

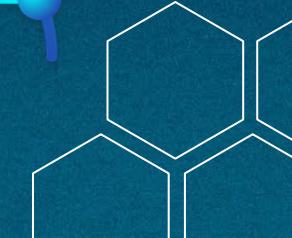
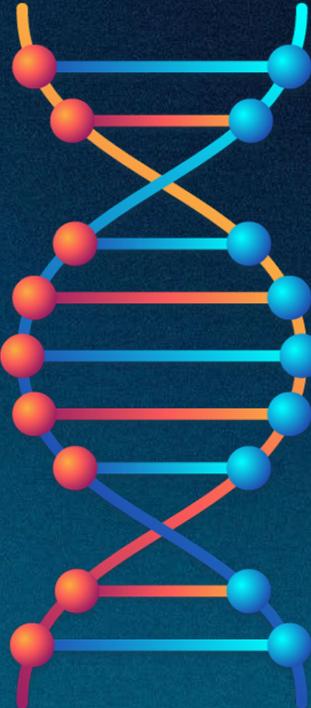


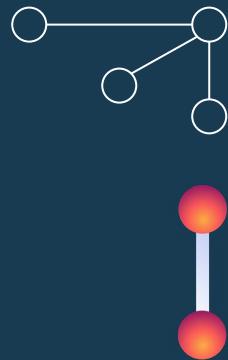
# DNA/RNA/Single -Cell Language Models

---

Student: Omnia Sarhan

Instructor: Yu Zhang





# Table of contents

01. \_\_\_\_\_

## Paper 1

DNABERT Model

02. \_\_\_\_\_

## Paper 2

5' UTR Model

03. \_\_\_\_\_

## Paper 3

scGPT Model

04. \_\_\_\_\_

## Summary

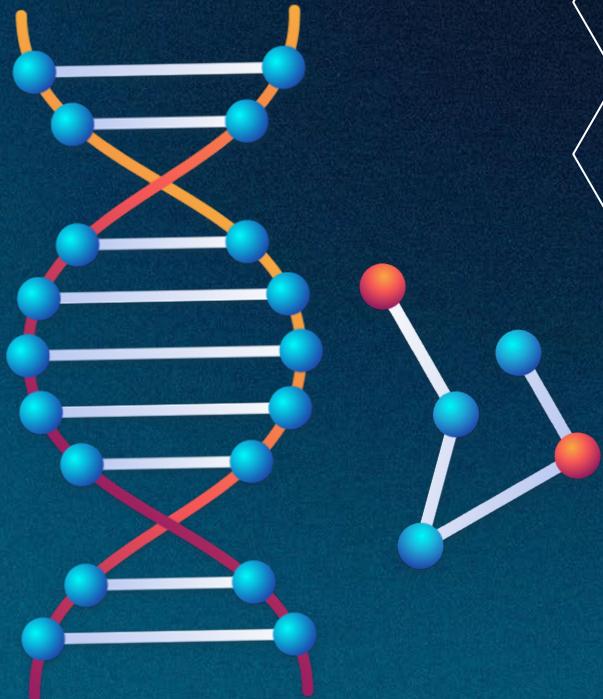
Conclusion & Questions  
session



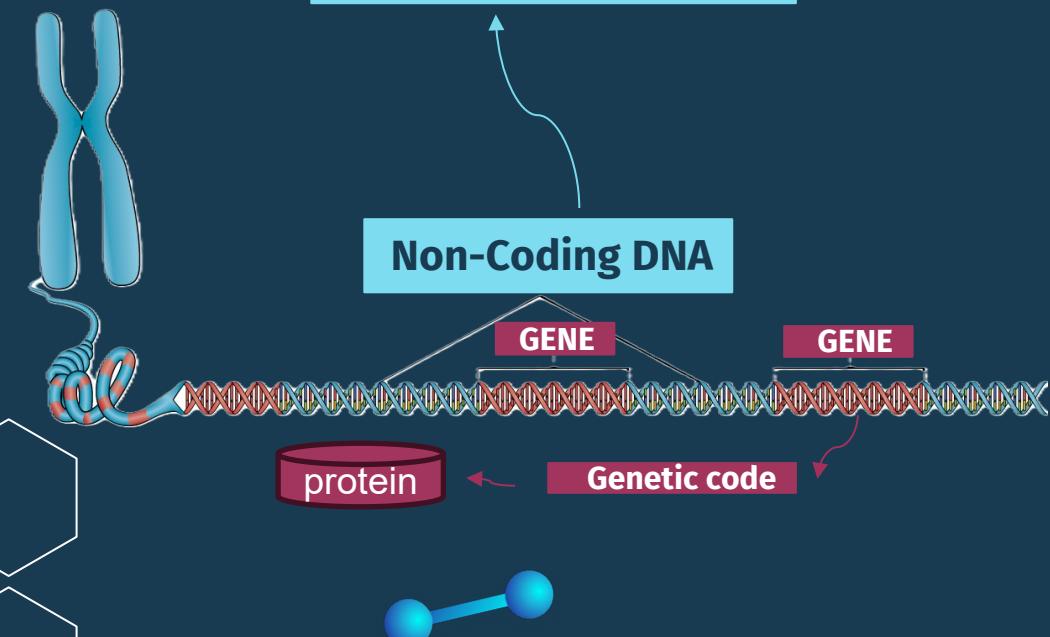
# 01.

# DNABERT

Pre-trained Bidirectional Encoder  
Representations from Transformers Model for  
DNA-Language in Genome



# Introduction



## Problem Statement

- Deciphering Non-Coding DNA for hidden instructions is challenging.
- Traditional models fail to capture long-range dependencies and polysemous relationships within DNA sequences.

# Objectives

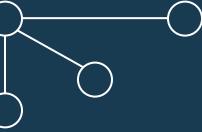
**Capture** global and transferable contextual information from DNA sequences

**Outperform** traditional deep learning models in various genomic tasks

**Provide** visualization mechanisms for interpretation of sequence motifs

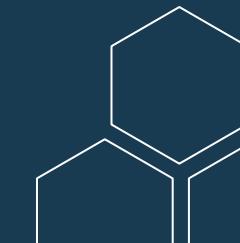
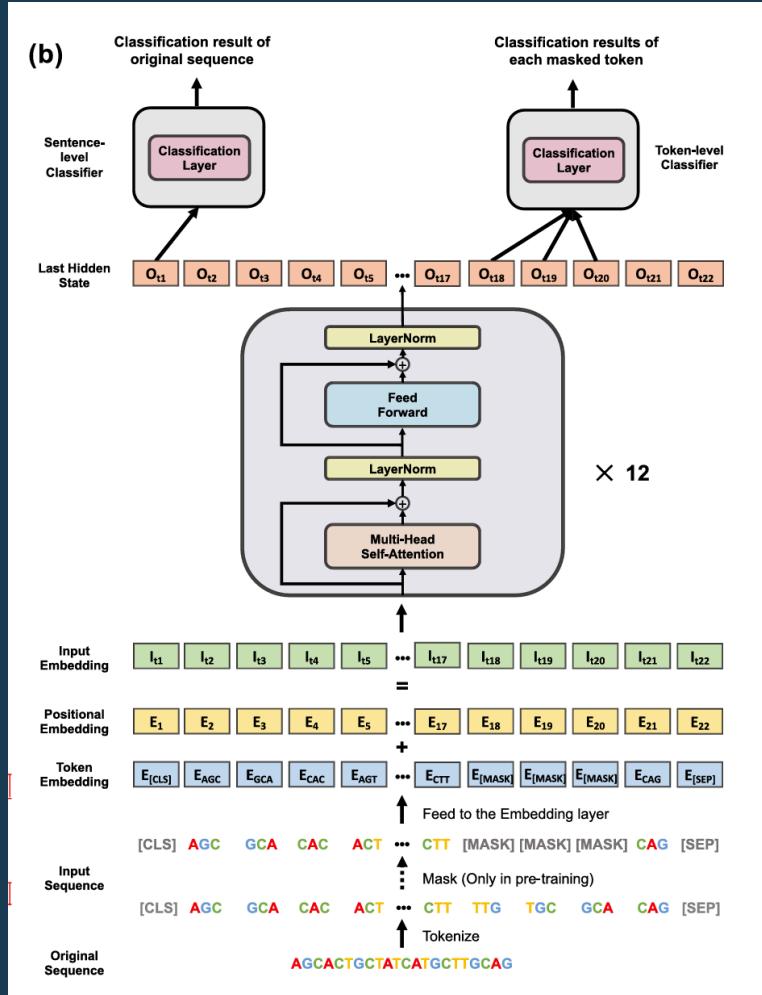
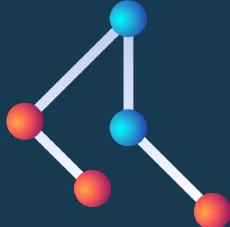
**Demonstrate** cross-organism applicability

**Facilitate** fine-tuning on task-specific datasets



# DNABERT Model

- BERT-based (same architecture)
- Adopts pre-training + fine-tuning





# Methodology



## Tokenization

- k-mer representation instead of single nucleotides.
- Different values of k (3, 4, 5, 6)
- Added special tokens like [CLS], [PAD], [UNK], [SEP], and [MASK]



## Pre-training

- masked language modeling (MLM) for random masking
- human genome (5-510 base pairs)
- 12 Transformer layers, 768 hidden units, and 12 attention heads



## Fine-tuning

- task-specific datasets
- Long sequences exceeding 512 tokens are split and processed as DNABERT-XL.
- Best = DNABERT-6

# Results

## DNABERT-Splice

Accurately recognizes canonical and non-canonical splice sites

## DNABERT-TF

Accurately identifies transcription factor binding sites

## DNABERT

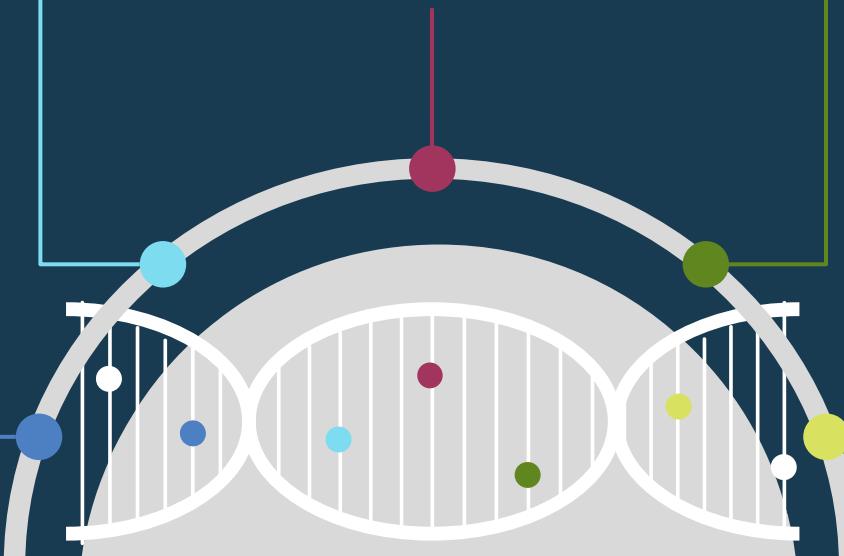
- Generalize over tasks
- Identifying functional genetic variants

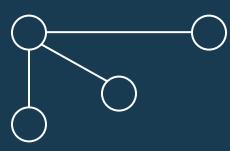
## DNABERT-viz

Allows visualization of important regions, contexts and sequence motifs

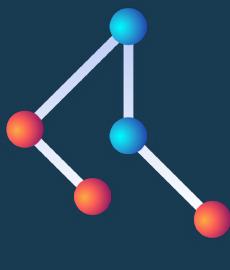
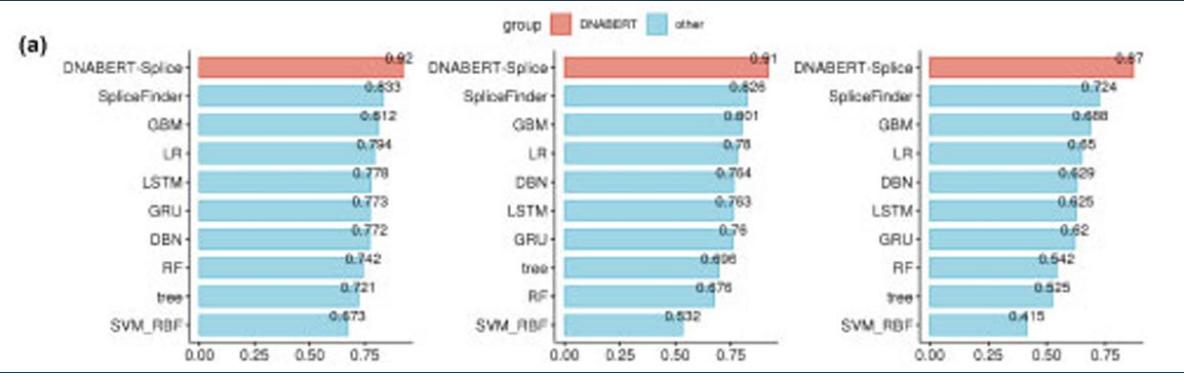
## DNABERT-Prom

Effectively predicts proximal and core promoter regions

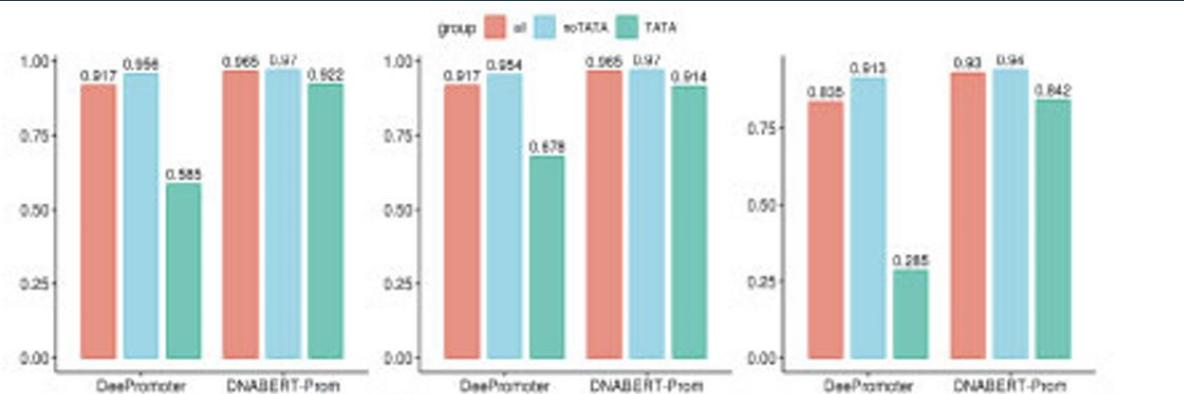




# Splice

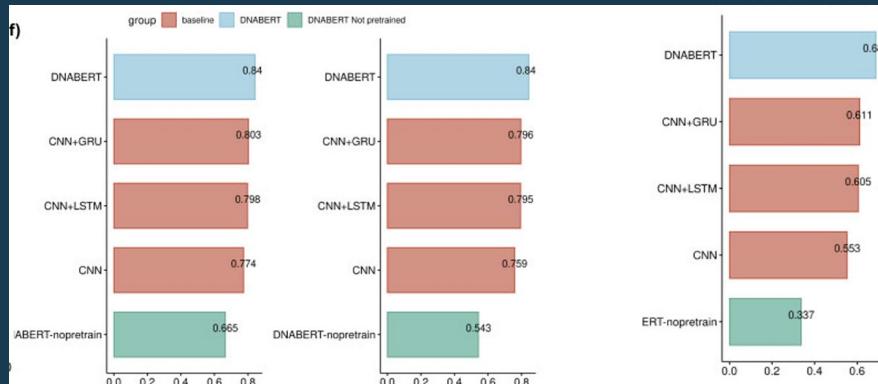
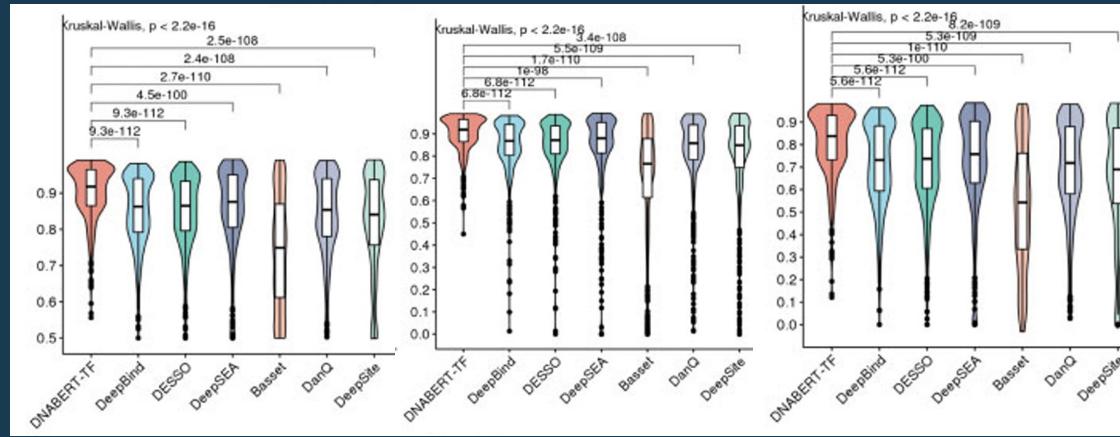


# Prom



# Results: (left to right) accuracy, F1 and MCC

TF

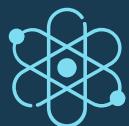


General

: : : :

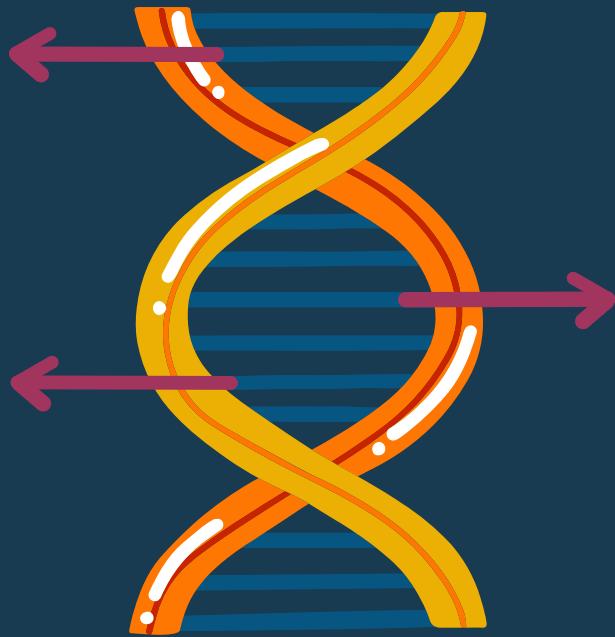
# Future Work

## 1. Other sequence prediction tasks



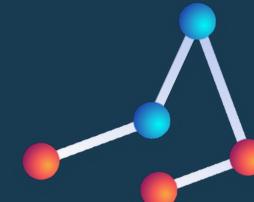
Determining CREs and enhancer regions from ATAC-seq and DAP-seq

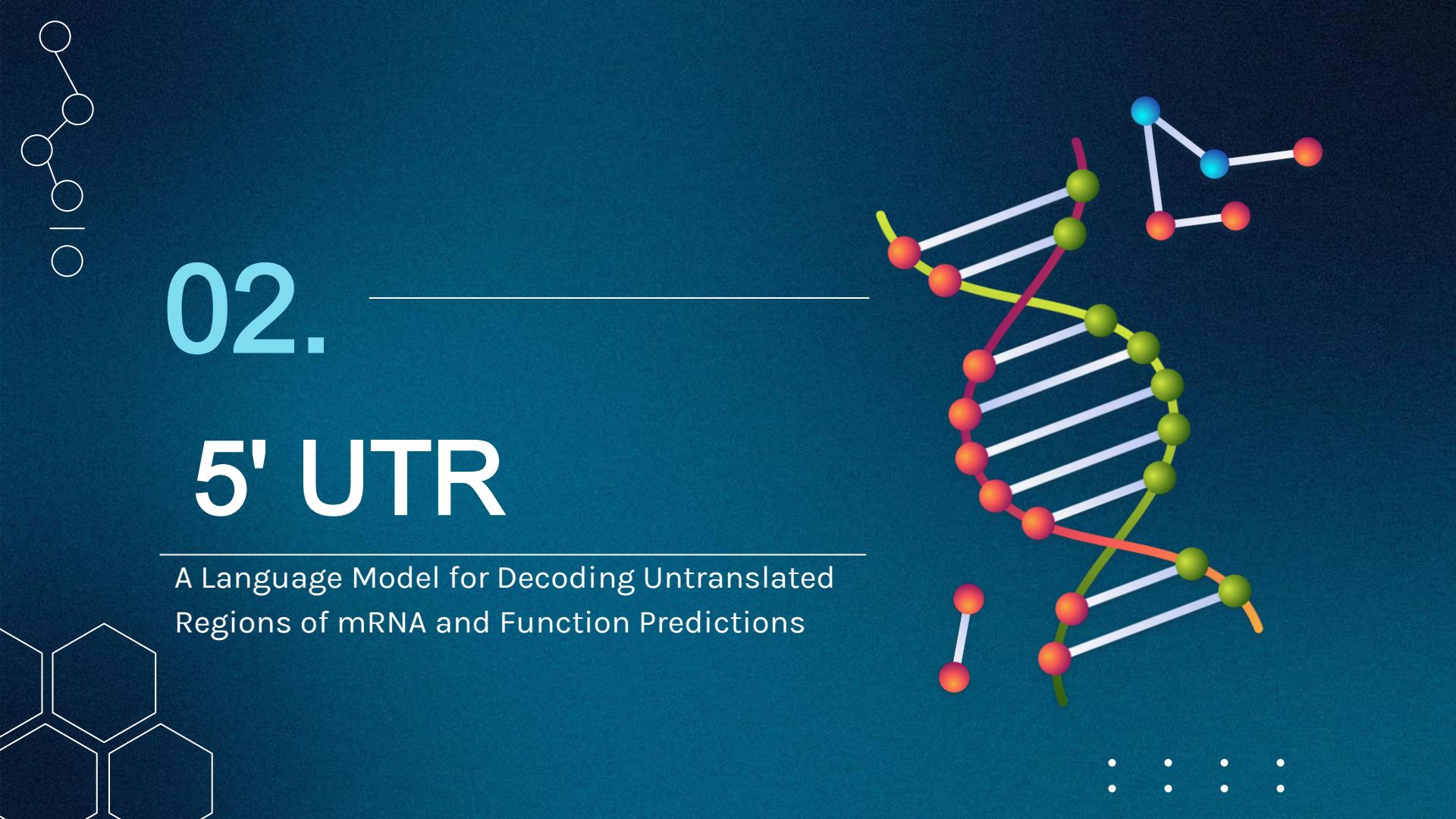
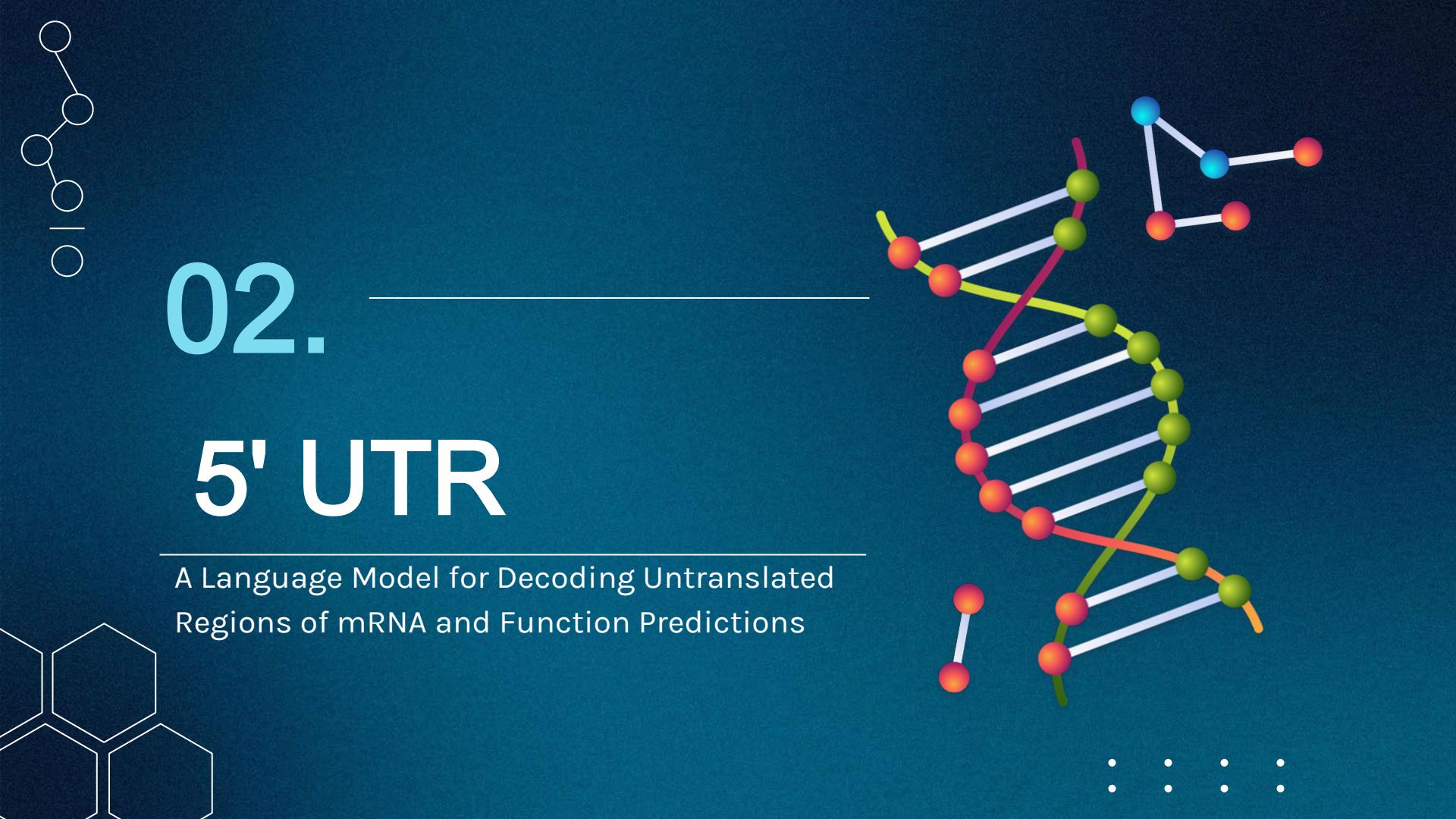
## 3. Direct machine translation on DNA



## 2. Prediction of binding preferences of RNA-binding proteins (RBPs)

s





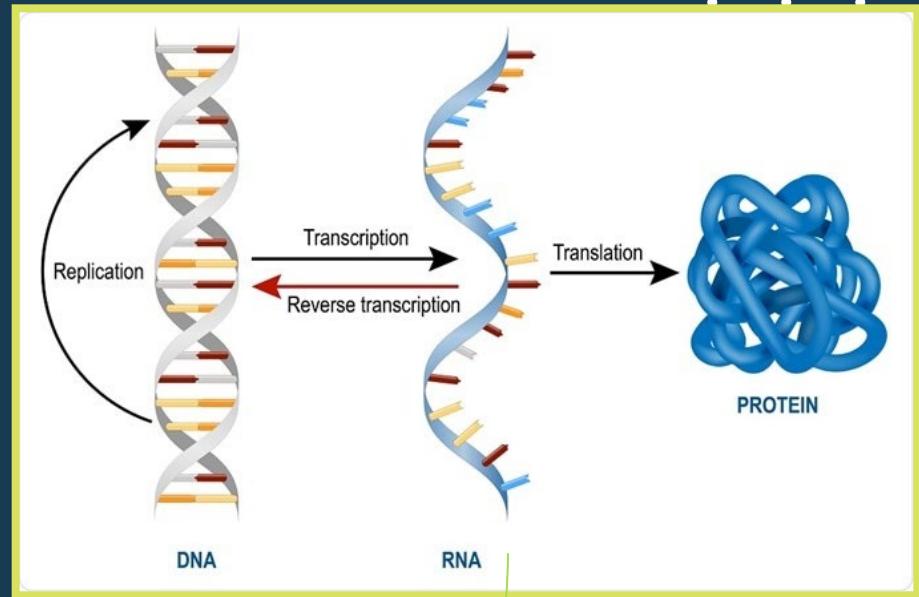
# 02. 5' UTR

---

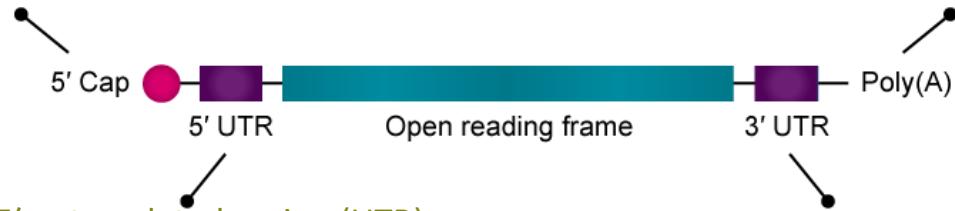
A Language Model for Decoding Untranslated  
Regions of mRNA and Function Predictions



# Background



5' Cap - plays a critical role in translational yield and nucleic acid stability *in vivo*



5' untranslated region (UTR)

5' UTR - regulates protein expression levels and translation initiation

Poly(A) tail - protects the mRNA from nuclease degradation

3' UTR - regulates protein expression levels and influences the stability and half-life of the mRNA

: : : :

# Introduction

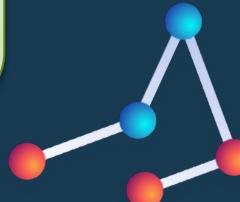


# Objectives

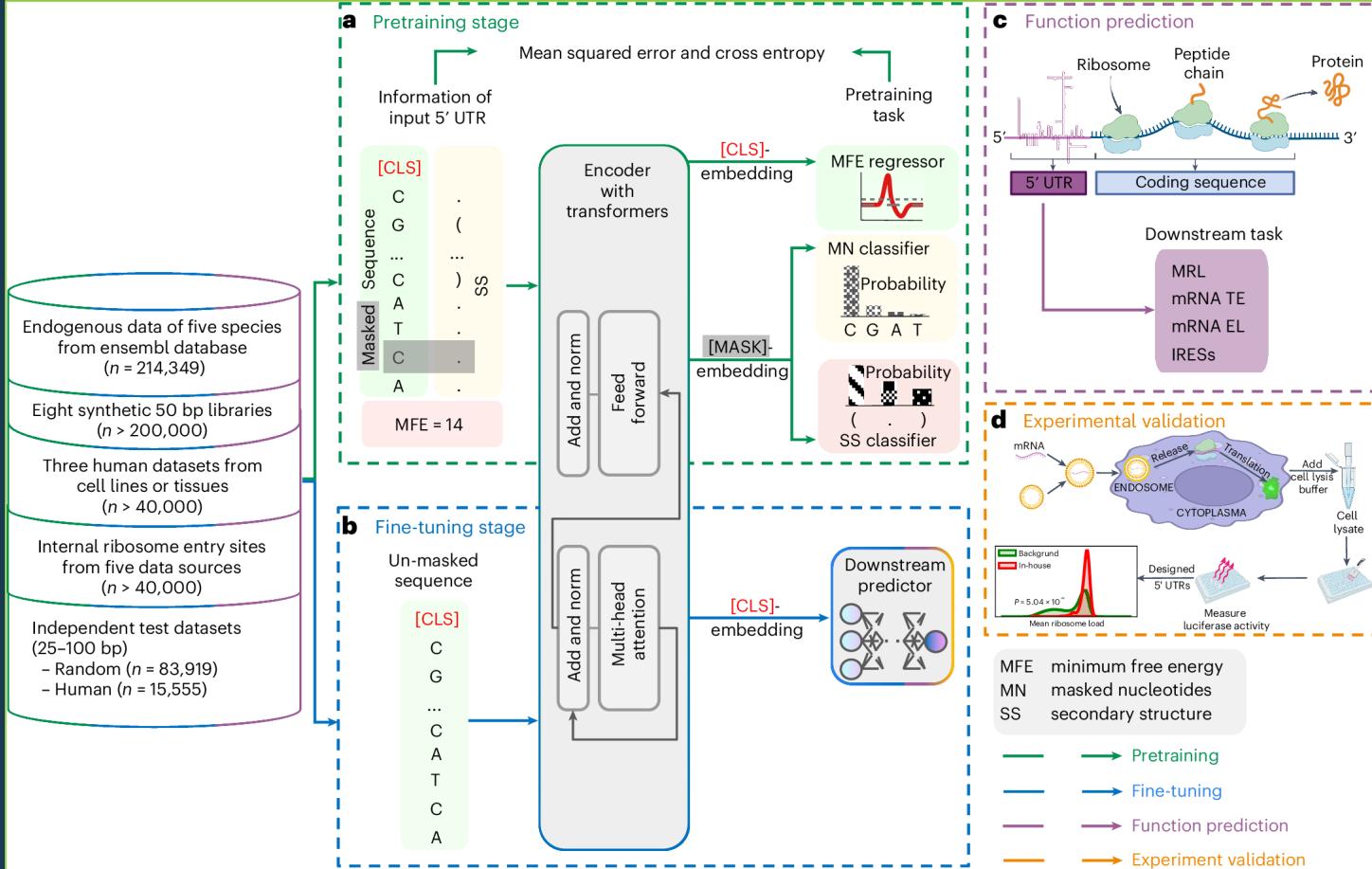
## Problem Statement

No unified foundation model to study function of 5'UTR

Use Language model to  
Extract meaningful  
semantic representations  
from UTRs of raw mRNA  
sequences and map them  
to predict functions of  
interest.



# 5'UTR-LM Model Overview



⋮ ⋮ ⋮ ⋮



# Results



UTR-LM predicts the  
ribosome loading



URR-LM identifies IRESSs

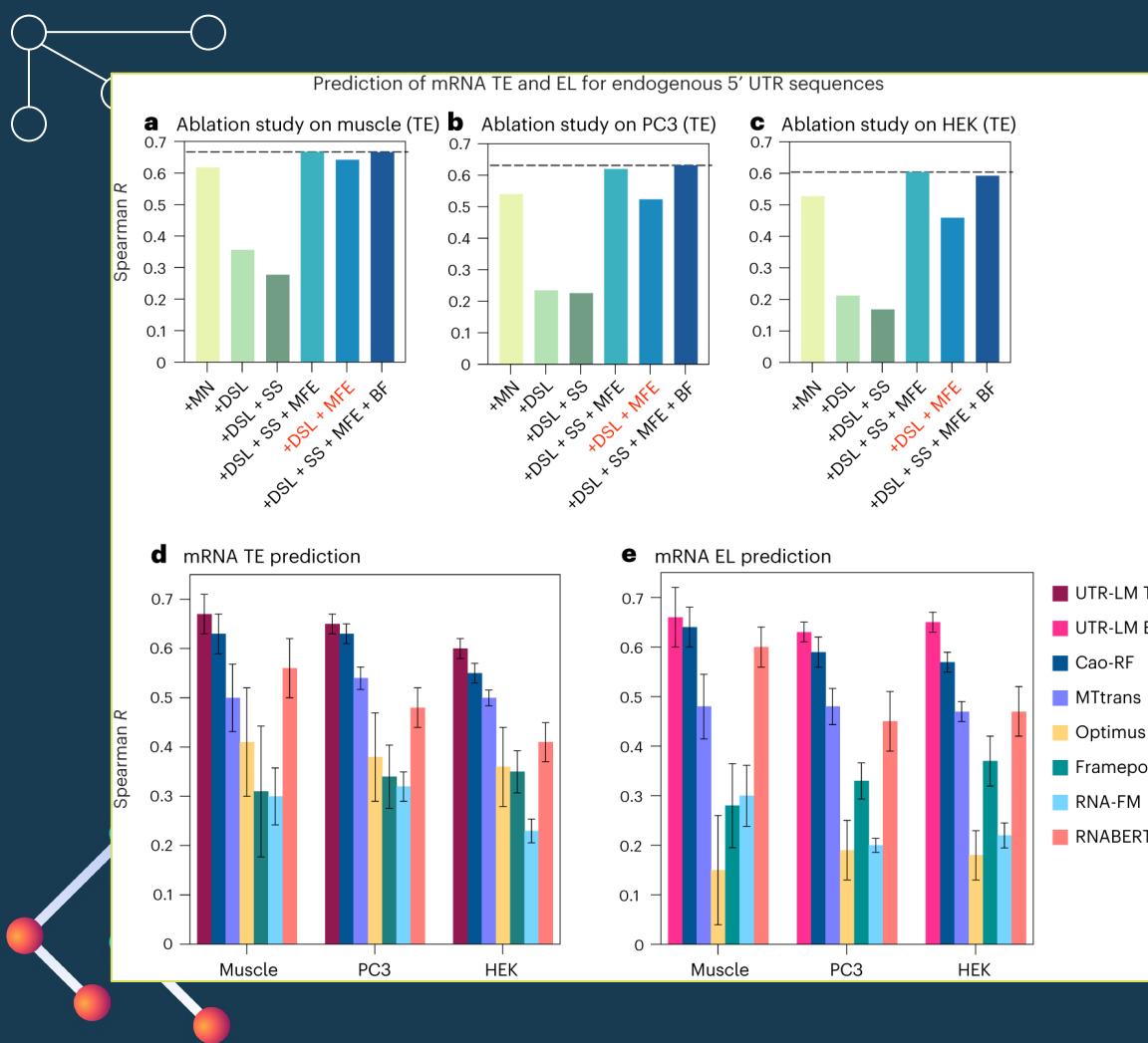


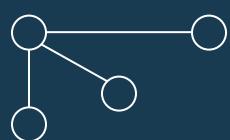
UTR-LM predicts mRNA  
TE and expression



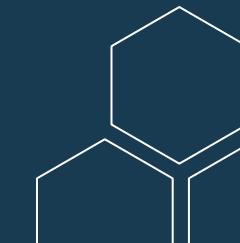
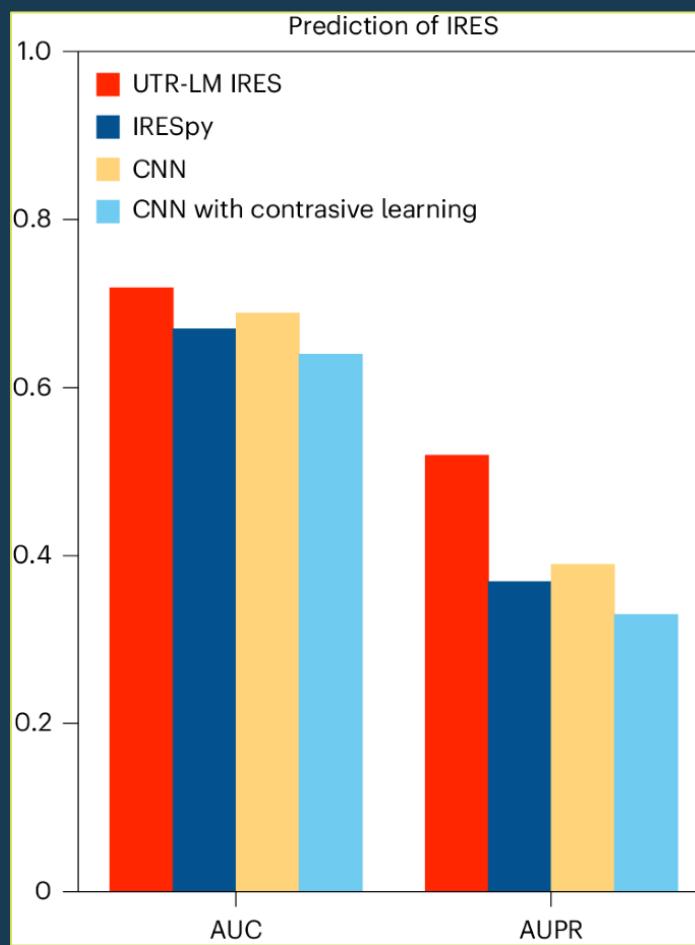
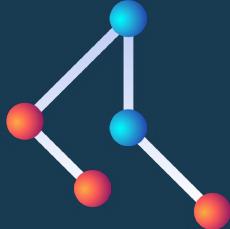
New designs validated in  
wet-lab experiments

# UTR-LM predicts mRNA TE and EL





# URR-LM identifies IRESs



: : : :

# Conclusion



## Conclusion

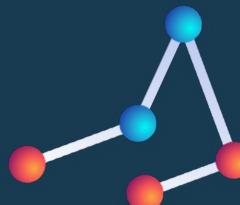
- Outperforms the best-known baseline in each task.
- Performance not limited by sequence length

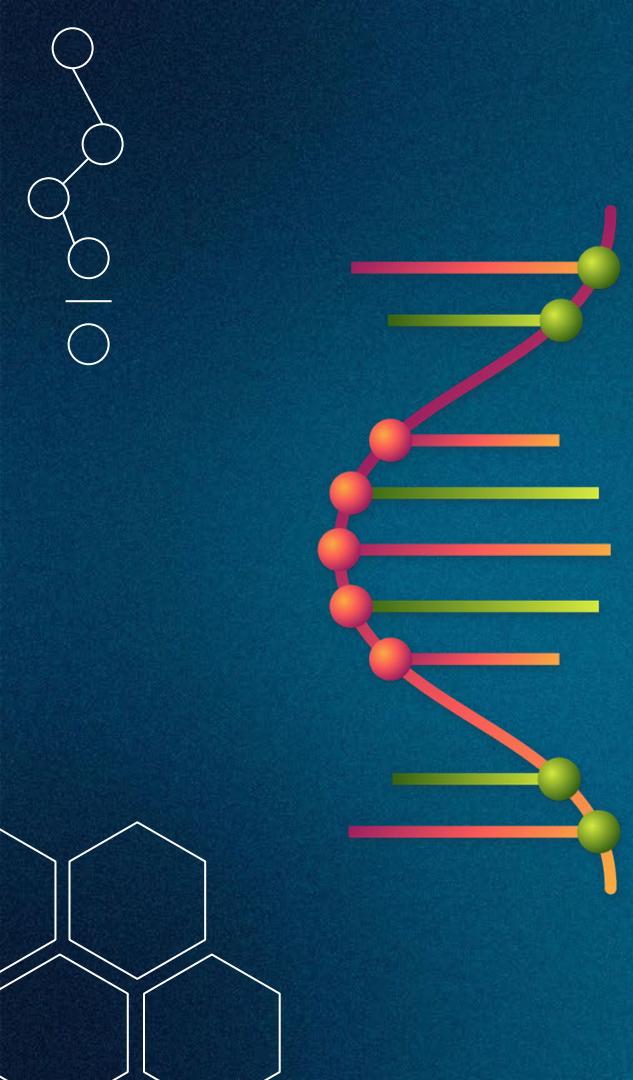
## limitations

- Computationally expensive

## Future

sparse transformers for modelling longer RNA sequences and more complex biological functions





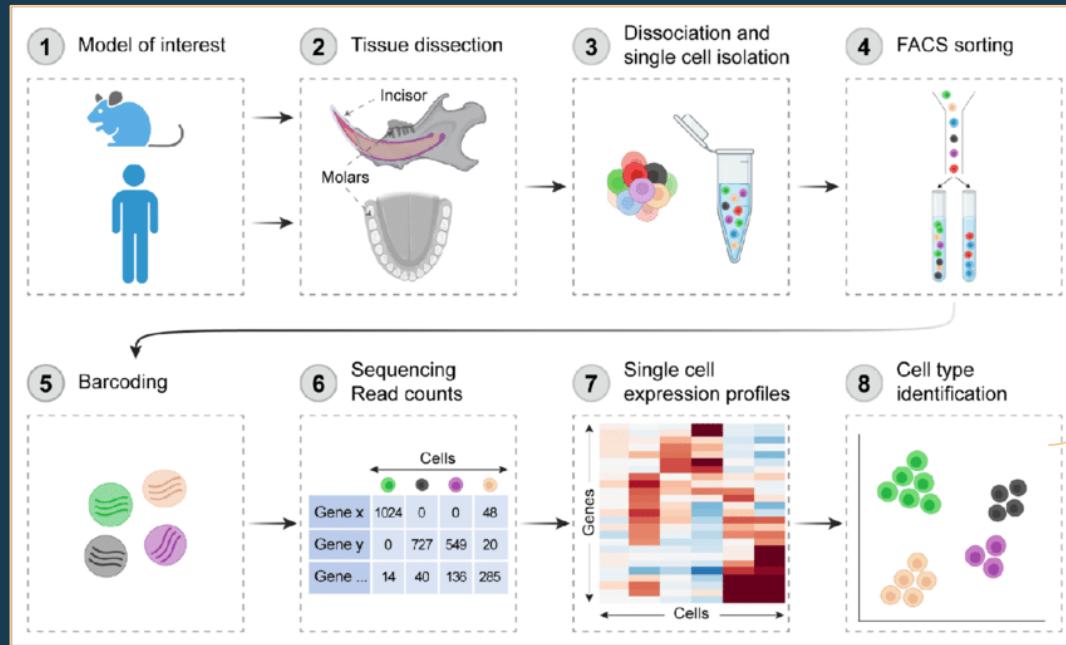
# 03. scGPT

---

Towards Building a Foundation Model for Single-  
Cell Multi-omics using Generative AI



# Single -cell RNA sequencing ( scRNA-seq)



personalized therapeutic strategies

cellular heterogeneity exploration

lineage tracking

pathogenic mechanism elucidation



# Introduction

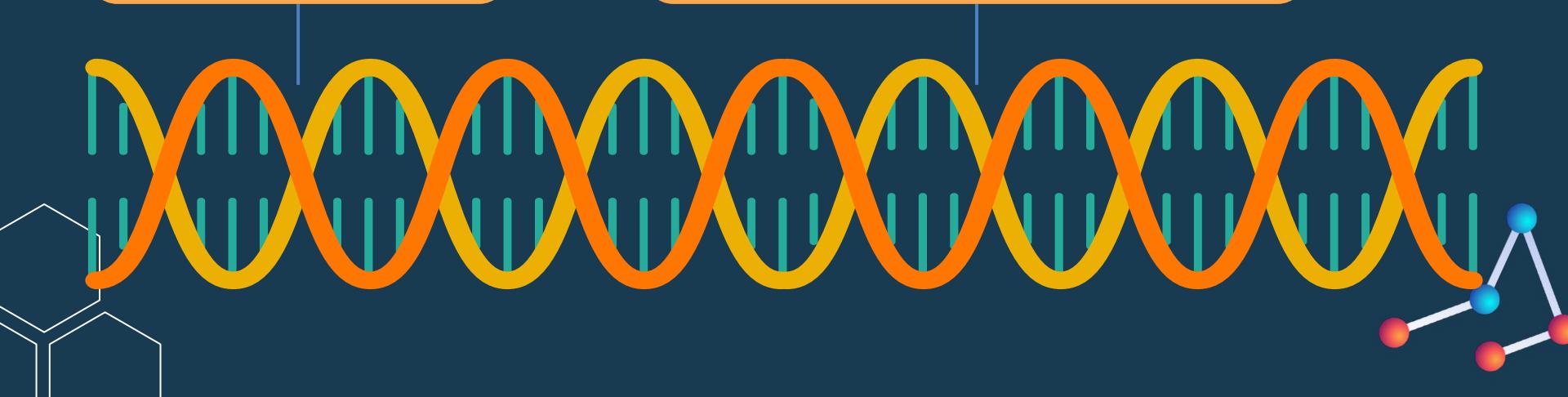


## Problem Statement

Current machine-learning-based methods in single-cell research are scattered

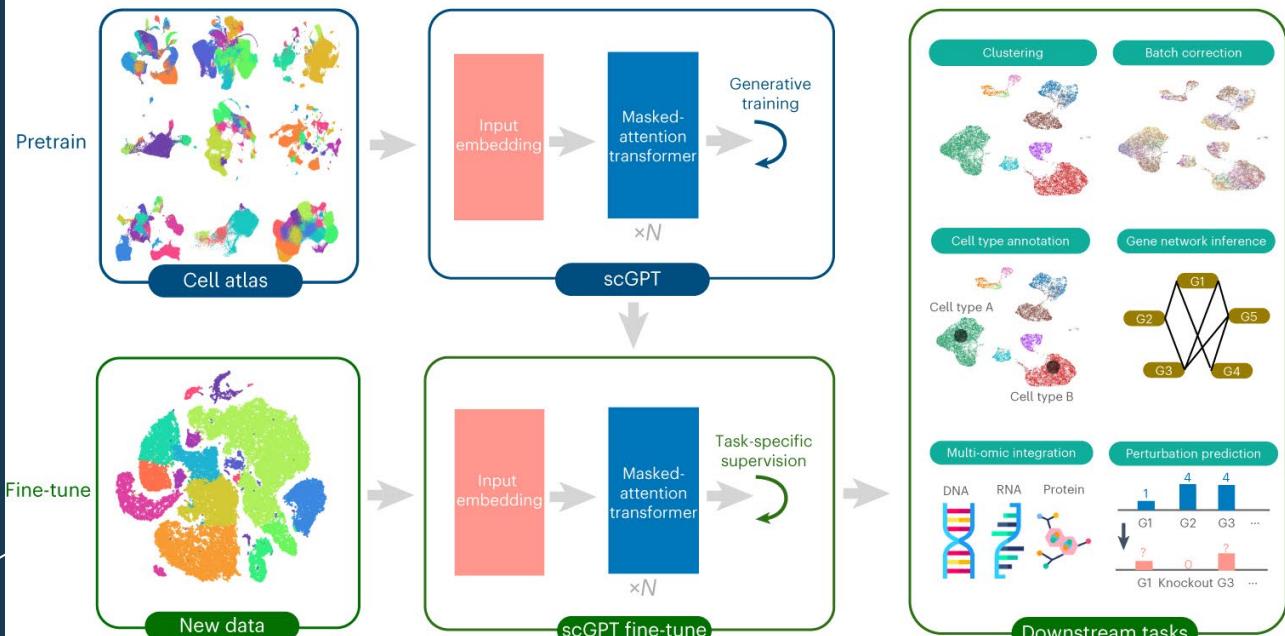
## Objectives

- Foundation model pretrained on large-scale data
- comprehend the complex interactions between genes across diverse tissues.

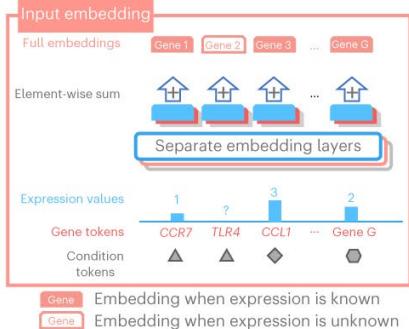


# scGPT Model Overview

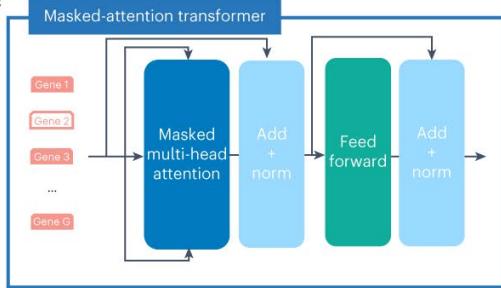
a



b



c



: : : :

# Results



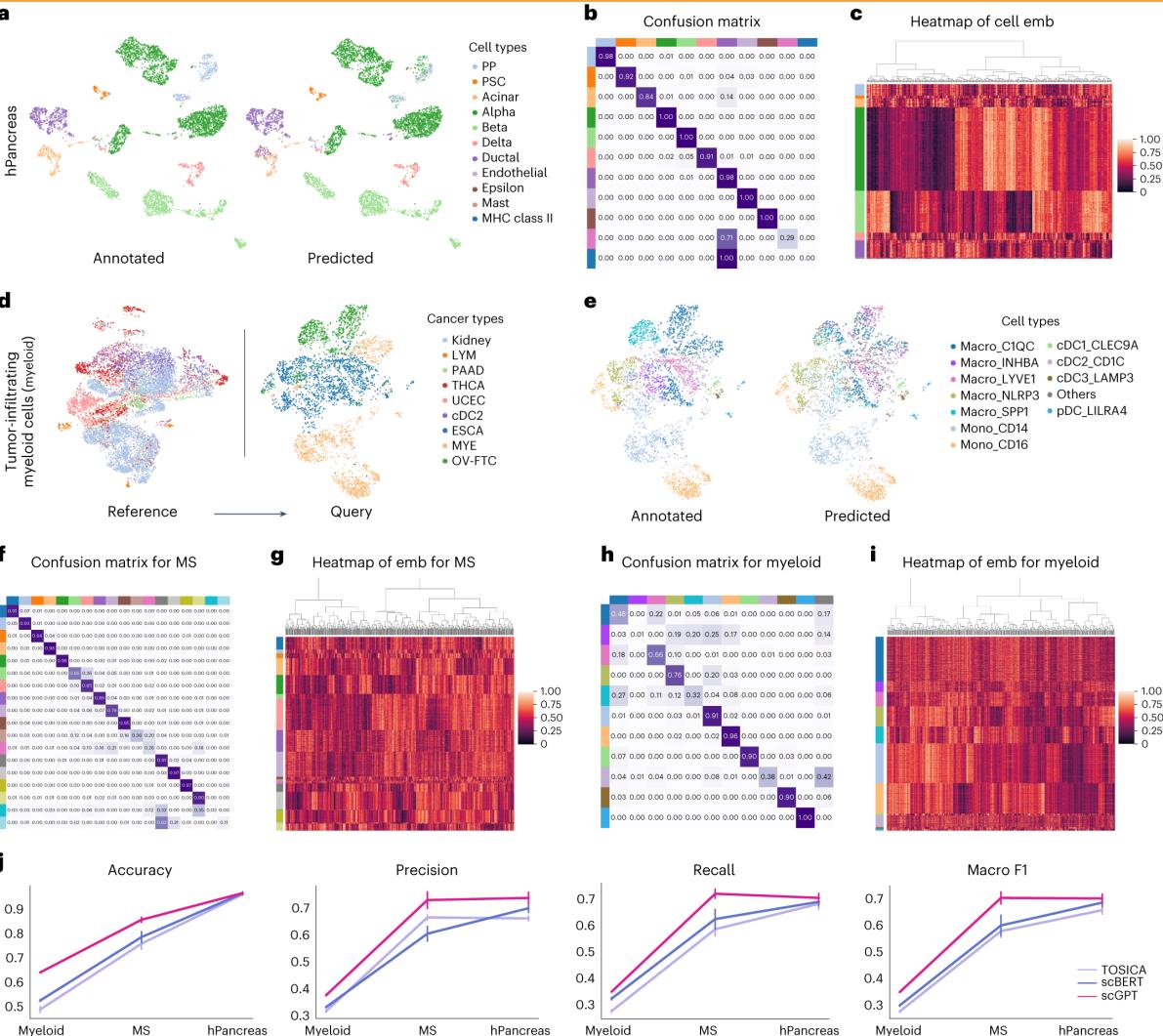
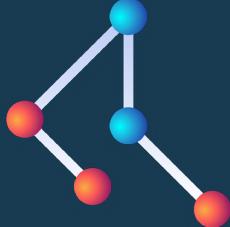
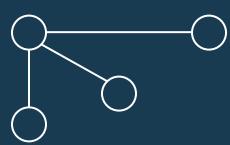
Improves the precision  
of cell type annotation

multi-batch and multi-omic  
integration

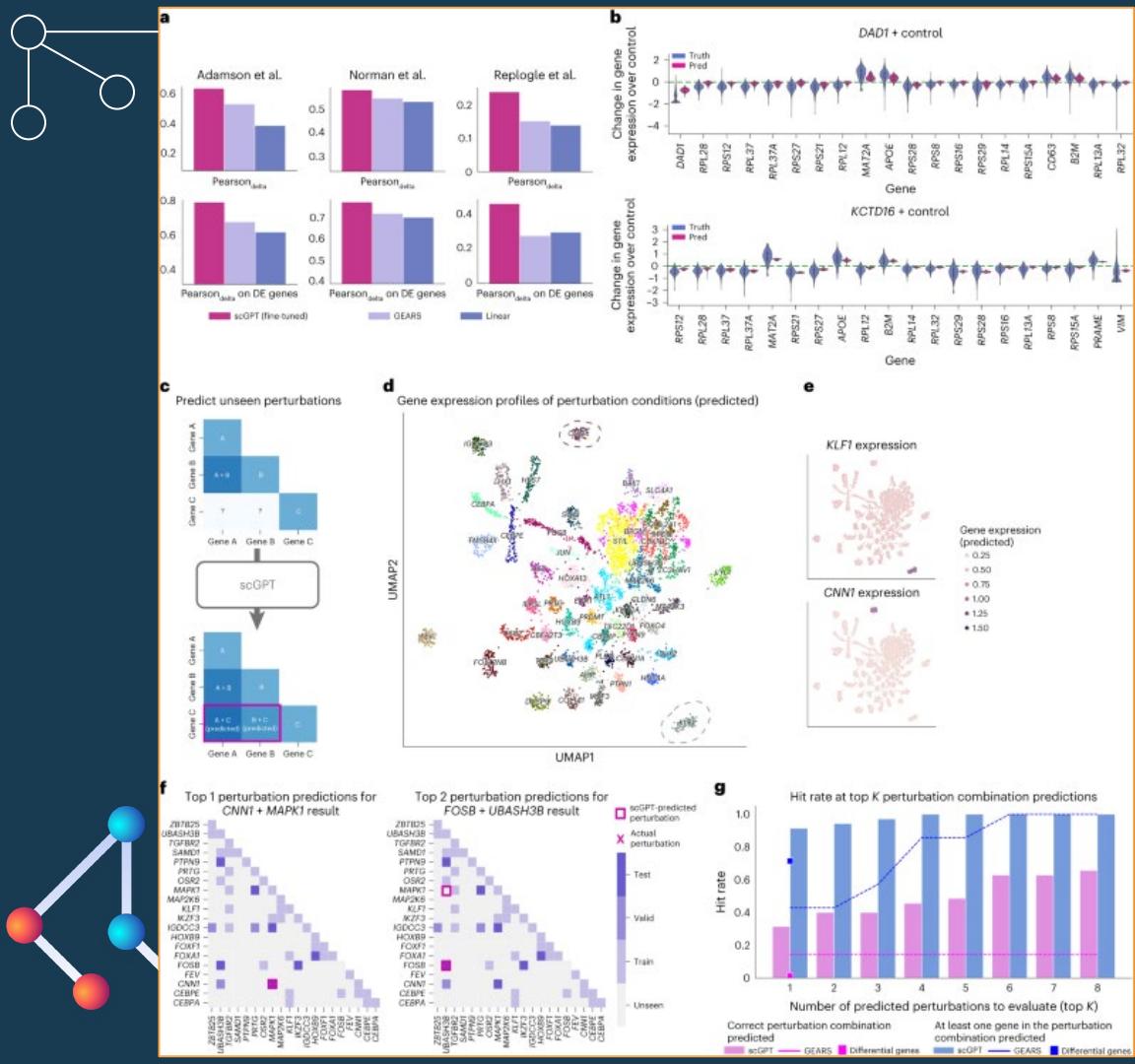
Predicting Unseen  
Genetic Perturbation  
Responses

Uncovers gene networks  
for specific cell states

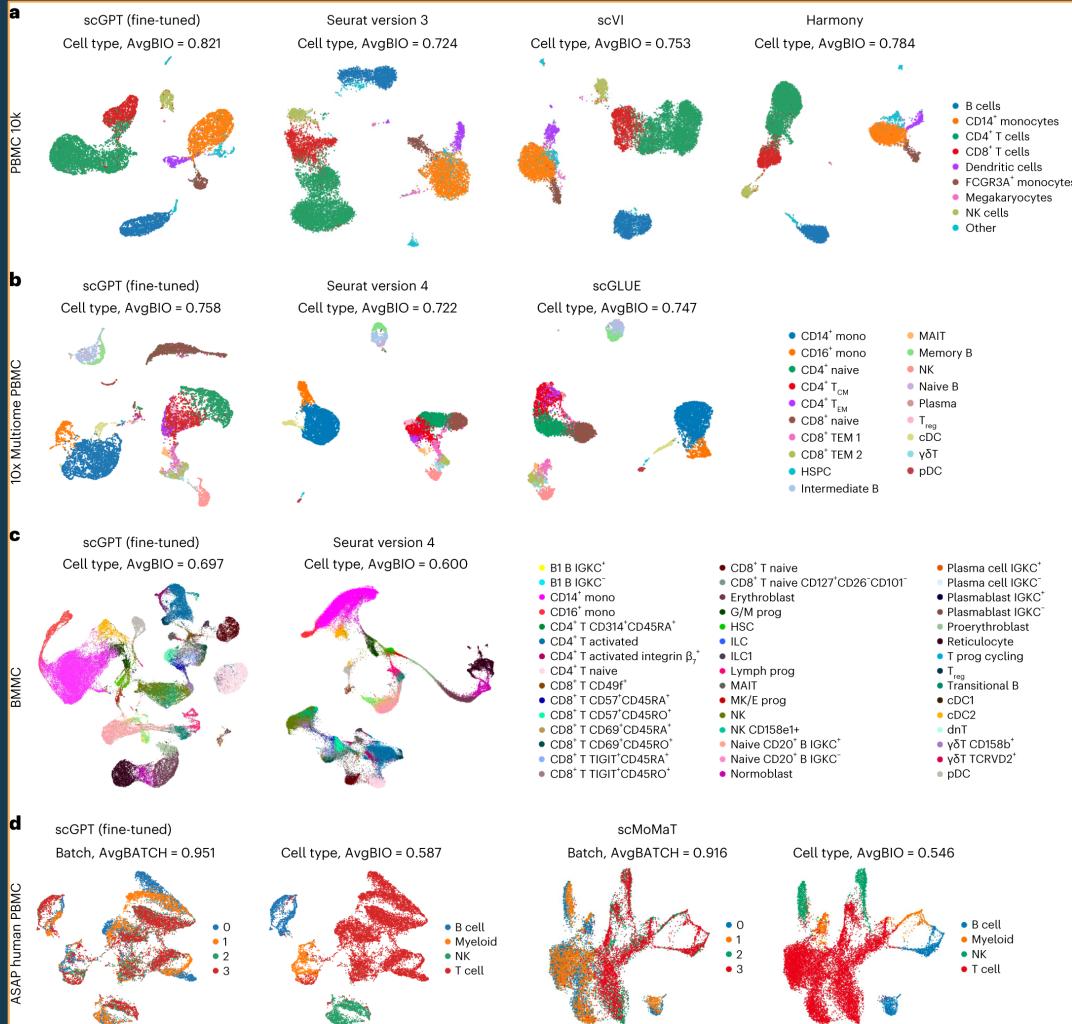
# Cell Type Annotation

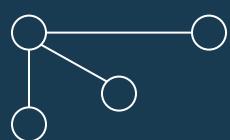


# Predicting Unseen Genetic Perturbation Responses



# Multi -Batch & Multi -Omic Integration

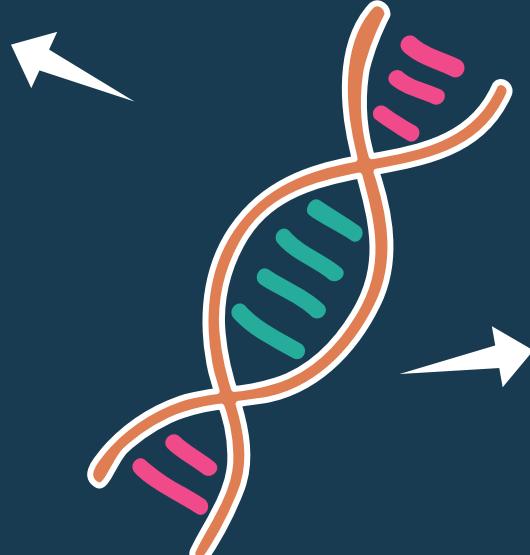




# Conclusion

## Limitations

- Pretraining does not mitigate batch effects.
- zero-shot performance could be constrained on datasets with technical variation
- Evaluating the model is also complex due to variation in data quality



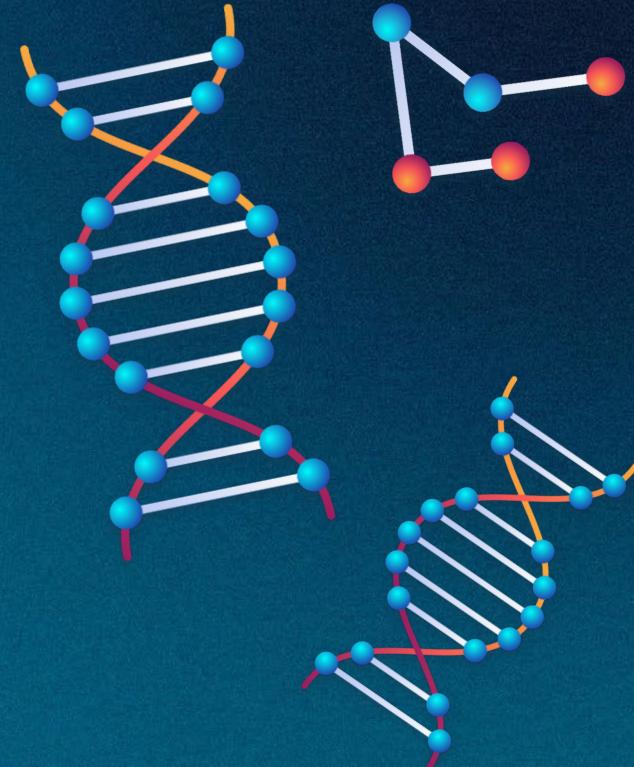
## Future Work

- pretrain on a larger-scale dataset with more diversity
  - explore in-context instruction learning for single-cell data.

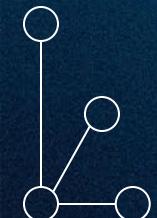
# Summary

---

Conclusion & Questions



⋮ ⋮ ⋮ ⋮



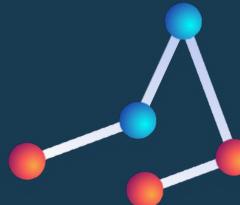
Model	DNABERT (BERT based)	UTR-LM	scGPT
Domain	DNA sequencing	5'UTR of mRNA	(scRNA-seq)
motive	Deciphering DNA sequences	Unified foundation model to study function of 5'UTR	Unified foundation model to study single-cell RNA functions
Method	<ul style="list-style-type: none"><li>BERT architecture</li><li>Tokenization with k-mer (6)</li><li>Modify pre-training process</li><li>Fine-tuned on 3 specific tasks</li><li>Benchmark with current tools</li></ul>	<ul style="list-style-type: none"><li>Transformer-based architecture</li><li>Masked nucleotide (MN) prediction</li><li>secondary structure (SS)</li><li>minimum free energy (MFE)</li><li>Fine-tuned on multiple downstream tasks</li></ul>	<ul style="list-style-type: none"><li>Transformer-based architecture</li><li>Pretrained on a large corpus of single-cell RNA data</li><li>tokenization of gene expression profiles.</li><li>Multi-task learning approach</li></ul>
Results	<ul style="list-style-type: none"><li>surpassing existing tools</li><li>Enhanced performance with limited data</li><li>No- separate training needed</li><li>Flexible learning of DNA in different situations</li></ul>	<ul style="list-style-type: none"><li>outperforms the best-known baseline in each task.</li><li>Performance not limited by sequence length</li><li>Validated through wet-laboratory experiments</li><li>Zero shot generalization</li></ul>	<ul style="list-style-type: none"><li>Pretrained model extrapolates to unseen datasets.</li><li>Outperform existing models</li><li>High accuracy in cell type annotation</li><li>strong scaling properties</li></ul>
limits	<ul style="list-style-type: none"><li>Sequence Length Limitation</li><li>Dependence on k-mer Tokenization</li></ul>	<ul style="list-style-type: none"><li>Computationally expensive</li></ul>	<ul style="list-style-type: none"><li>Pretraining does not mitigate batch effects.</li><li>zero-shot performance could be constrained on datasets with technical variation</li></ul>

• • • :



# Questions

Paper	Question
Paper 1	
Paper 2	
Paper 3	





Thank you  
for listening

---

Any More Questions?

