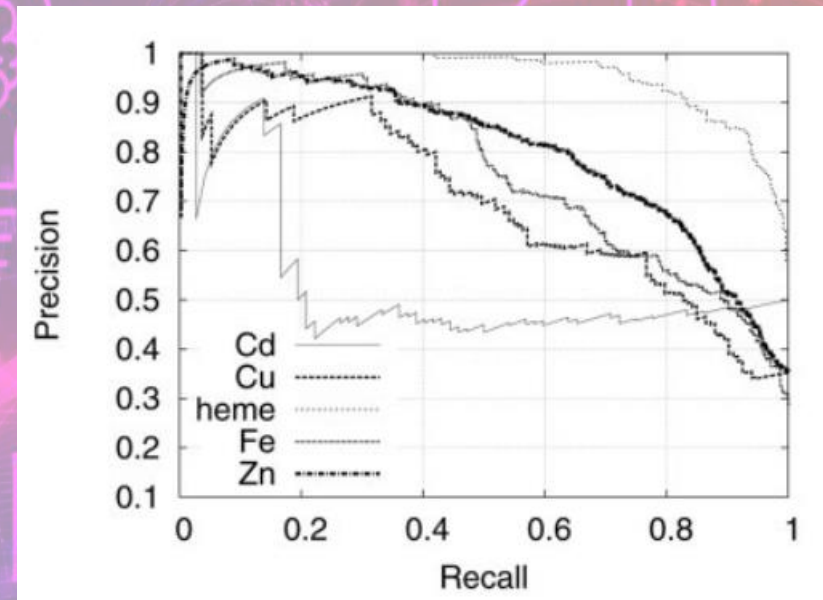
Metal	CYS	HIS
Zn	46 (508/1115)	24 (374/1562)
Heme	50 (115/230)	34 (151/450)
Fe/S	63 (205/326)	3 (10/329)
Cu	33 (36/108)	32 (86/269)
Cd	62 (48/77)	32 (25/79)
Fe	13 (16/122)	18 (59/325)
Ni	4 (2/46)	16 (18/112)
Any	48 (930/1923)	25 (723/2942)





Identifying Cys and His in Transition-Metal-Binding Sites Using-SVM

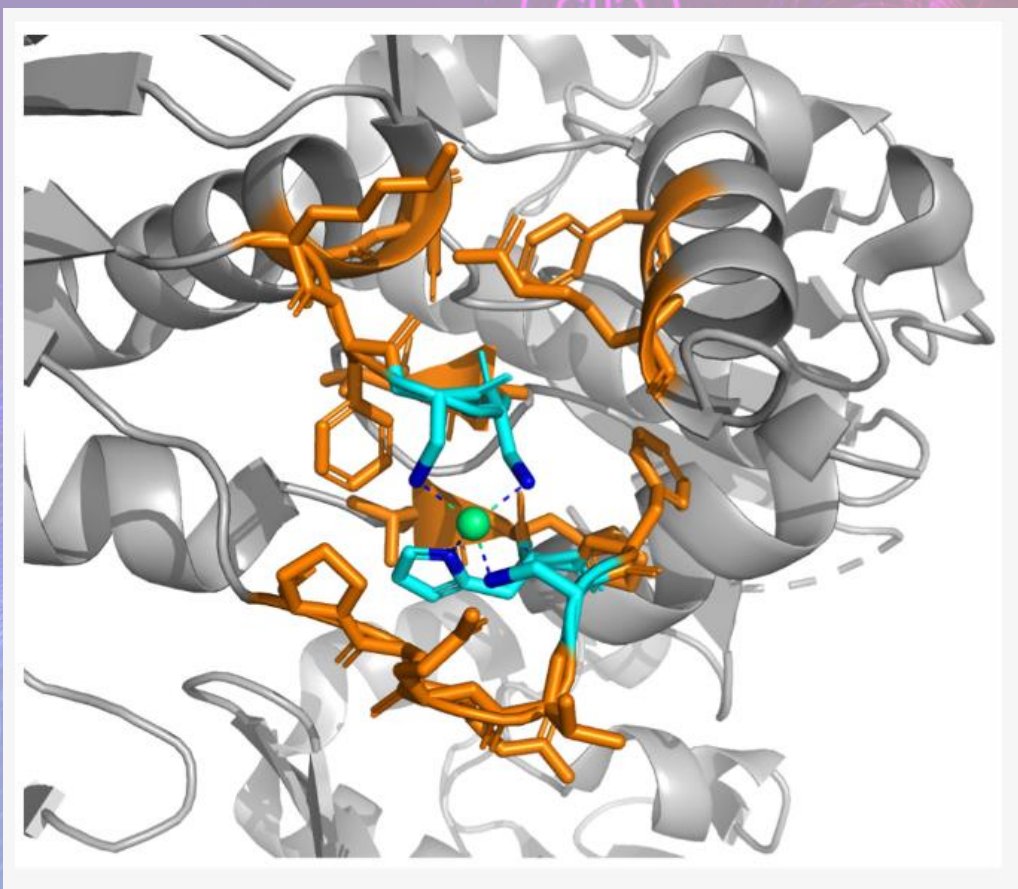
Experiment Details	Results and Performance
Subset of 2982 protein chains analyzed using UniProt	Overall AUC: 0.959
SVM Tools Used	Precision (MBS): 73%, Recall (MBS): 61%
SVMLighty for binary classification (HIS)	Precision (Disulfide bridges): 86%, Recall (Disulfide bridges): 87%
bsvm for multiclass classification (CYS)	Performance loss without descriptors: 0.918 ± 0.004 AUC





Using Simplified Amino Acid Alphabets and Random Forest Model

- Challenges associated with protein 3D structure determination
- Clustering the 20 amino acids into a simplified amino acid alphabet
- Employment of a random forest algorithm





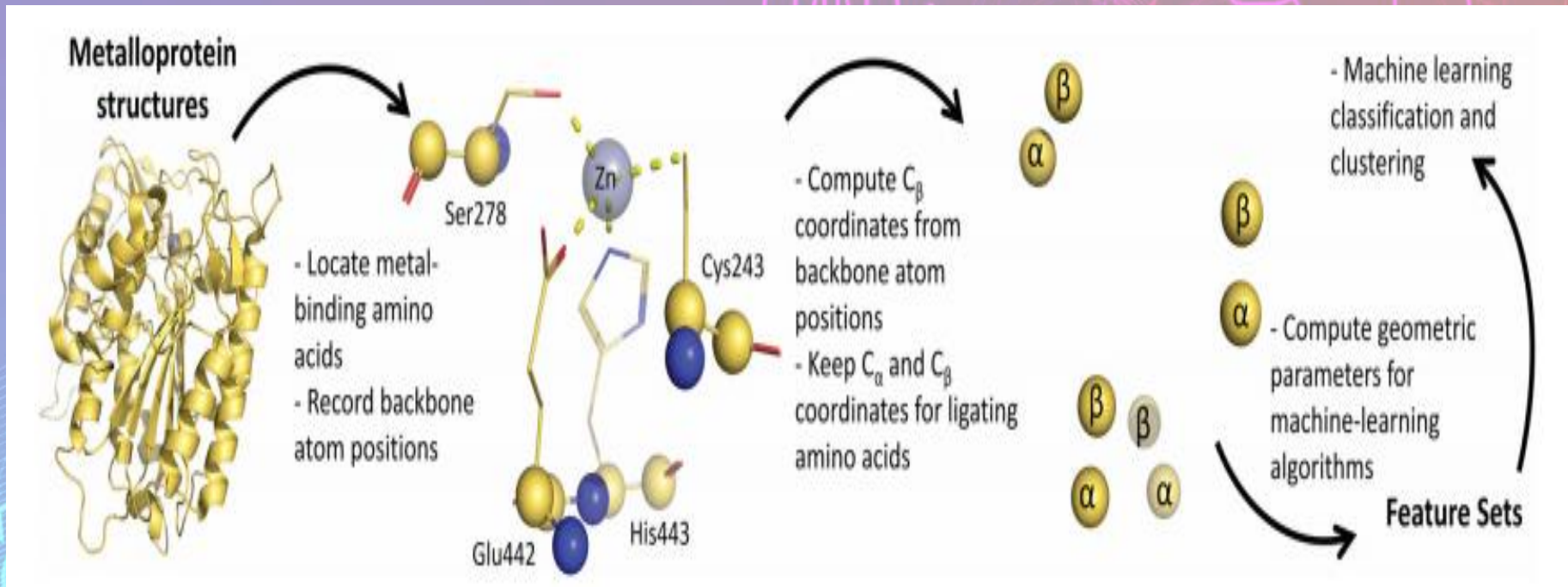
Using Simplified Amino Acid Alphabets and Random Forest Model

None
P
EDNQ
EDRK
PH
CILMV
AG
CFILMVW
NQSTY
CMQLEKRA

Metal Ion	Prediction Accuracy
Iron	69% ★
Copper	75%
Manganese	82%
Magnesium	80%
Nickel	90% ★
Calcium	78%
Cobalt	72%
Zinc	85%

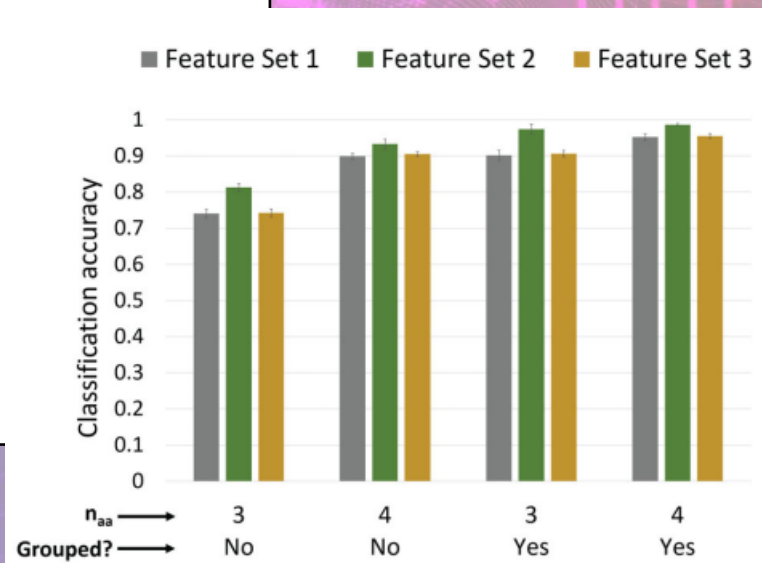
Metal binding amino acids based on backbone geometries

- Decision tree machine-learning algorithm to analyze entire protein structures



Metal binding amino acids based on backbone geometries

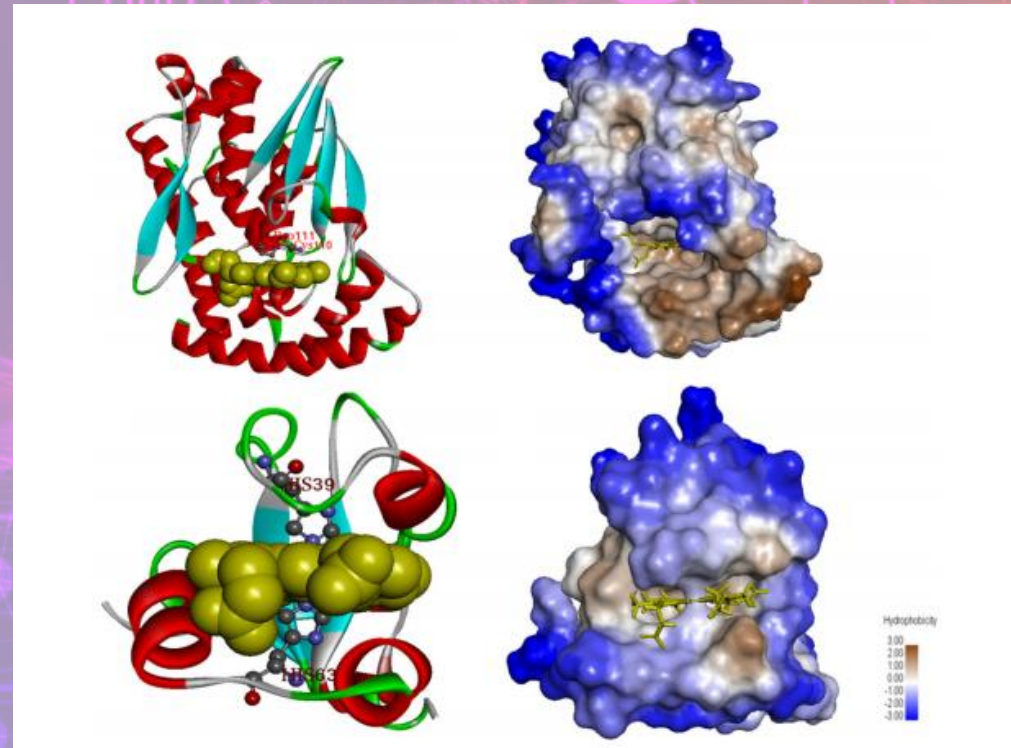
Feature Set	Description	Performance
Feature Set 1	13 independent features representing the backbone geometry of coordinating amino acids	97% accuracy
Feature Set 2	Incorporates order-independent features along with the count of each type of amino acid binding to metals	
Feature Set 3	Considers all possible orderings of the coordinating amino acids, resulting in a larger dataset requiring more computational resources	



Heme binding proteins

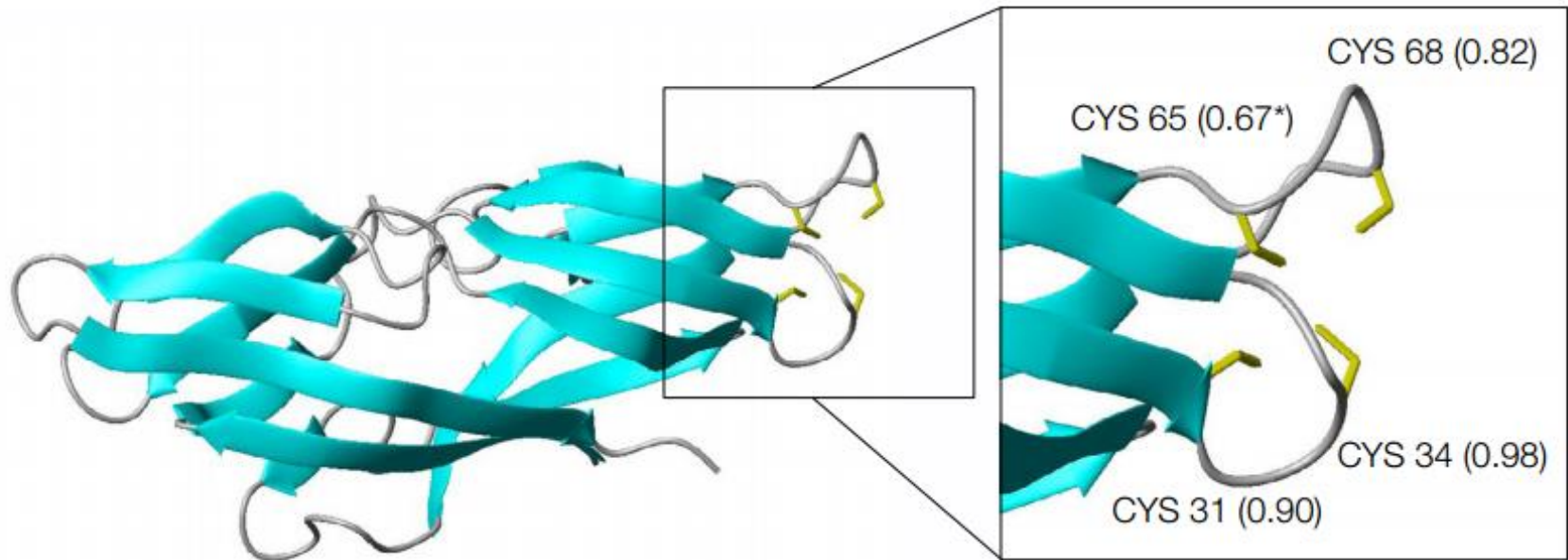
- Heme binding proteins are metalloproteins containing heme ligands.
- Computational methods for predicting heme binding residues are limited.

Attribute	Value/Type
Model	SVM
Dataset Size (HBPs)	747
Dataset Size (Non-HBPs)	91,414
Training Accuracy	85.90%



Predicting zinc binding at the proteome level

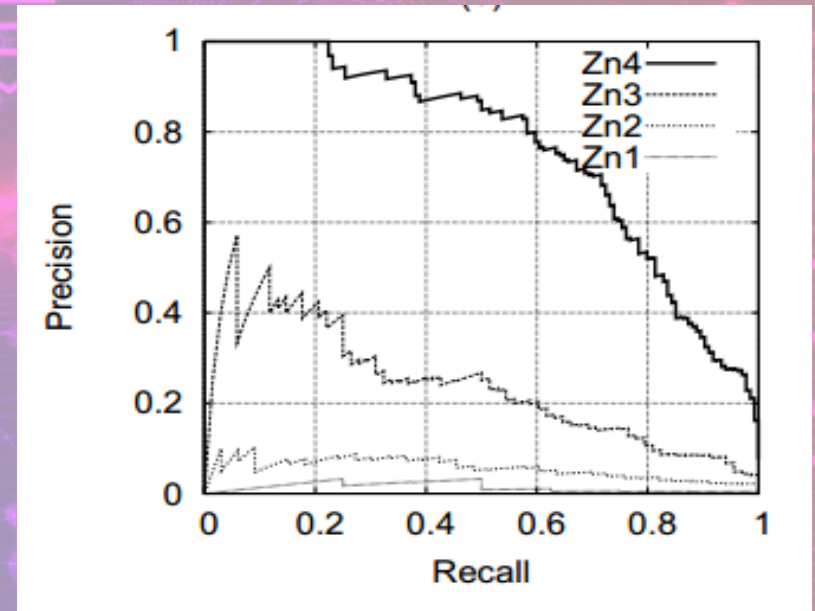
- Zinc has crucial roles in catalytic and structural functions in living organisms
- A SVM approach was developed to predict zinc-binding attitudes of sequential pairs of residues



Predicting zinc binding at the proteome level

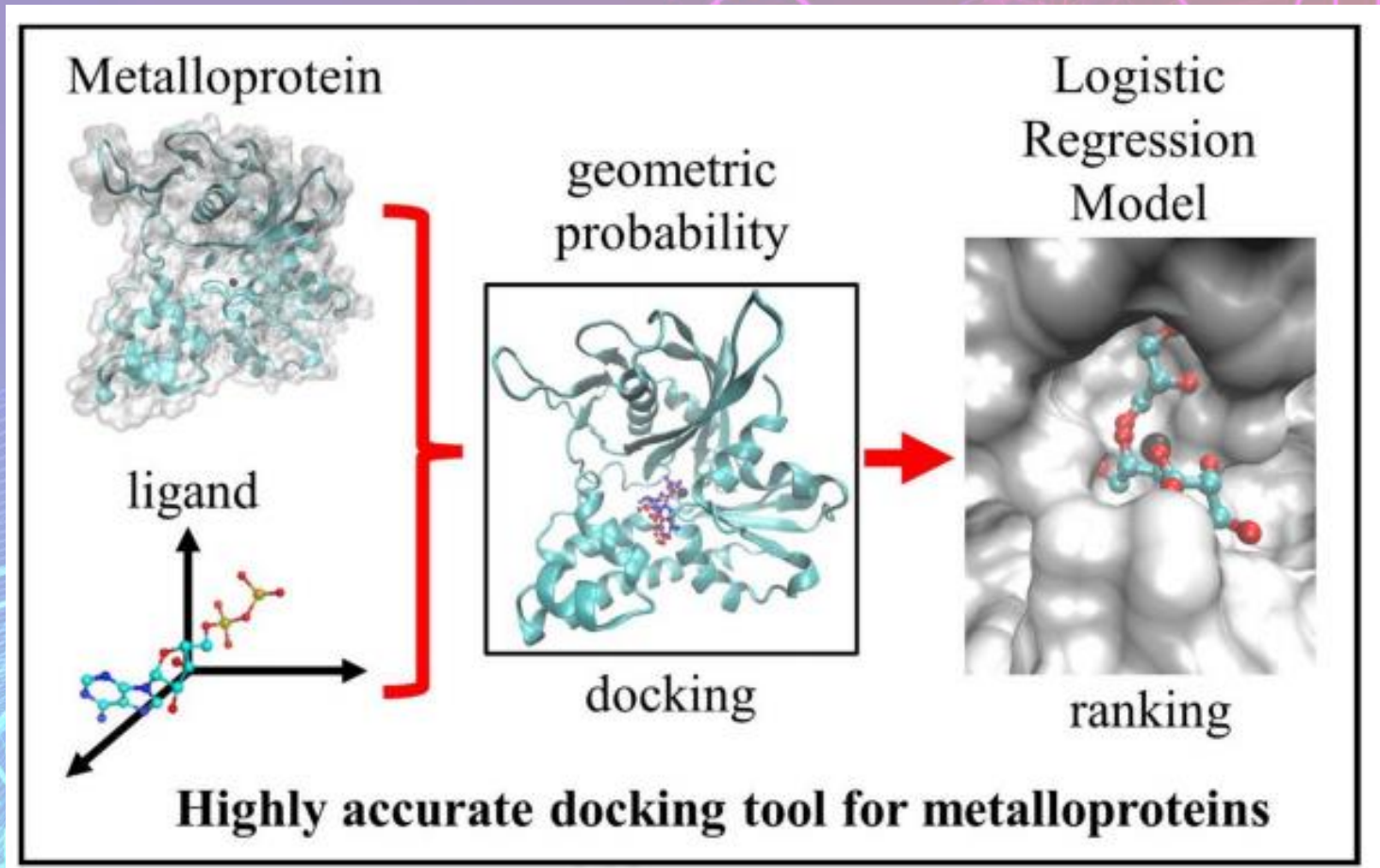
Procedure	Evaluation Metric
Model Selection: SVM	precision of 78% and a recall of 89%

Value
Local Predictor AURPC: 0.428, AUC: 0.867 \pm 0.007
Cross-validation AUC: 0.867 \pm 0.007
2,833 putative zinc-binding human chains



GPDOCK

- A docking method called GPDOCK (Geometric Probability Docking) is introduced, boasting unprecedented accuracy.



GPDOCK

Aspect	Description
Docking Accuracy	GPDOCK achieves 94.3% accuracy in predicting binding poses for 10 metal ions and 9360 complexes.
Docking Capability	Accurately docks metalloproteins with ligands, even with water molecules in metal ion coordination.
Dependency	Relies solely on protein and ligand structures, boosting computational efficiency.
Computational Efficiency	Employs a machine learning model for efficient scoring of binding poses.
Effectiveness	Effective and efficient for drug design and studying metalloprotein binding mechanisms.

GPDOCK

Method	Description
Datasets	Analyzed 48,184 protein structures from PDB (December 2020) with 10 metal ions.
	Coordination information obtained from crystal structure files.
Test Datasets	SM dataset: Metalloproteins with one metal ion in docking pocket.
	MM dataset: Metalloproteins with two to four metal ions in docking pocket.
	SW dataset: Metalloproteins with one metal ion and water coordination in docking pocket.

Future Directions

- **Incorporation of ligand features in articles:** Many existing studies neglect to account for crucial ligand features in their analyses, hindering accurate prediction of binding sites in metalloproteins.
- **Integration of ion coordination:** Current models often overlook the coordinated fixation of ions, which is vital for understanding metalloprotein function and should be incorporated into future machine learning approaches.
- **Machine learning-based design of new protein sequences:** Future research should explore machine learning techniques to design novel protein sequences with optimized binding capabilities for specific ligands, thus advancing protein engineering efforts in the development of metalloenzymes.



PubMed Keywords

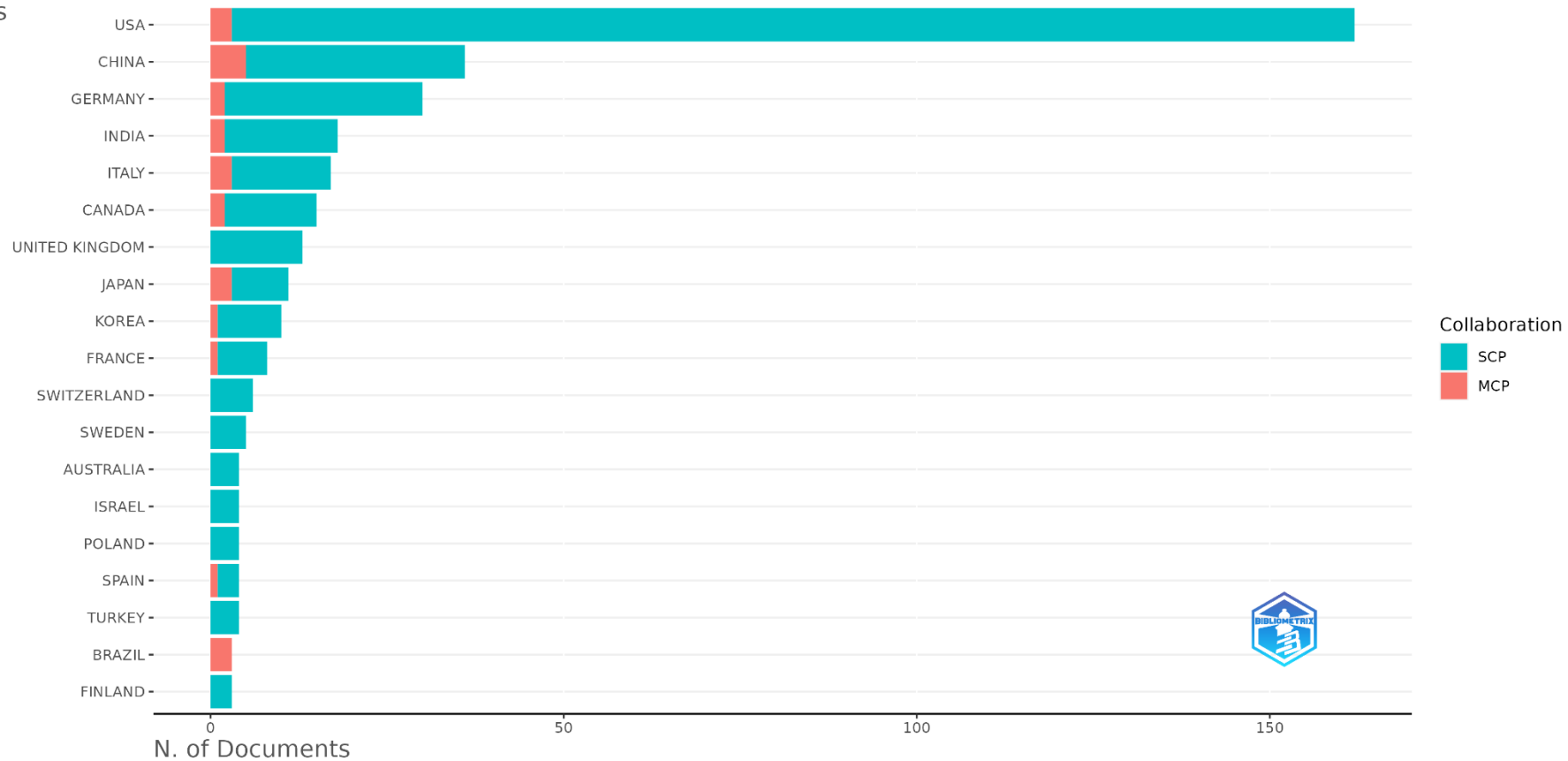
- **Keywords Searched:**
 - **Machine Learning**
 - **Metalloproteins**
 - **Binding Sites**
 - **Forecasting**
 - **Machine Learning AND Metalloproteins AND (Binding Sites AND (Forecasting OR Prediction))**



Research Landscape

Corresponding Author's Countries

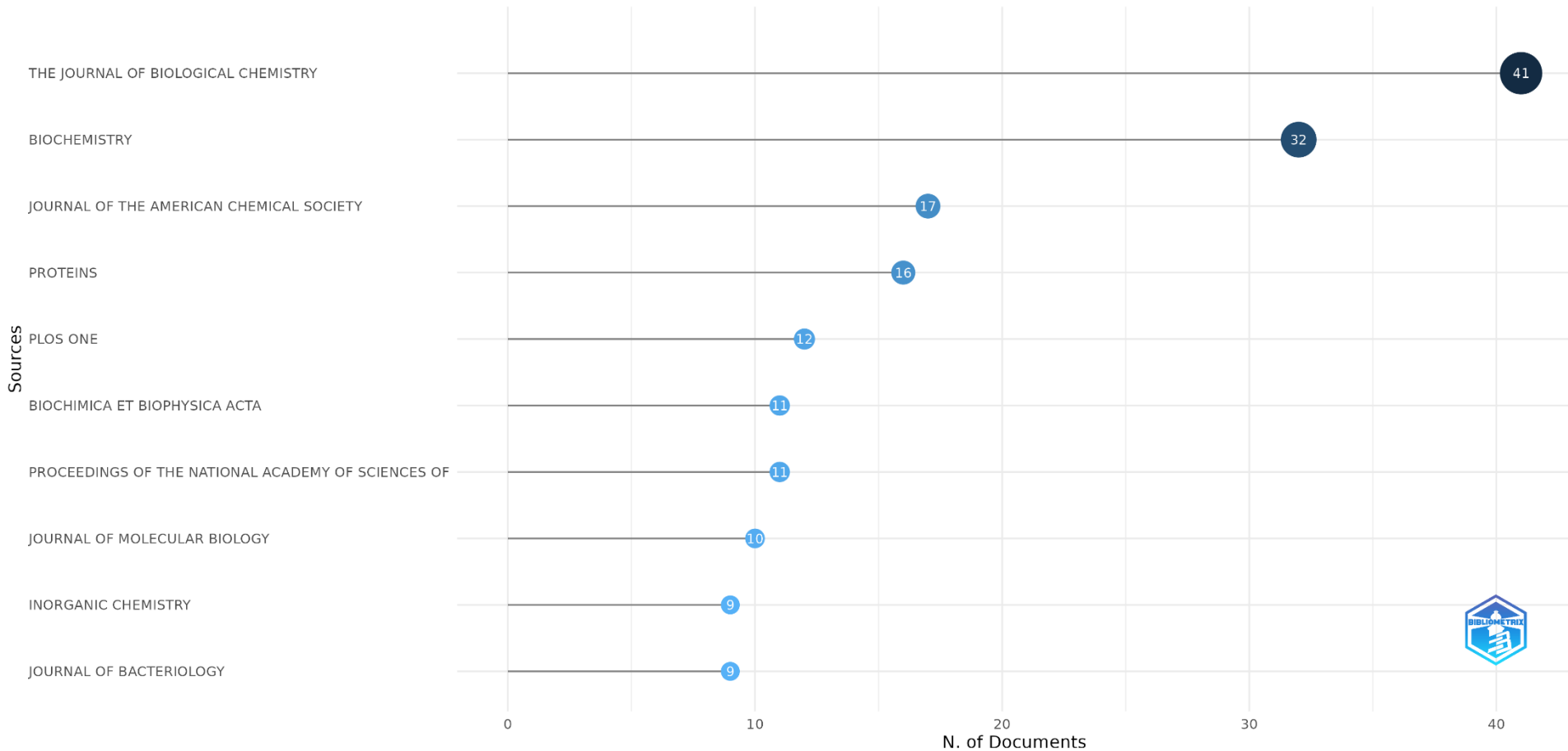
Countries

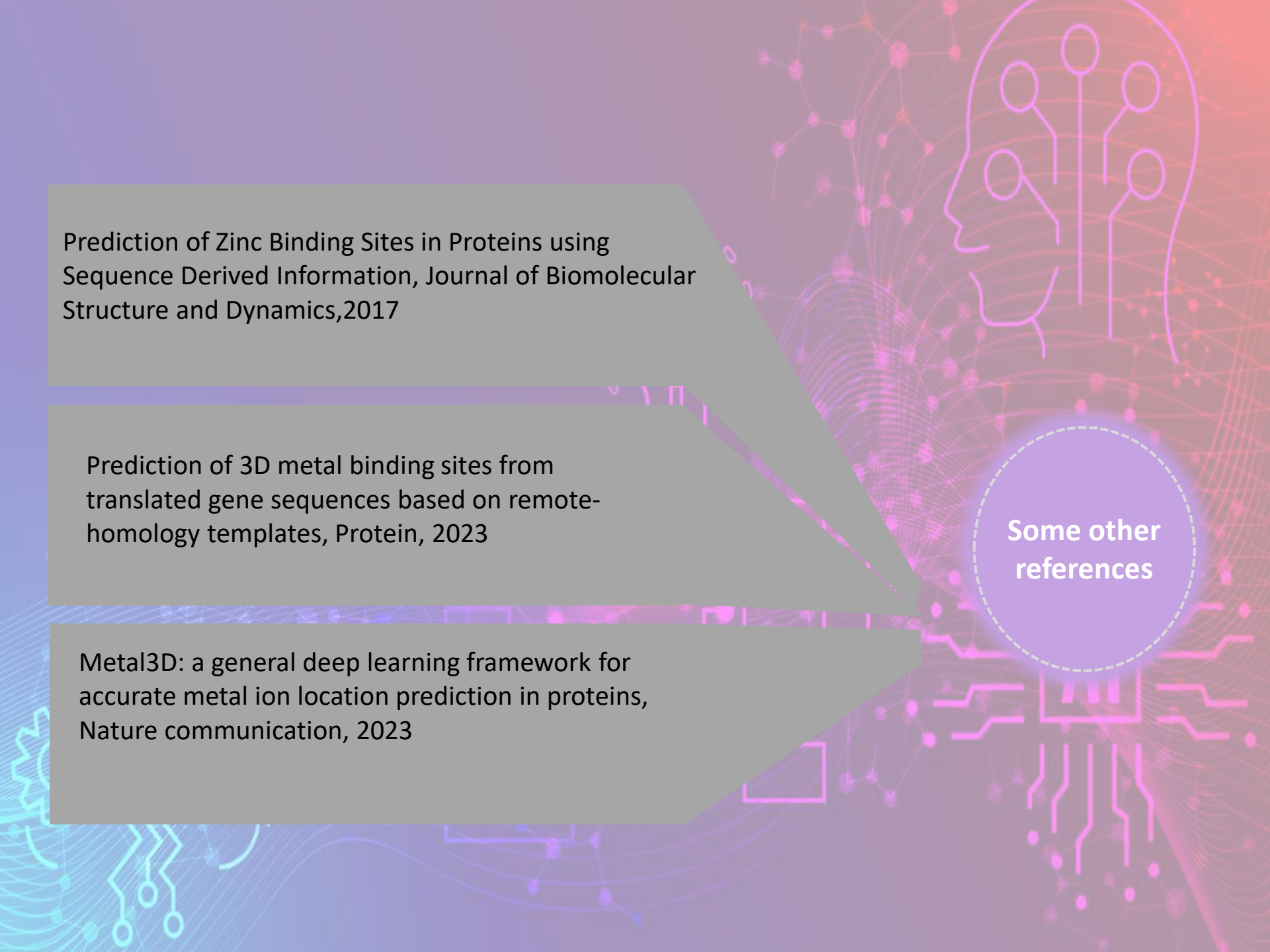


SCP: Single Country Publications, MCP: Multiple Country Publications

Research Landscape

Most Relevant Sources





Prediction of Zinc Binding Sites in Proteins using Sequence Derived Information, Journal of Biomolecular Structure and Dynamics, 2017

Prediction of 3D metal binding sites from translated gene sequences based on remote-homology templates, Protein, 2023

Metal3D: a general deep learning framework for accurate metal ion location prediction in proteins, Nature communication, 2023

**Some other
references**

