# Identification of catalytic residues in protein structures using a novel feature that integrates the micro-environment and geometrical location properties of residues

Lei Han[1], Jiangning Song[2,3], Ming S. Liu[4,*] and Ziding Zhang[1,*]

[1]State Key Laboratory of Agrobiotechnology, College of Biological Sciences, China Agricultural University, Beijing 100193, China  [2]National Engineering Laboratory for Industrial Enzymes and Key Laboratory of Systems Microbial Biotechnology, Tianjin Institute of Industrial Biotechnology, Chinese Academy of Sciences, Tianjin 300308, China [3]Department of Biochemistry and Molecular Biology, Faculty of Medicine, Monash University, Melbourne, VIC 3800, Australia [4]CSIRO - Mathematics, Informatics & Statistics, Private Bag 33, Clayton South 3169, Australia

## Introduction

Enzymes help to fulfill diverse biochemical functions and play critical roles in almost all cellular processes. Although different mechanisms of some enzymes have been characterized，there is a difficulty in rationalizing the available enzyme sequences/structures with their annotated functions. Identification of catalytic residues (CRs) is the first and important step to characterize the catalytic mechanism and function of an enzyme. Since experimental determination of CRs on large-scale proteome data is still a costly and daunting task, computational methods to predict CRs from enzyme sequences/structures are playing an increasingly important role in complementing the experimental identification. Though there are lots of features for predicting CRs, advancing novel features does not only increase the prediction accuracy but also deepen our understanding of catalytic mechanisms.

In this work, we developed a novel structural feature called MEDscore to determine the CRs. Firstly, a residue's micro-environment (ME) was converted into a series of spatially neighboring residue pairs and a ME-based score (i.e. MEscore) was proposed to quantify a residue's ME information. Secondly, a parameter named Dscore was set up to measure a residue's global positional information. Finally, MEDscore was defined from an effective nonlinear integration of MEscore and Dscore.

## Materials and Methods

### Benchmark Enzyme Dataset

The enzyme dataset used in this study was based on the Catalytic Site Atlas (CSA) database (version 2.2.12).After filtering homology sequences and based on other criteria, 223 enzyme structure domains were retained in our final dataset. These structure domains were downloaded from ASTRAL.
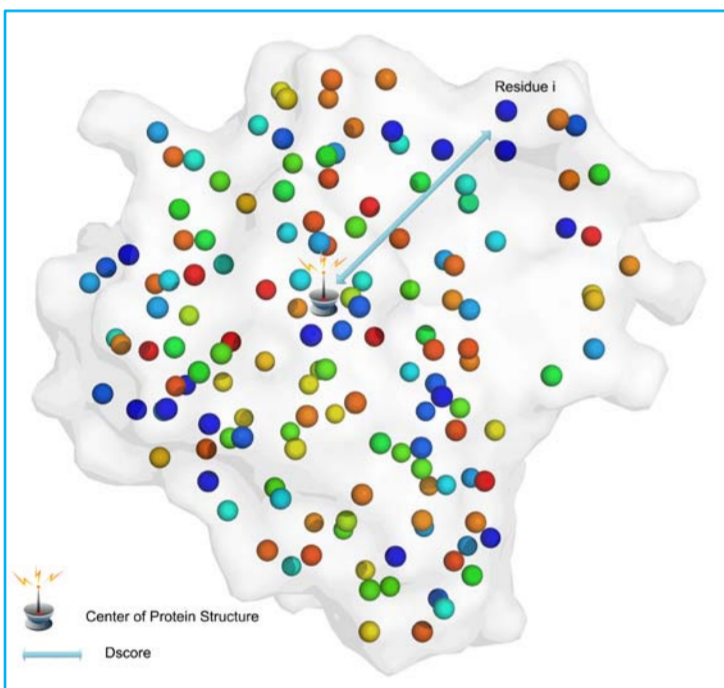
### Definition and Calculation of Dscore
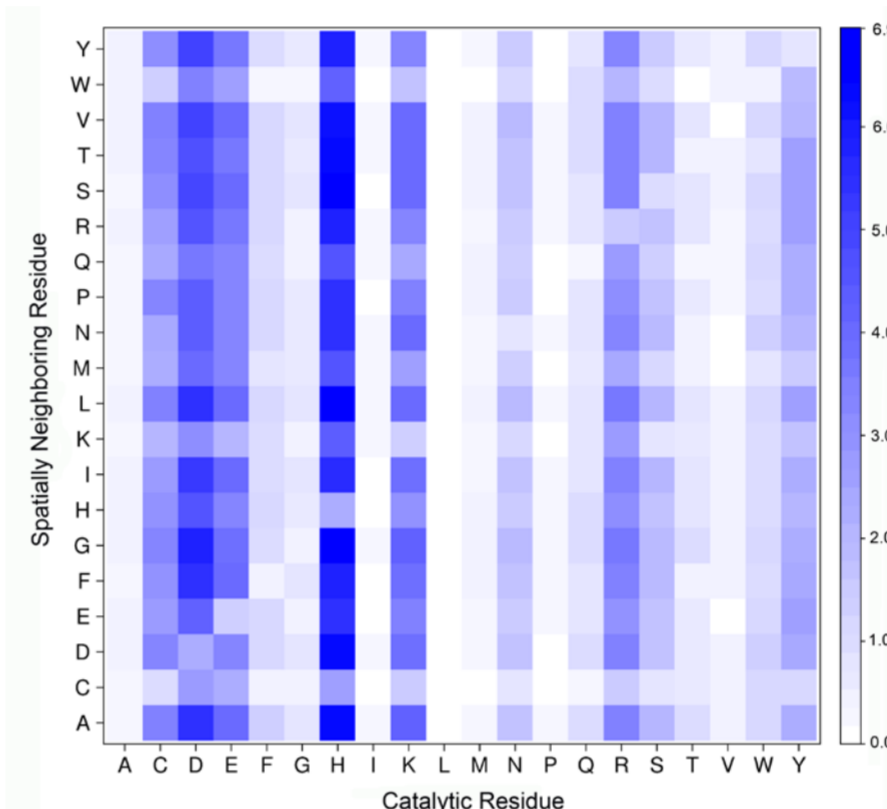


Geometrical center of one protein structure:

$$(c_x, c_y, c_z) = (\frac{\sum_{i=1}^{N} x_i}{N}, \frac{\sum_{i=1}^{N} y_i}{N}, \frac{\sum_{i=1}^{N} z_i}{N})$$

The distance between a residue i and the center of the structure (i.e. $Dscore_i$)

$$Dscore_i = \sqrt{(c_x - x_i)^2 + (c_y - x_i)^2 + (c_z - x_i)^2}$$

### Definition and Calculation of MEDscore

MEDscore is a feature which integrates the micro-environment (ME) and geometrical properties of amino acid residues. The flow-chart of calculating MEDscore is shown at right.

The ME of a residue (MEscore) is represented as a series of spatially neighboring residue pairs in the local of the query residue. The details of calculation MEDscore is similar as the process of measuring MEDscore.



## Results and Discussion

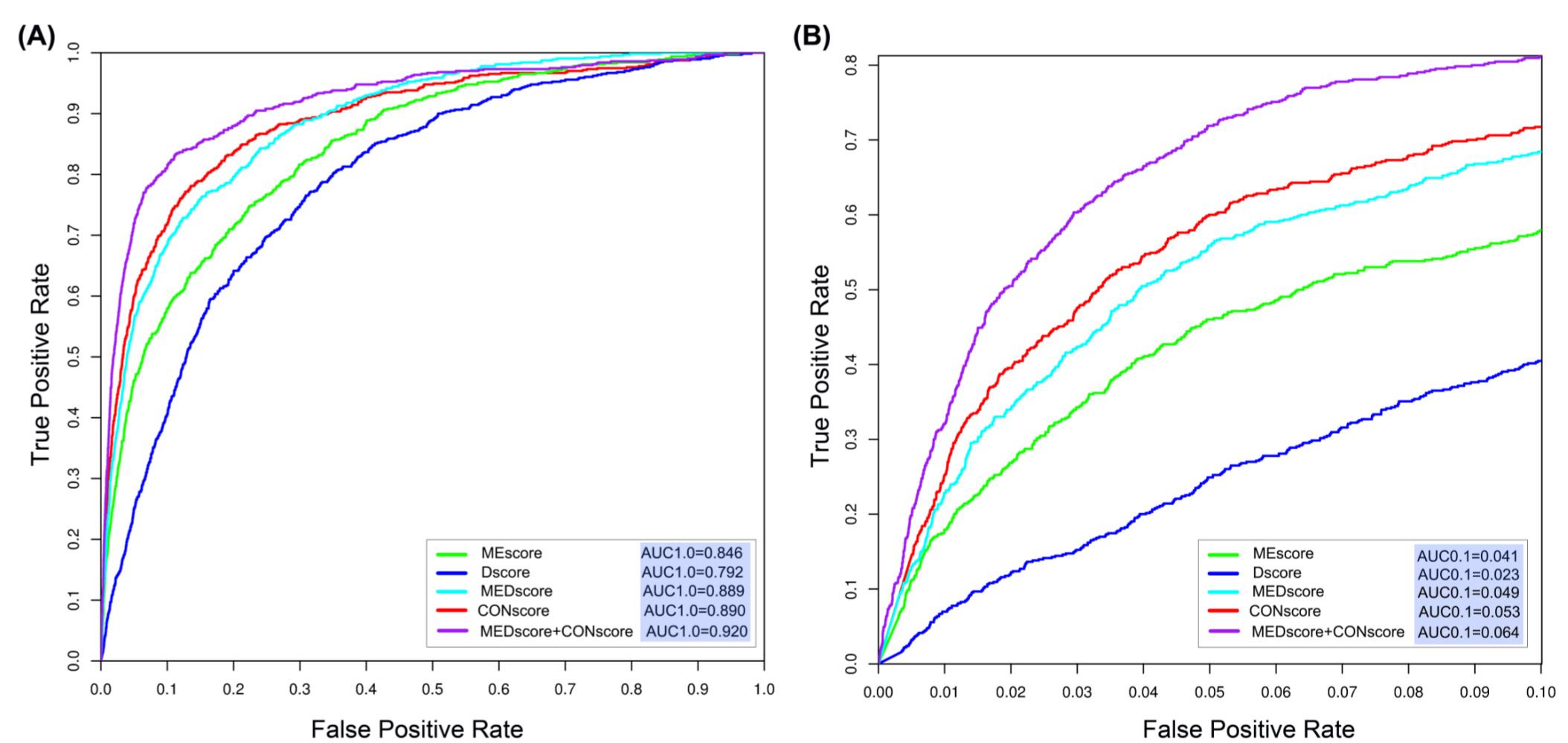### Propensities of residues in the ME surrounding the CRs



### The weight coefficients of spatially neighboring residue pairs in the MEs of CRs



As shown in the left figure, different residue pairs exhibit the scaled propensities in the ME of catalytic residues, providing important insights into the molecular mechanism of enzymatic catalysis.
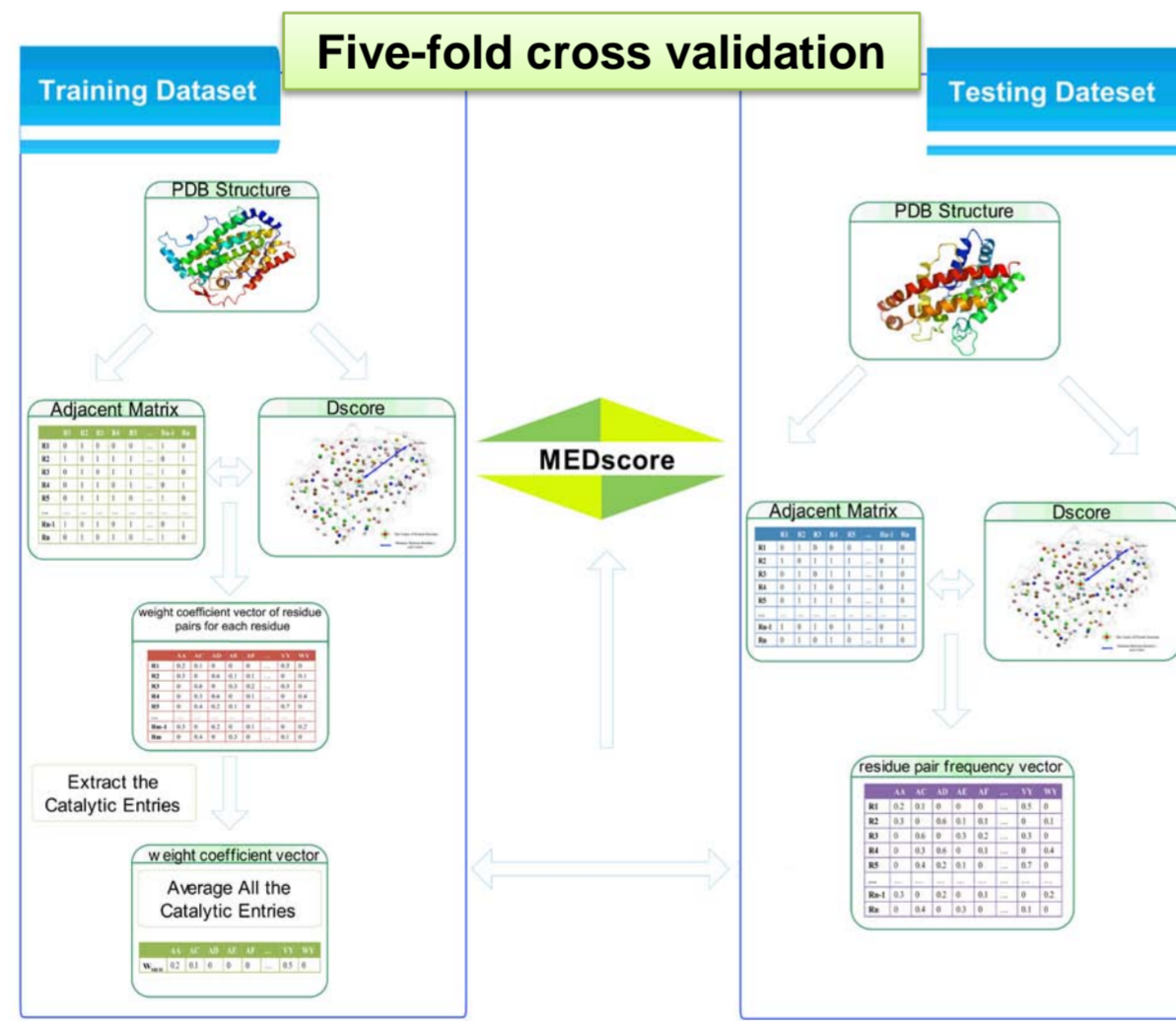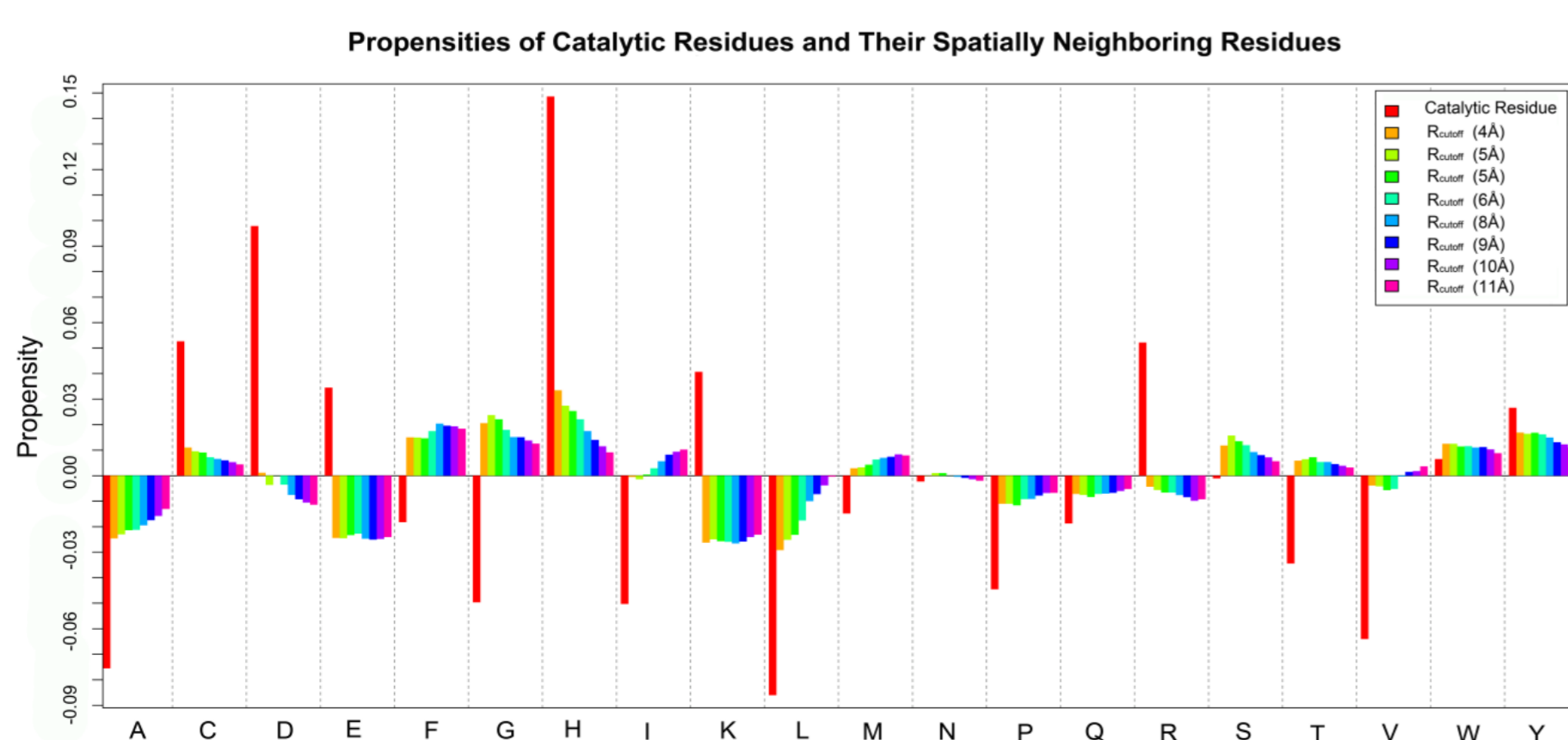
### Performance comparison of each feature



### Web server for CRs prediction (http://protein.cau.edu.cn/mepi/)



## Acknowledgements

Ziding Zhang's Lab of Protein Bioinformatics  http://protein.cau.edu.cn  Lei HAN's email : hanlei_45@126.com