# 大数据时代下基于人工智能算法的蛋白质功能预测研究

南京农业大学 人工智能学院

汇报人：朱一亨

2024年05月31日

# 01 Part one

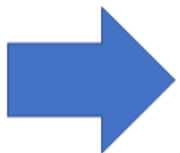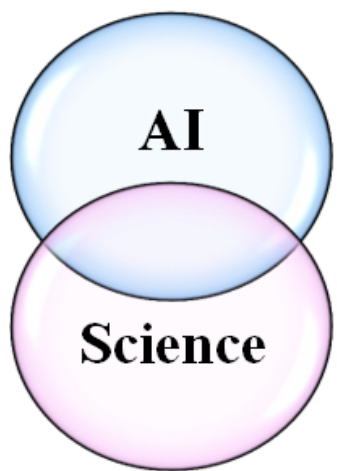# 研究背景

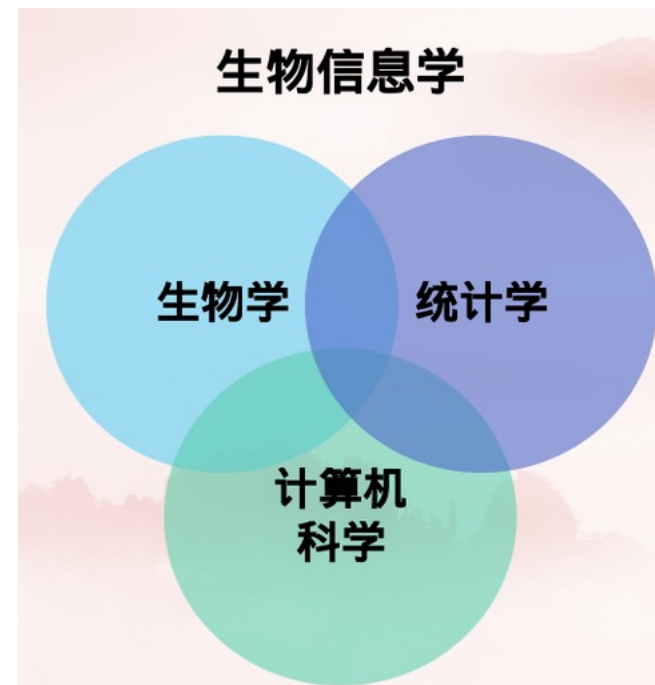➢ 生物信息学：生物学＋统计学＋计算机科学 <span style="color:red">揭示</span>⟶ 生物数据中所蕴含的生物学奥秘

➢ 国家发改委（2022）：首部"十四五"生物经济发展规划

➢ 中国科协（2022）：重大前沿科学问题之一



AI / Science ⟶
- 探究生命科学中的基础问题
- 促进生物医药产业技术升级
- 创新人工智能领域前沿算法



生物信息学
- 生物学
- 统计学
- 计算机科学

➤ 生物信息学的体系分类

```
基因组学

生物信息学          蛋白质结构分析

              蛋白质组学

                        蛋白质功能分析
```

后基因时代的研究重点

当前的研究课题：蛋白质功能预测

➢ 蛋白质是生命现象的物质基础之一

➢ 蛋白质参与了生物体内几乎全部的生命过程，并发挥着各种重要的功能

➢ 蛋白质是生命活动的主要承担者

DNA —Transcription 转录→ RNA —Translation 翻译→ Protein 蛋白质 —调控→
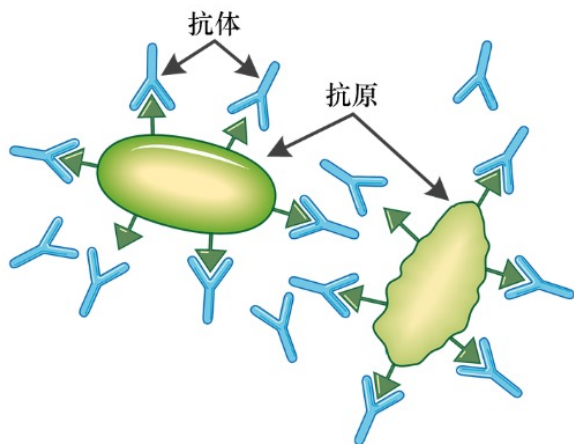
Reverse Transcription 逆转录

生物中心法则

人体代谢活动

➢ 识别和分析蛋白质的功能有助于解释各种生命活动现象，并阐明相关疾病的发病机理，进而指导相应的药物设计，以期推动智能医疗的发展。

➢ 蛋白质功能注释是后基因时代的首要任务之一。



(a) 催化反应　　(b) 免疫保护　　(c) 运输载体　　(d) 基因调控
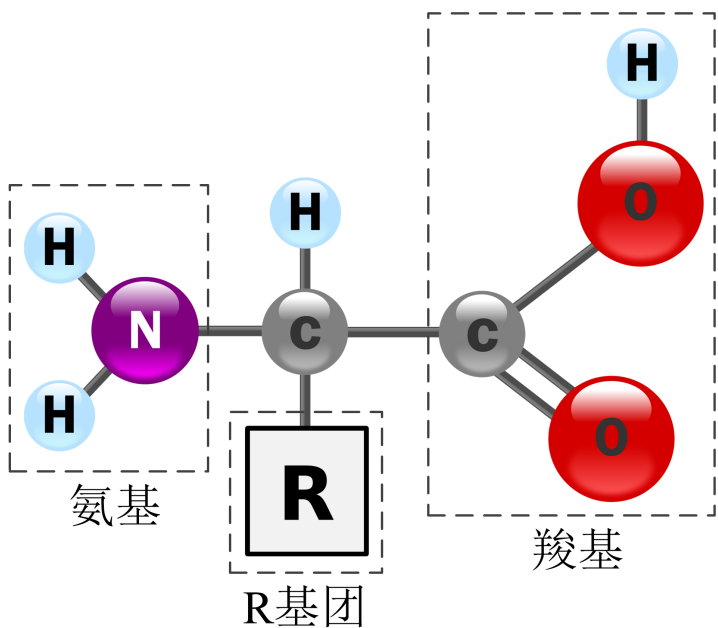
蛋白质的生物功能

➤ **蛋白质的基本组成单位是氨基酸**

氨基

R基团

羧基

H O

H C C O

H N R

H H

**氨基酸的化学分子式**

氨基酸信息汇总表

■ 碱性氨基酸　　■ 极性氨基酸（不电离）
■ 非极性氨基酸　　■ 酸性氨基酸

H 155.16 137.14 C$_6$H$_7$N$_3$O$_2$ **His** 组氨酸 Histidine

D 133.10 115.09 C$_4$H$_5$NO$_4$ **Asp** 天门冬氨酸 Aspartic Acid

R 174.20 156.19 C$_6$H$_{14}$N$_4$O$_2$ **Arg** 精氨酸 Arginine

F 165.19 147.18 C$_9$H$_{11}$NO$_2$ **Phe** 苯丙氨酸 Phenylalanine

A 89.09 71.08 C$_3$H$_7$NO$_2$ **Ala** 丙氨酸 Alanine

C 121.16 103.14 C$_3$H$_7$NO$_2$S **Cys** 半胱氨酸 Cysteine

G 75.07 57.05 C$_2$H$_5$NO$_2$ **Gly** 甘氨酸 Glycine

Q 146.15 128.13 C$_5$H$_{10}$N$_2$O$_3$ **Gln** 谷氨酰胺 Glutamine

E 147.13 129.11 C$_5$H$_9$NO$_4$ **Glu** 谷氨酸 Glutamic Acid

K 146.19 128.17 C$_6$H$_{14}$N$_2$O$_2$ **Lys** 赖氨酸 Lysine

L 131.17 113.16 C$_6$H$_{13}$NO$_2$ **Leu** 亮氨酸 Leucine

M 149.21 131.20 C$_5$H$_{11}$NO$_2$S **Met** 甲硫氨酸 Methionine

N 132.12 114.10 C$_4$H$_8$N$_2$O$_3$ **Asn** 天门冬酰胺 Asparagine

S 105.09 87.08 C$_3$H$_7$NO$_3$ **Ser** 丝氨酸 Serine

Y 181.19 163.17 C$_9$H$_{11}$NO$_3$ **Tyr** 酪氨酸 Tyrosine

T 119.12 101.10 C$_4$H$_9$NO$_3$ **Thr** 苏氨酸 Threonine

I 131.18 113.16 C$_6$H$_{13}$NO$_2$ **Ile** 异亮氨酸 Isoleucine

W 204.23 186.21 C$_{11}$H$_{12}$N$_2$O$_2$ **Trp** 色氨酸 Tryptophan

P 115.13 97.12 C$_5$H$_9$NO$_2$ **Pro** 脯氨酸 Proline

V 117.15 99.13 C$_5$H$_{11}$NO$_2$ **Val** 缬氨酸 Valine

单字母缩写 — S — 三字母缩写
分子量
残基分子量 105.09 87.08 C$_3$H$_7$NO$_3$ **Ser**
分子式 — — 化学结构式
丝氨酸 Serine — 化学名

**20种常见氨基酸**

- 氨基酸脱水缩合组成肽链

- 肽链相互缠绕组成蛋白质分子



肽链 ——相互缠绕——>

蛋白质分子

➢ 氨基酸四级结构

（1）一级结构：氨基酸的线性序列

（2）二级结构：肽链上的局部几何构象

（3）三级结构：肽链上所有原子的空间位置

（4）四级结构：肽链在空间上的相对位置

一级结构
氨基酸顺序

二级结构
常规次结构

α螺旋

β折叠

三级结构
三维结构

P13蛋白

四级结构
蛋白质分子复合体

血红蛋白

➢ **蛋白质功能注释的语义词汇标准**

——基因本体论

**(1) 分子功能**

Molecular Function

**(2) 生物过程**

Biological Process

**(3) 细胞组件**

Cellular Component



(A) 分子功能　　　　(B) 生物过程　　　　(C) 细胞组件
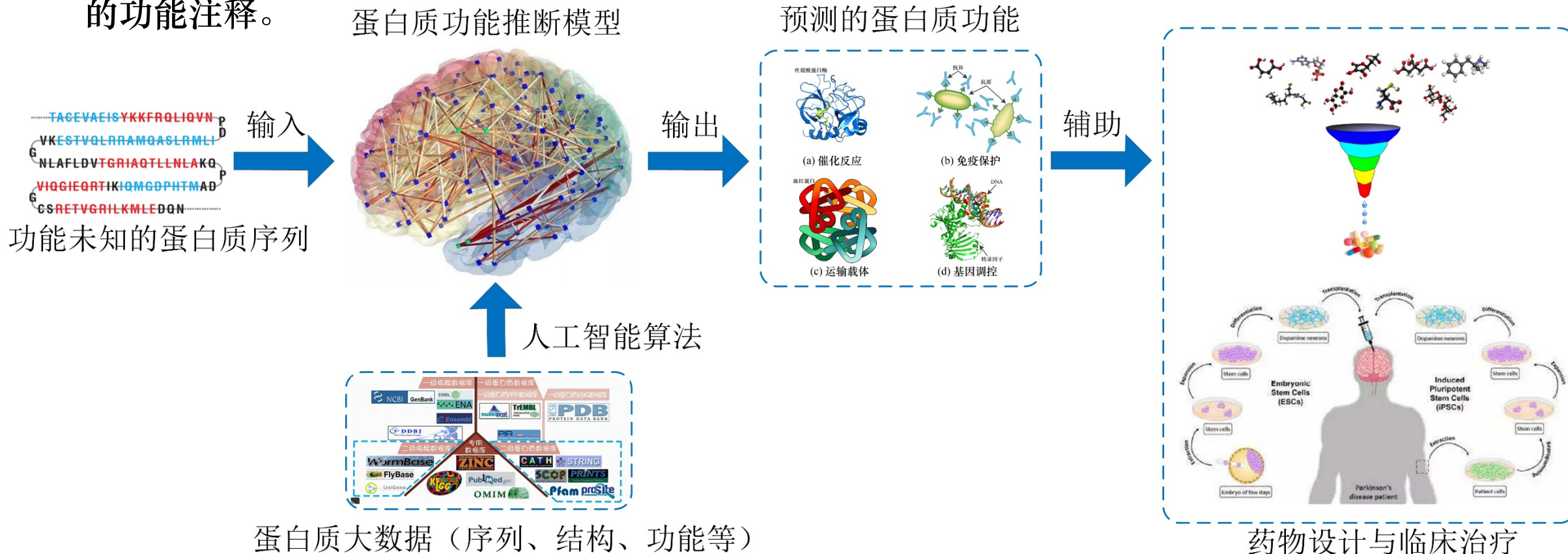
**谷氨酰环化转移酶在分子功能、生物过程和细胞组件分支下的功能注释图**

# 蛋白质功能注释面临的挑战

➤ 蛋白质功能注释最可靠的途径是生物实验，但它存在周期长、成本高等缺陷。

➤ 截止2024年1月，UniProt中已累积约2.52亿条序列，但具有生物实验功能注释的序列数目不足序列总数的0.1%。

➤ 研发高效的生物计算方法来预测蛋白质功能已迫在眉睫。

➤ 蛋白质功能预测目标：利用生物计算方法准确地推断出查询蛋白质在基因本体论三个分支下的功能注释。



蛋白质功能推断模型

预测的蛋白质功能

输入

输出

辅助

功能未知的蛋白质序列

人工智能算法

蛋白质大数据（序列、结构、功能等）

药物设计与临床治疗

➢ 数学

　高等数学、线性代数、概率论等

➢ 计算机科学

　编程语言(python、C++、R)、数据结构、操作系统等

➢ 人工智能算法

　机器学习、深度学习等

➢ 生物学

　蛋白质 (序列、结构和功能)、基因等基础知识

# 02 Part two
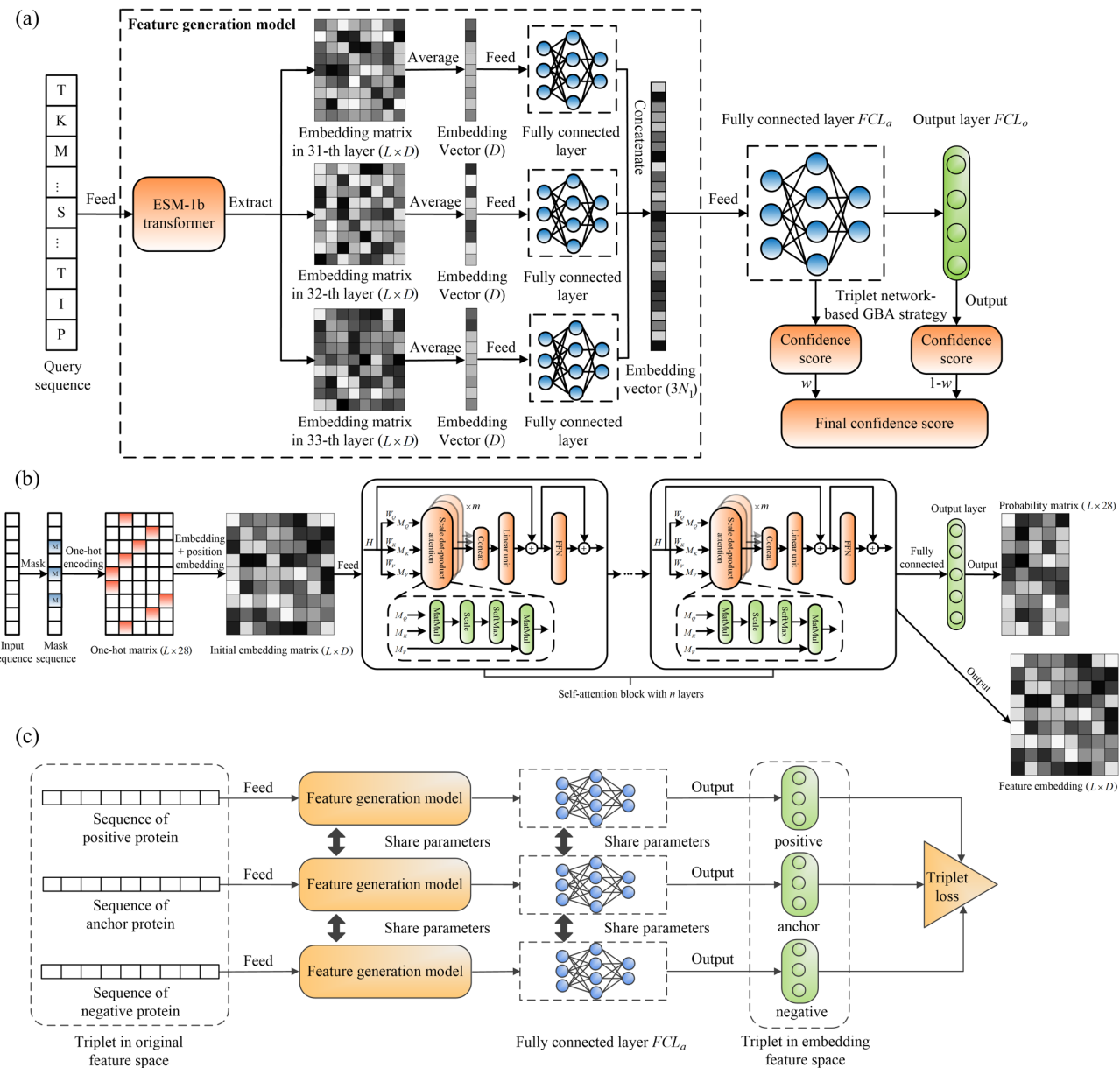
# 研究内容

蛋白质功能预测

从蛋白质视角出发预测功能

从基因视角出发预测功能

从配体视角出发预测功能

➢ 基于注意力机制与三元组神经网络的
  预测方法 ATGO

➢ 主要贡献：首次将计算机视觉领域的
  无监督语言模型迁移到蛋白质功能预
  测领域

Yi-Heng Zhu, Chengxin Zhang, Dong-Jun Yu, Yang Zhang. Integrating Unsupervised Language Model with Triplet Neural Networks for Protein Gene Ontology Prediction. **PLOS Computational Biology**. 2022, 18(12): e1010793.

# 从蛋白质视角出发预测功能

➤ 基于度量学习与多源信息融合的功能预测
方法 TripletGO

➤ 主要贡献：首次将基于基因视角的方法和
基于蛋白质视角的相结合，为后续的研究
开辟了新的思路

Yi-Heng Zhu, Chengxin Zhang, Yan Liu, Gilbert Omenn, Peter Freddolino, Dong-Jun Yu, Yang Zhang. TripletGO: Integrating Transcript Expression Profiles with Protein Homology Inferences for Gene Function Prediction. Genomics, **Proteomics & Bioinformatics**. 2022, 20(5): 1013-1027.

# 从基因视角出发预测功能

➤ 基于多粒度支持向量机集成与序列特征的蛋白质-DNA绑定位点预测方法 DNAPred

➤ 主要贡献：提出了新的类不平衡学习算法E-HDSVM，显著地提升了蛋白质-DNA绑定位点预测精度。

Yi-Heng Zhu, Jun Hu, Xiao-Ning Song, Dong-Jun Yu. DNAPred: Accurate Identification of DNA-binding Sites from Protein Sequence by Ensembled Hyperplane-Distance-Based Support Vector Machines. **Journal of Chemical Information and Modeling**. 2019, 59:3057-3071.


DNA
蛋白质-DNA绑定位点
蛋白质



(1) Training stage

PSI-BLAST → PSSM
PSIPRED → PSS
SANN → PRSA
AAFD-BN

Training sequences

Serial concatenation + Slide window

Training feature vectors for residues

(2) Prediction stage

Query sequence with length $L$
V S R T I ⋮ Y F Q D

PSI-BLAST → PSSM
PSIPRED → PSS
SANN → PRSA
AAFD-BN

Serial concatenation + Slide window

Feature vectors for residues

E-HDSVM

HD-US

Subset1 Subset2 ⋯ Subset $N$

$HDSVM_1$ $HDSVM_2$ ⋯ $HDSVM_N$

AdaBoost

Feed

$AdaHDSVM$

Predict

DNA-binding
Non-DNA binding
1
$T$
0
1 2 3 4 ⋯ $L$

**DNAPred**: Identifying DNA-Binding Sites from Protein Sequence by Ensemble Hyperplane-Distance-Based Support Vector Machine

| Read Me | Dataset | Citation | Large-Scale Test |

**Input query protein sequence(s) in FASTA format:**

```
>2XTNA
MDQNEHSHWGPHAKGQCASRSELRIILVGKTGTGKSAAGNSILRKQAFESKLGS
QTLTKTCSKSQGSWGNREIVIIDTPDMFSWKDHCEALYKEVQRCYLLSAPGPHV
LLLVTQLGRYTSQDQQAAQRVKEIFGEDAMGHTIVLFTHKEDLNGGSLMDYMH
DSDNKALSKLVAACGGRICAFNNRAEGSNQDDQVKELMDCIEDLLMEKNGDHY
TNGLYSLIQRSKCGPVGSDE
```

[Example]  [Reset Sequence(s)]

**Choose a prediction model**

◉ Model constructed on PDNA-543      ○ Model constructed on PDNA-335

**Choose a threshold**

◉ Threshold 1 (*Max MCC*)      ○ Threshold 2 (*FPR≈5%*)      ○ Threshold 3 (*Sen≈Spe*)

**Email Address (For receiving your prediction results)\***

[                                    ]

[Submit]  [Clear All]

**Reference:**
Yi-Heng Zhu, Jun Hu, Xiao-Ning Song and Dong-Jun Yu *. DNAPred: Identifying DNA-Binding Sites from Protein Sequence by Ensemble Hyperplane-Distance-Based Support Vector Machine. Journal of Chemical Information and Modeling, 2019.

---

**RESULTS PAGE**

Predicting Protein-DNA Binding Sites

**Protein Name**

2XTNA

**Model constructed on Dataset**

PDNA-543

**Threshold**

0.265 (*Max MCC*)

**Prediction Summary**

Number of predicted DNA-binding residues in protein **2XTNA**: **4**

Specific position: **58 T  117 R  119 T  147 H**

**Predicted Results**

| Residue # | Amino Acid Type | Probability | Binding Residue |
|-----------|-----------------|-------------|-----------------|
| 0001 | M | 0.058 | N |
| 0002 | D | 0.028 | N |
| 0003 | Q | 0.024 | N |
| 0004 | N | 0.049 | N |
| 0005 | E | 0.013 | N |
| 0006 | H | 0.063 | N |
| 0007 | S | 0.008 | N |
| 0008 | H | 0.037 | N |
| 0009 | W | 0.095 | N |
| 0010 | G | 0.009 | N |
| 0011 | P | 0.019 | N |
| 0012 | H | 0.081 | N |
| 0013 | A | 0.017 | N |
| 0014 | K | 0.080 | N |
| 0015 | G | 0.006 | N |
| 0016 | Q | 0.013 | N |

- 基于无监督语言模型与多源信息融合的蛋白质-DNA绑定位点预测方法 ULDNA

- 主要贡献：融合多种无监督蛋白质语言模型，显著地提升了蛋白质-DNA绑定位点预测精度。



Yi-Heng Zhu, Zi Liu, Zhiwei Ji, Dong-Jun Yu. ULDNA: Integrating Unsupervised Multi-Source Language Models with LSTM-Attention Network for High-Accuracy Protein-DNA Binding Site Prediction. **Briefings in Bioinformatics**. 2024, 25(2):bbae040.

**ULDNA: Integrating Unsupervised Multi-Source Language Models with LSTM-Attention Network for Protein-DNA Binding Site Prediction**

| Read Me | Dataset | Citation |

**Input query protein sequence(s) in FASTA format:**

>2XTNA
MDQNEHSHWGPHAKGQCASRSELRIILVGKTGTGKSAAGNSILRKQAFESKLGS
QTLTKTCSKSQGSWGNREIVIIDTPDMFSWKDHCEALYKEVQRCYLLSAPGPHV
LLLVTQLGRYTSQDQQAAQRVKEIFGEDAMGHTIVLFTHKEDLNGGSLMDYMH
DSDNKALSKLVAACGGRICAFNNRAEGSNQDDQVKELMDCIEDLLMEKNGDHY
TNGLYSLIQRSKCGPVGSDE

[ Example ]  [ Reset Sequence(s) ]

**Choose a prediction model**
- Model constructed on PDNA-543
- Model constructed on PDNA-335

**Choose a threshold**
- Threshold 1 (*Max MCC*)
- Threshold 2 (*FPR≈5%*)
- Threshold 3 (*Sen≈Spe*)

**Email Address (For receiving your prediction results)***

[ Submit ]  [ Clear All ]

**Reference:**
Yi-Heng Zhu, Zi Liu, Zhiwei Ji*, Dong-Jun Yu*. ULDNA: Integrating Unsupervised Multi-Source Language Models with LSTM-Attention Network for High-Accuracy Protein-DNA Binding Site Prediction. Briefings in Bioinformatics. 2024, 25(2):bbae040.

Contact @ Dong-Jun Yu
Programmed by Yi-Heng Zhu

**RESULTS PAGE**

Predicting Protein-DNA Binding Sites

**Protein Name**
2XTNA

**Model constructed on Dataset**
PDNA-543

**Threshold**
0.265 (*Max MCC*)

**Prediction Summary**
Number of predicted DNA-binding residues in protein **2XTNA: 2**
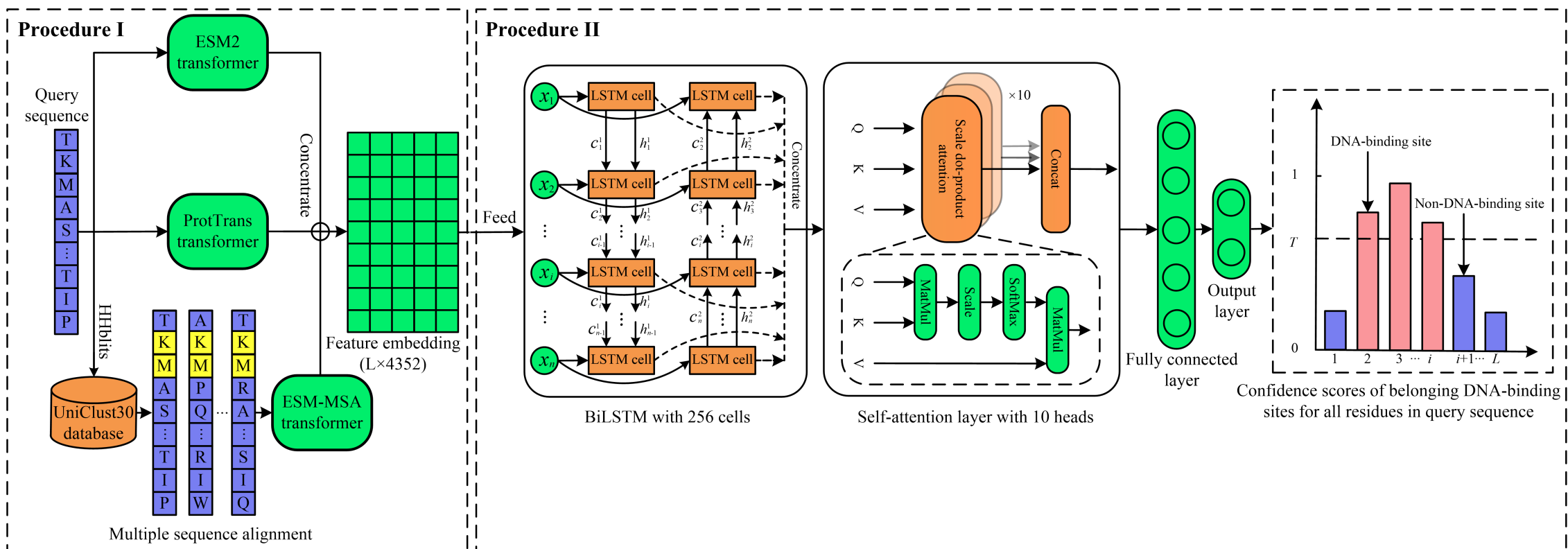Specific position: **58 T  117 R**

**Predicted Results**

| Residue # | Amino Acid Type | Probability | Binding Residue |
|---|---|---|---|
| 0001 | M | 0.046 | N |
| 0002 | D | 0.016 | N |
| 0003 | Q | 0.010 | N |
| 0004 | N | 0.013 | N |
| 0005 | E | 0.007 | N |
| 0006 | H | 0.079 | N |
| 0007 | S | 0.006 | N |
| 0008 | H | 0.067 | N |
| 0009 | W | 0.079 | N |
| 0010 | G | 0.005 | N |
| 0011 | P | 0.012 | N |
| 0012 | H | 0.116 | N |
| 0013 | A | 0.028 | N |
| 0014 | K | 0.090 | N |
| 0015 | G | 0.006 | N |
| 0016 | Q | 0.013 | N |
| 0017 | C | 0.010 | N |
| 0018 | A | 0.004 | N |
| 0019 | S | 0.006 | N |
| 0020 | R | 0.010 | N |

➢ 基于无监督语言模型与多视角多序列联配的蛋白质-蛋白质相互作用预测方法ICCPred

Zi Liu#, Yi-Heng Zhu#, Long-Chen Shen, Xuan Xiao, Wang-Ren Qiu, Dong-Jun Yu. Integrating Unsupervised Language Model with Multi-View Multiple Sequence Alignments for High-Accuracy Inter-Chain Contact Prediction. **Computers in Biology and Medicine**. 2023, 166: 107529

# 03　Part three

# 未来展望

蛋白质功能预测研究有助于推动智能医疗的发展

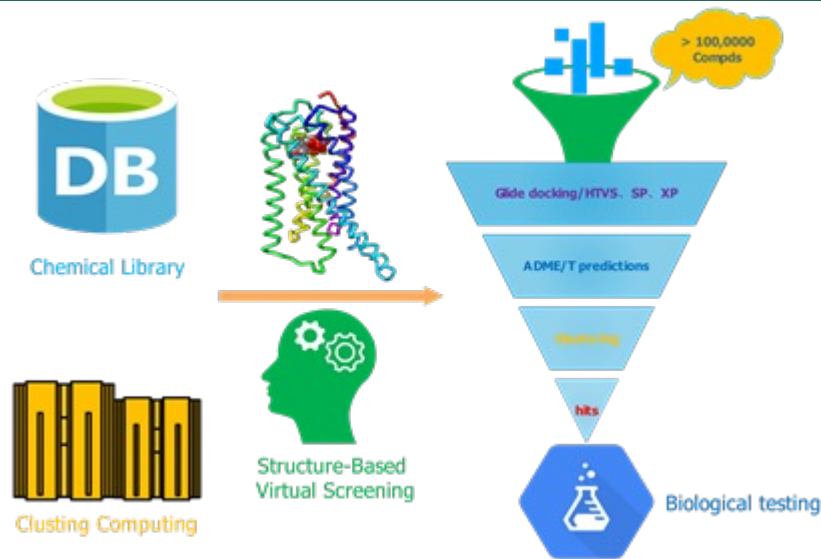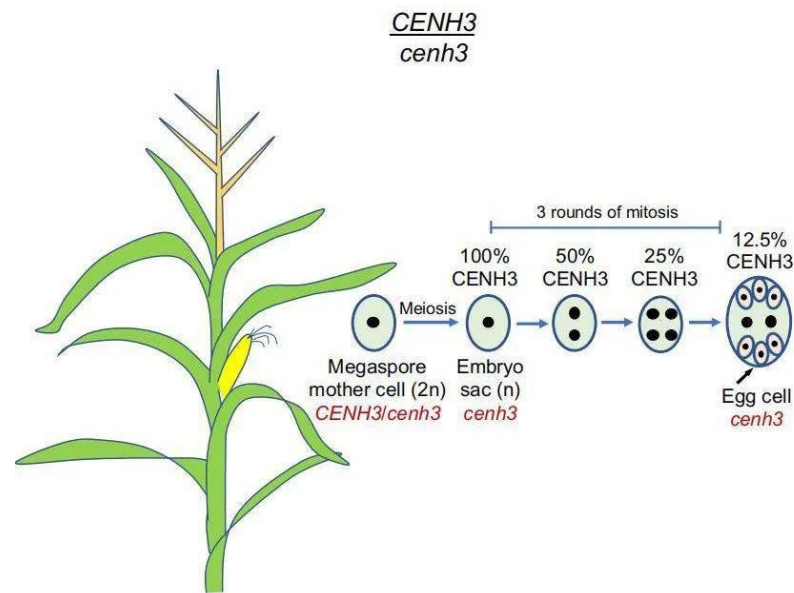(1) 辅助疾病分析和诊断 (推断关键致病蛋白质)

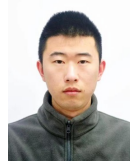(2) 辅助药物设计（药物分子筛选）

蛋白质功能预测在农业领域的应用前景

(1)  植物遗传育种

(2)  植物与微生物的相互作用

(3)  植物蛋白组学



Works as a vigorous hybrid

**研究方向**

**人工智能与模式识别**

人工智能的理论及应用
大数据计算与模式识别

**生物信息与系统生物学**

多组学数据整合分析与计算
复杂生物系统的数学建模与预测

团队组成：
教授1人
副教授1人
讲师3人
博士生3人
硕士生8人
已毕业研究生4人

◆ 时间序列数据挖掘与异常模式发现

◆ 高维复杂数据的维度约简和模型优化



time-series data

characterization methods
- distribution of values
- entropy
- correlation properties
- nonlinear time series analysis
- stationarity



**Z Ji***, Y Wang, X Xie, et al., *Expert Systems with Applications*, 2022.
N Jin, Y Zeng, K Yan, **Z Ji**, *IEEE Transactions on Industrial Informatics*, 2021.
M Hu, X Feng, **Z Ji***, et al., *Information Sciences*, 2019.
K Yan#, **Z Ji**#, et al., *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, 2019.
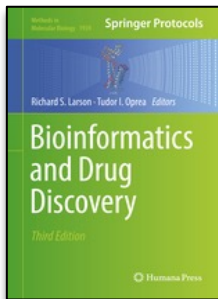K Yan, **Z Ji***, et al., *Neurocomputing*, 2017.
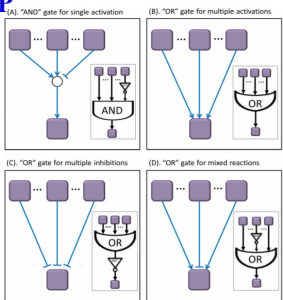
X Xie, F Xia, K Yan, H Xu, **Z Ji***, *Plant Phenomics*, 2023.
F Xia, X Xie , S Jin, K Yan, **Z Ji***, *Frontiers in Plant Science*, 2021.
K Yan, …, **Z Ji**., et al., *IEEE/ACM Trans on Computational Biology and Bioinformatics,* 2021.
X Xie, X Gu, Y Li, **Z Ji***, *Knowledge-based Systems*, 2021.
**Z Ji**, …, B Wang*, *Computational and mathematical methods in medicine*, 2015.

◆ 创立了一套独特的**生物分子网络**建模方法 | ◆ 创立了**分子-细胞-组织**的3D多尺度建模方法 | ◆ 建立了**生物组学大数据挖掘**的计算框架
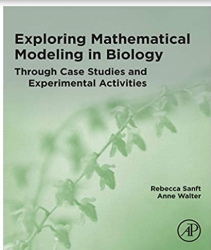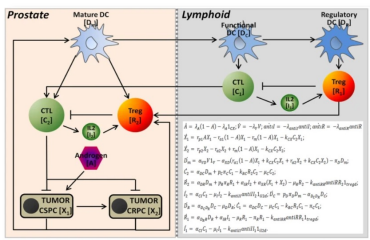
**基于线性规划的离散时间建模方法BLP,DILP, TILP, MIP**

BLP模型



✓ **BLP**作为**经典模型**被写入了Springer教材Methods in Molecular Biology丛书之一《**Bioinformatics and Drug Discovery**》（第16章第287页）
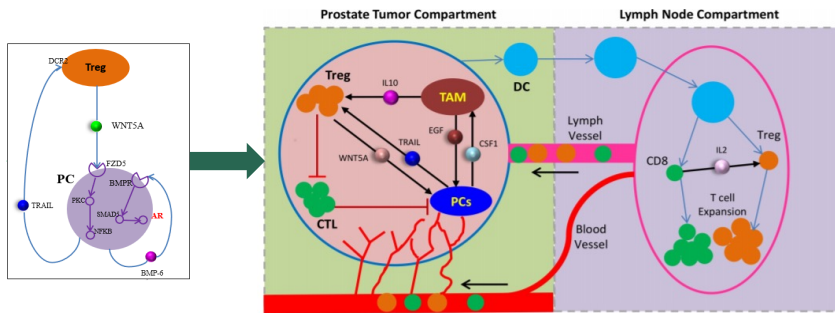
**基于微分方程的连续时间建模方法**

Cell-cell interaction模型
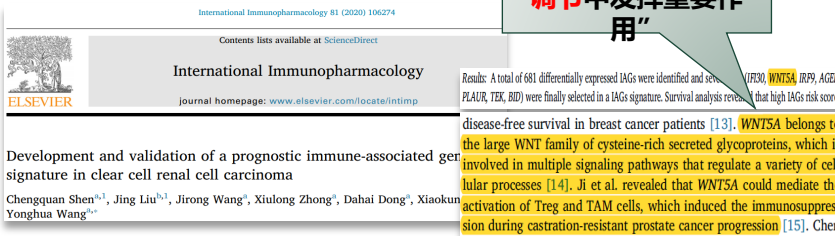
✓ 该模型被Matlab工具包**CRA**收录，并在BMC Bioinformatics进行长篇报道

✓ 作为**经典模型**被写入**Elsevier**教材(2020年)：《**Exploring Mathematical Modeling in Biology**》（第2章第54页）

**发现了WNT5A调控CRPC（前列腺癌）进展的新机制 构建了面向分子-细胞-组织的多尺度3D模型HABM**



"**WNT5A在免疫调节中发挥重要作用**"

International Immunopharmacology 81 (2020) 106274

Development and validation of a prognostic immune-associated gene signature in clear cell renal cell carcinoma
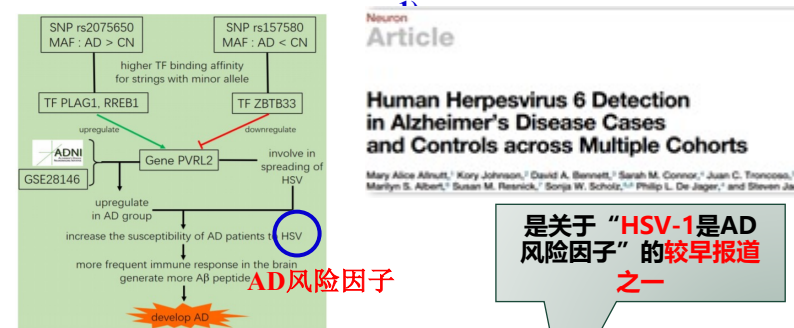
✓ **HABM是十几年来首个**对肿瘤生长-免疫反应-血管生成进行**3D时空建模**的**数学模型**

Digital Pathology Analysis Quantifies Spatial Heterogeneity of CD3, CD4, CD8, CD20, and FoxP3 Immune Markers in Triple-Negative Breast Cancer

**组学大数据挖掘发现AD潜在风险因子 (TREM2, HSV-1)**



Human Herpesvirus 6 Detection in Alzheimer's Disease Cases and Controls across Multiple Cohorts

**是关于"HSV-1是AD风险因子"的较早报道之一**

**AD风险因子**

or plasma (Lövheim et al., 2018). In addition, several groups have identified overlap between AD genetic risk factors and genes affected by viral infection, such as a receptor involved in spreading HSV-1 (Liu et al., 2018) and a human leukocyte antigen (HLA) subtype associated with increased susceptibility to HHV-6A infection (Rizzo et al., 2019).

✓ **时序组学大数据挖掘，首次解析了TREM2调控Microglia表型转换的分子机制**

谢谢各位老师和同学观看
请批评指正！