



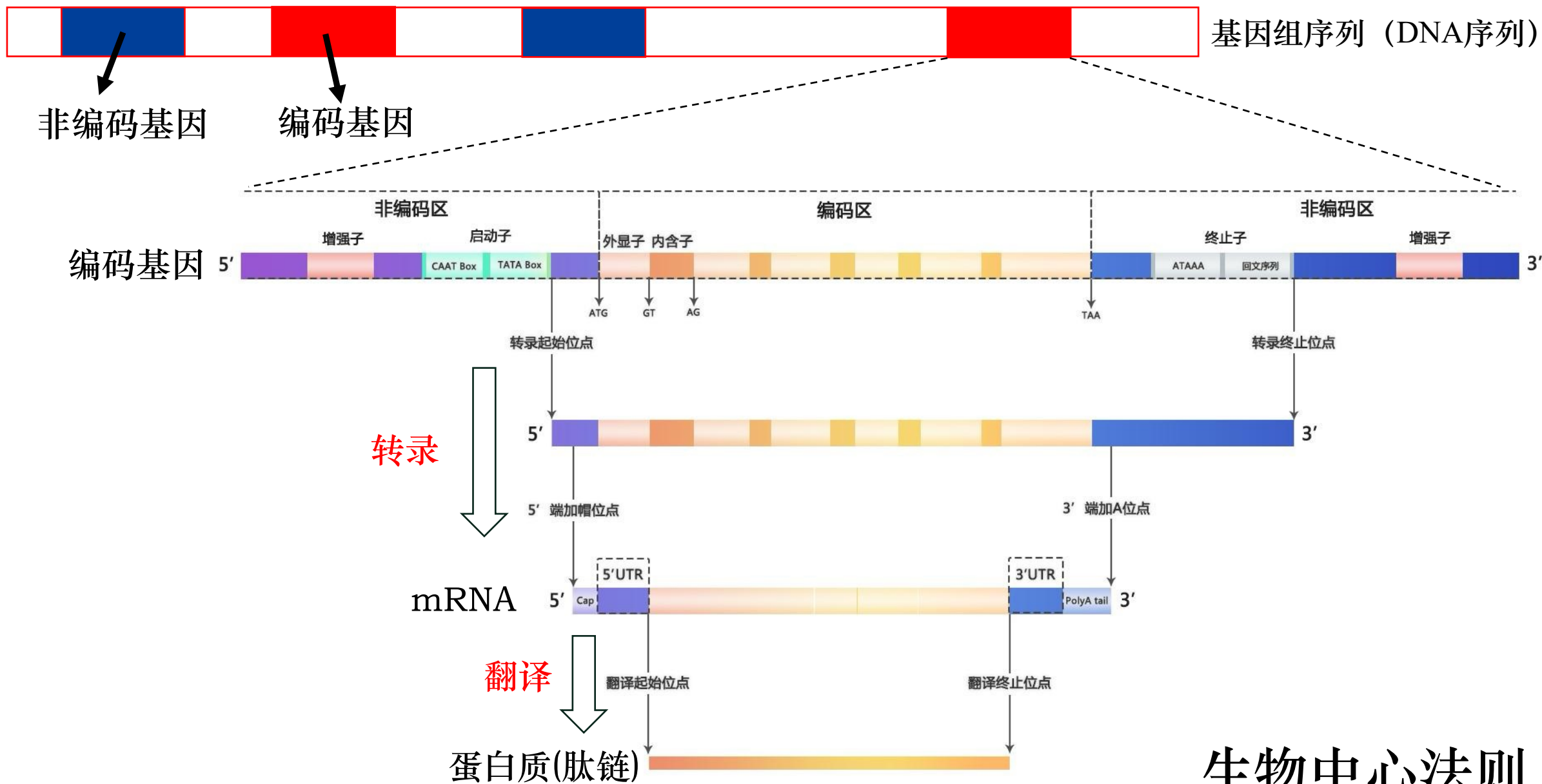
# 多模态数据融合的 蛋白质功能预测

南京农业大学 人工智能学院

汇报人：朱一亨

2024年06月20日

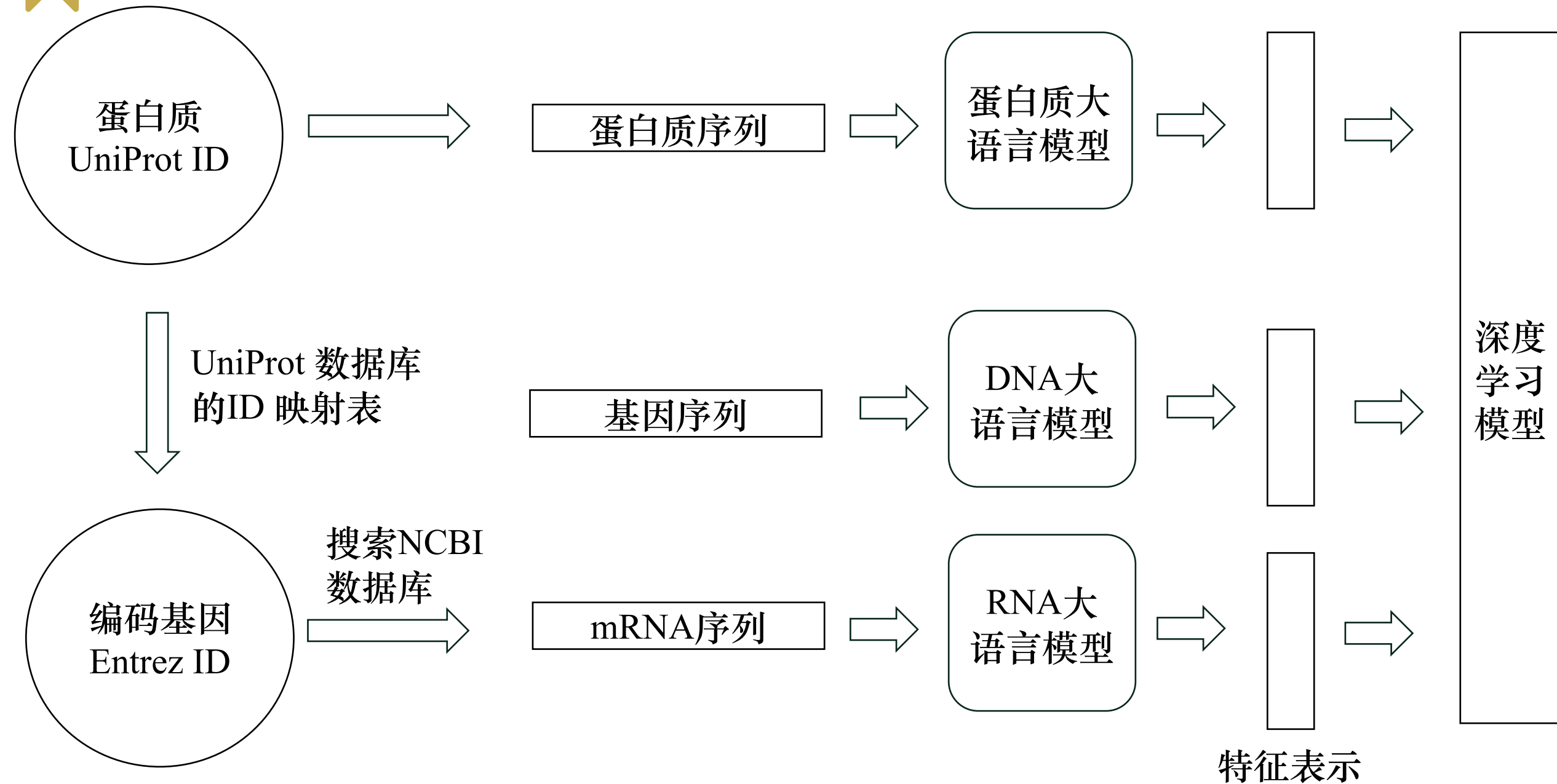
# 1 研究背景



## 2 现有方法的不足和挑战

- 现有的蛋白质功能预测方法只注重从蛋白质自身挖掘知识，忽略了编码基因和mRNA中的知识。
- 基因组数据： 基因序列（DNA序列） 和表达数据
- 转录组数据： mRNA序列

### 3 初步的方法和设想



## 4 数据集构建

- (1) 从蛋白质功能注释数据库GOA中，下载全部的蛋白质，共计134778个。
- (2) 部分蛋白质，在UniProt中找不到对应的编码基因记录，故丢弃。剩余100947个蛋白质。
- (3) Training dataset: 81657 proteins, before 2021-12-31
- (4) Validation dataset: 588 proteins, 2022-01-01 between 2022-06-30
- (5) Test dataset: 1723 proteins, 2022-07-01 between 2023-06-30
- (6) CD-HIT: (sequence identity<30%, training, validation, test datasets)

# 初步的实验结果

Dataset	Method	F <sub>max</sub>			AUPR		
		MF	BP	CC	MF	BP	CC
Validation Dataset (588 proteins)	DIOMAND	0.609	0.370	0.499	0.321	0.207	0.279
	BLAST	0.631	0.363	0.520	0.434	0.236	0.360
	Naive	0.432	0.256	0.449	0.244	0.155	0.331
	ESM2	<b>0.689</b>	0.411	0.589	<b>0.657</b>	0.292	0.563
	DNABERT2	0.461	0.319	0.520	0.357	0.219	0.468
	ESM2 + DNABERT2	0.681	<b>0.427</b>	<b>0.615</b>	0.651	<b>0.310</b>	<b>0.606</b>
Test Dataset (1723 proteins)	DIOMAND	0.632	0.340	0.474	0.386	0.193	0.249
	BLAST	0.630	0.331	0.490	0.445	0.207	0.323
	Naïve	0.367	0.218	0.389	0.196	0.112	0.257
	ESM2	0.665	0.396	0.552	0.627	0.266	0.521
	DNABERT2	0.444	0.356	0.483	0.339	0.266	0.427
	ESM2 + DNABERT2	<b>0.677</b>	<b>0.433</b>	<b>0.585</b>	<b>0.661</b>	<b>0.324</b>	<b>0.562</b>

## 目前存在的问题

➤ 部分基因序列过长，不好处理。

<10000: 61945

10000-20000: 8954

20000-30000: 5322

30000-40000: 3432

40000-50000: 2562

50000-60000: 1872

60000-70000: 1403

70000-80000: 1075

80000-90000: 961

>90000: 7112

Average: 28668

Go to [reference sequence details](#)

Primary Assembly

Go to nucleotide: [Graphics](#) [FASTA](#) [GenBank](#)

Tools | Tracks | Download |

7 M 7,500 K 8 M

**RBFOX1**

**Gene:** RBFOX1  
**Name:** RNA binding fox-1 homolog 1  
**Location:** 5,239,721..7,713,340  
**Length:** 2,473,620 nt  
[Positional Info]  
**NC\_000016.10 position:** 6,923,633  
**Gene position:** 1,683,913

**Links & Tools**  
**GeneID:** [54715 \(RBFOX1\)](#)  
**HGNC:** [18222](#)  
**MIM:** [605104](#)

**BLAST nr:** [NC\\_000016.10 \(5,239,721..7,713,340\)](#)  
**BLAST to Genome:** [NC\\_000016.10 \(5,239,721..7,713,340\)](#)  
**FASTA record:** [NC\\_000016.10 \(5,239,721..7,713,340\)](#)  
**GenBank record:** [NC\\_000016.10 \(5,239,721..7,713,340\)](#)

[GTR](#)  
[Nucleotide](#)  
[OMIM](#)  
[Probe](#)  
[Protein](#)  
[PubChem Comp](#)  
[PubChem Substa](#)  
[ed](#)  
[ed \(GeneR](#)  
[ed \(OMIM\)](#)  
[ed\(nucleot](#)  
[q Proteins](#)  
[q RNAs](#)  
[qGene](#)  
[d gene-spe](#)

## 下一步的计划

- 测试mRNA序列的性能
- 将DNABERT2替换成其他语言模型，如NC-Transformer。
- 在蛋白质通道加入AF2预测的结构，采用图注意力网络，使模型更复杂。



## 数据统计

- 目前的数据集共计83968 proteins (Dataset I)

Training dataset: 81657 proteins

Validation dataset: 588 proteins

Test dataset: 1723 proteins

- 只有62683proteins在NCBI中找到mRNA sequences (Dataset II)


Training dataset: 60994 proteins

Validation dataset: 453proteins

Test dataset: 1236 proteins

Dataset I	Method	F <sub>max</sub>			AUPR		
		MF	BP	CC	MF	BP	CC
Validation Dataset (588 proteins)	DIOMAND	0.609	0.370	0.499	0.321	0.207	0.279
	BLAST	0.631	0.363	0.520	0.434	0.236	0.360
	Naïve	0.432	0.256	0.449	0.244	0.155	0.331
	ESM2	<b>0.689</b>	0.411	0.589	<b>0.657</b>	0.292	0.563
	DNABERT2	0.461	0.319	0.520	0.357	0.219	0.468
	ESM2 + DNABERT2 (Gene Sequence)	0.681	<b>0.427</b>	<b>0.615</b>	0.651	<b>0.310</b>	<b>0.606</b>
	ESM2 + DNABERT2 (mRNA Sequence)	0.678	0.407	0.590	0.655	0.294	0.579
Test Dataset (1723 proteins)	DIOMAND	0.632	0.340	0.474	0.386	0.193	0.249
	BLAST	0.630	0.331	0.490	0.445	0.207	0.323
	Naïve	0.367	0.218	0.389	0.196	0.112	0.257
	ESM2	0.665	0.396	0.552	0.627	0.266	0.521
	DNABERT2	0.444	0.356	0.483	0.339	0.266	0.427
	ESM2 + DNABERT2	<b>0.677</b>	<b>0.433</b>	<b>0.585</b>	<b>0.661</b>	<b>0.324</b>	<b>0.562</b>
	ESM2 + DNABERT2 (mRNA Sequence)	0.671	0.420	0.579	0.622	0.300	0.544

Dataset II	Method	$F_{\max}$			AUPR		
		MF	BP	CC	MF	BP	CC
Validation Dataset (453 proteins)	ESM2	<b>0.698</b>	0.390	0.596	0.667	0.269	0.580
	DNABERT2 (Gene Sequence)	0.504	0.317	0.543	0.387	0.218	0.523
	DNABERT2 (mRNA Sequence)	0.531	0.297	0.518	0.411	0.198	0.492
	ESM2 + DNABERT2 (Gene Sequence)	0.691	0.405	<b>0.618</b>	<b>0.670</b>	<b>0.298</b>	<b>0.601</b>
	ESM2 + DNABERT2 (mRNA Sequence)	0.687	0.393	0.594	0.661	0.272	0.586
Test Dataset (1236 proteins)	ESM2	0.659	0.392	0.550	0.621	0.258	0.525
	DNABERT2 (Gene Sequence)	0.464	0.364	0.489	0.341	0.271	0.450
	DNABERT2 (mRNA Sequence)	0.470	0.330	0.484	0.347	0.225	0.439
	ESM2 + DNABERT2 (Gene Sequence)	<b>0.677</b>	<b>0.434</b>	<b>0.588</b>	<b>0.637</b>	<b>0.326</b>	<b>0.564</b>
	ESM2 + DNABERT2 (mRNA Sequence)	0.667	0.414	0.568	0.620	0.293	0.545



➤ mRNA不起作用的原因：mRNA中只有外显子序列，它们直接通过密码子表翻译成氨基酸序列，因此mRNA序列中的大部分知识和蛋白质序列中的知识重复。

>NC\_004327.3:c489287-488579 Plasmodium falciparum 3D7 genome assembly, chromosome: 6

ATGGTAATAAAAAGAAGAAAAAGAAAAAACAACACTTATAAATATATAAATAATTTATATATATA  
AATCAGCGTTATTATAAAATGAAAATAAATAAATATATAAATATATATAAATATATATAATGAATAGA  
TTAATAAAAATAAGAAAATCAAAAGAAATATATACCCCTTTAATATTTTAGTTAACATATATATATATAT  
ATATATATATATATATATATATGAATATACATTAGTTAAAATATTTTCCTTGATTATGTTTATATTTATG  
AATTTTTTTTTTTTTTTTTTTGAAGGGATCAATTAAACGTTTTAGATTAAAACAAAGACTTGGAATGCAG  
AAGGCAAATAGGCCTGTACCCATTGGTATAGATTAAAGAAAGATACAAAAATAAGGTAAGGCTTAAAC  
AAAAATATATTTATATATATATGTATGTCAATGAATATATGCTATGCTATGAATAAAAAGAAAAAATTATAA  
AAAAACGAGTACTTTTTATTTTATGCAAAATAACTTAATGTTGTTATAATATATAACATGGATTTGTTTCG  
TATTTGTATATAATATTATATTATATTATATTATATTATATTATATTATATTATAATATGTTTTT  
TTTTTTTTTTTAATTTTTTTTTTTTATTAGATATAATACAAAAGAAGACACTGGAGAAGAACCAAATTA  
GGATTATAA

← Gene sequence

>XM\_002808677.1 Plasmodium falciparum 3D7 60S ribosomal protein L39 (PF3D7\_0611700), partial mRNA

ATGGGATCAATTAAACGTTTTAGATTAAAACAAAGACTTGGAATGCAGAAGGCAAATAGGCCTGTAC  
CCCATTTGGTATAGATTAAAGAAAGATACAAAAATAAGATATAATACAAAAGAAGACACTGGAGAAGAAC  
CAAATTAGGATTATAA

mRNA sequence (cDNA)

>C0H4H3  
MGSIKRFRLKQRLGKCRRQNRPVPHWYRLKKDTKIRYNTKRRHWRRTKLGL

← Protein sequence

		密码子的第二位					
		T	C	A	G		
密码子的第一位 (5'端)	T	TTT: Phe F TTC: Phe F TTC: LeI L TTC: LeI L	TCT: Ser S TCC: Ser S TCA: Ser S TCG: Ser S	TAT: Tyr Y TAC: Tyr Y TAA: Ter * TAG: Ter *	TGT: Cys C TGC: Cys C TGA: Ter * TGG: Trp W	T C A G	密码子的第三位 (3'端)
	C	CTT: LeI L CTC: LeI L CTA: LeI L CTG: LeI L	CCT: Pro P CCC: Pro P CCA: Pro P CCG: Pro P	CAT: His H CAC: His H CAA: Gln Q CAG: Gln Q	CGT: Arg R CGA: Arg R CGC: Arg R CGG: Arg R	T C A G	
	A	ATT: Ile I ATC: Ile I ATA: Ile I ATG: Met M	ACT: Thr T ACC: Thr T ACA: Thr T ACG: Thr T	AAT: Asn N AAC: Asn N AAA: Lys K AAG: Lys K	AGT: Ser S AGC: Ser S AGA: Arg R AGG: Arg R	T C A G	
	G	GTT: Val V GTC: Val V GTA: Val V GTG: Val V	GCT: Ala A GCC: Ala A GCA: Ala A GCG: Ala A	GAT: Asp D GAC: Asp D GAA: GlT E GAG: GlT E	GGT: Gly G GGC: Gly G GGA: Gly G GGG: Gly G	T C A G	