

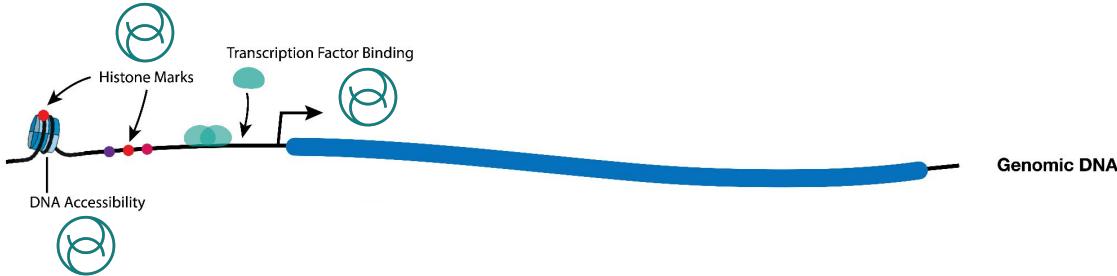
# *Foundation models for genome biology and drug discovery*

Presented by Albi Celaj for Deep Genomics 

# The Deep Genomics team is growing!



# Genome biology is complicated - RNA plays a central part and ML models are needed!



## Deep Genomics Models

**Algorithms:** Multi-task learning, Pretraining, Fine tuning, Zero-shot learning

**Architectures:** Transformers, Structured state space models, Convolutional layers

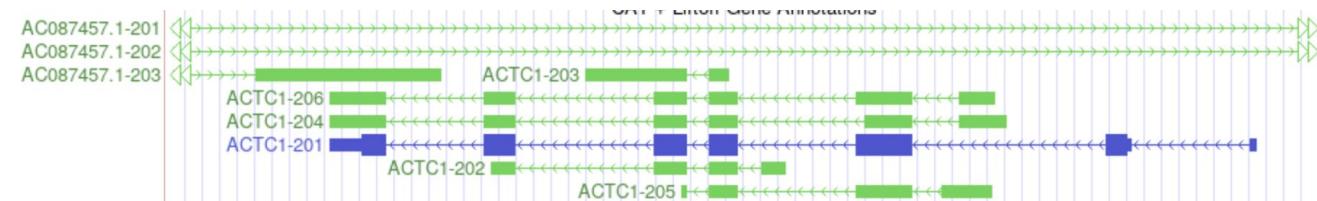
# Encoding genome biology through genome tracks: A treasure trove of experimental information for the human genome



## Raw data: DNA Sequence

chr15: | | | | | 32,596,000 | | | | 32,597,000 | | | | 32,598,000 | | | | 32,599,000 | | | | 32,600,000 | | | | 32,601,000 |  
CAGTACGTCAGTCAGTACGT...

## Human labels: Genes and Transcripts

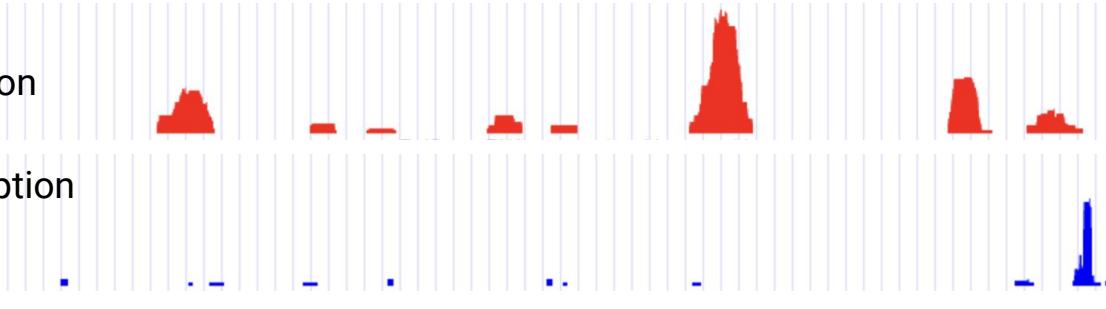


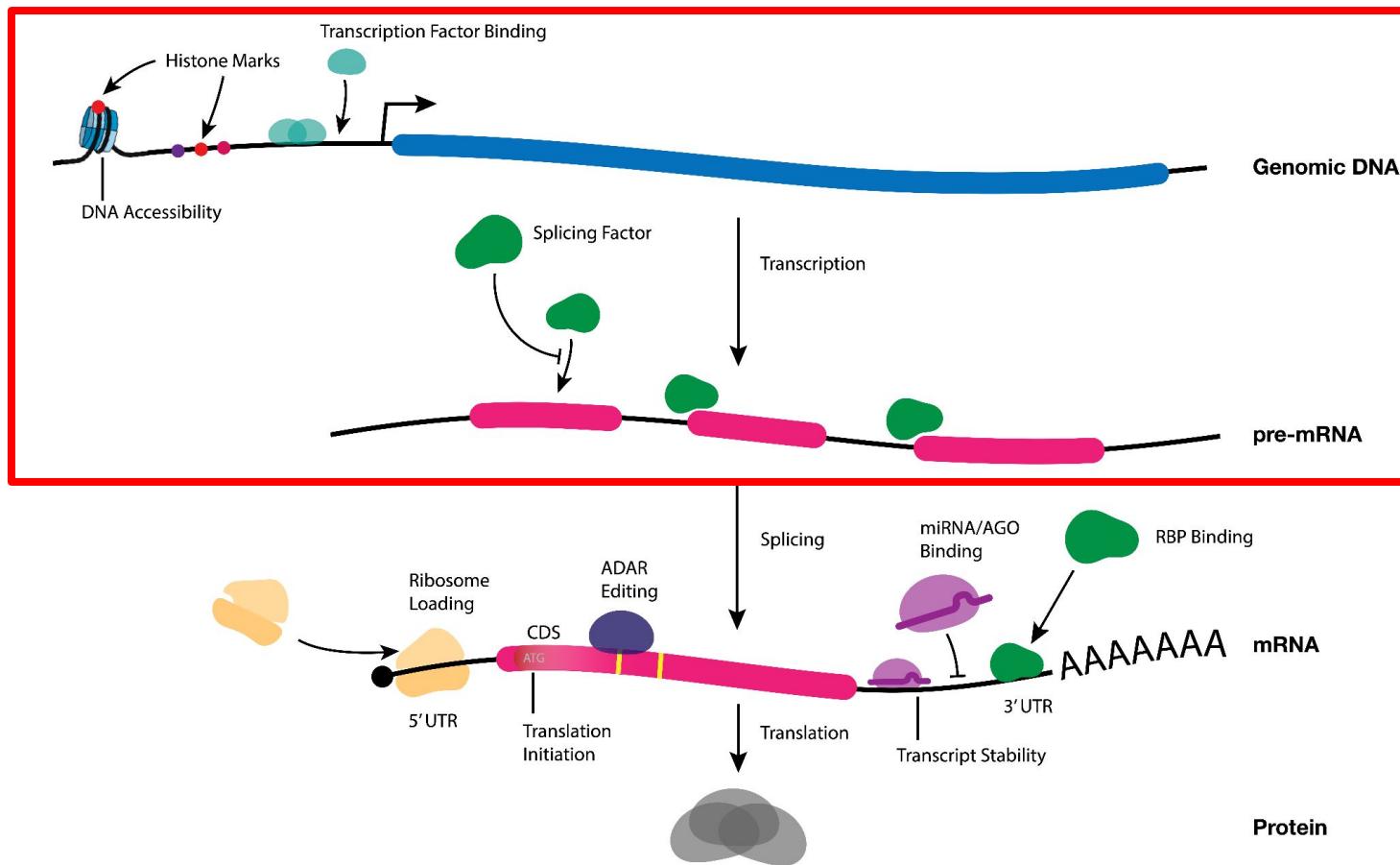
## Experimental labels: Many experiments on the human genome!

RNA  
expression

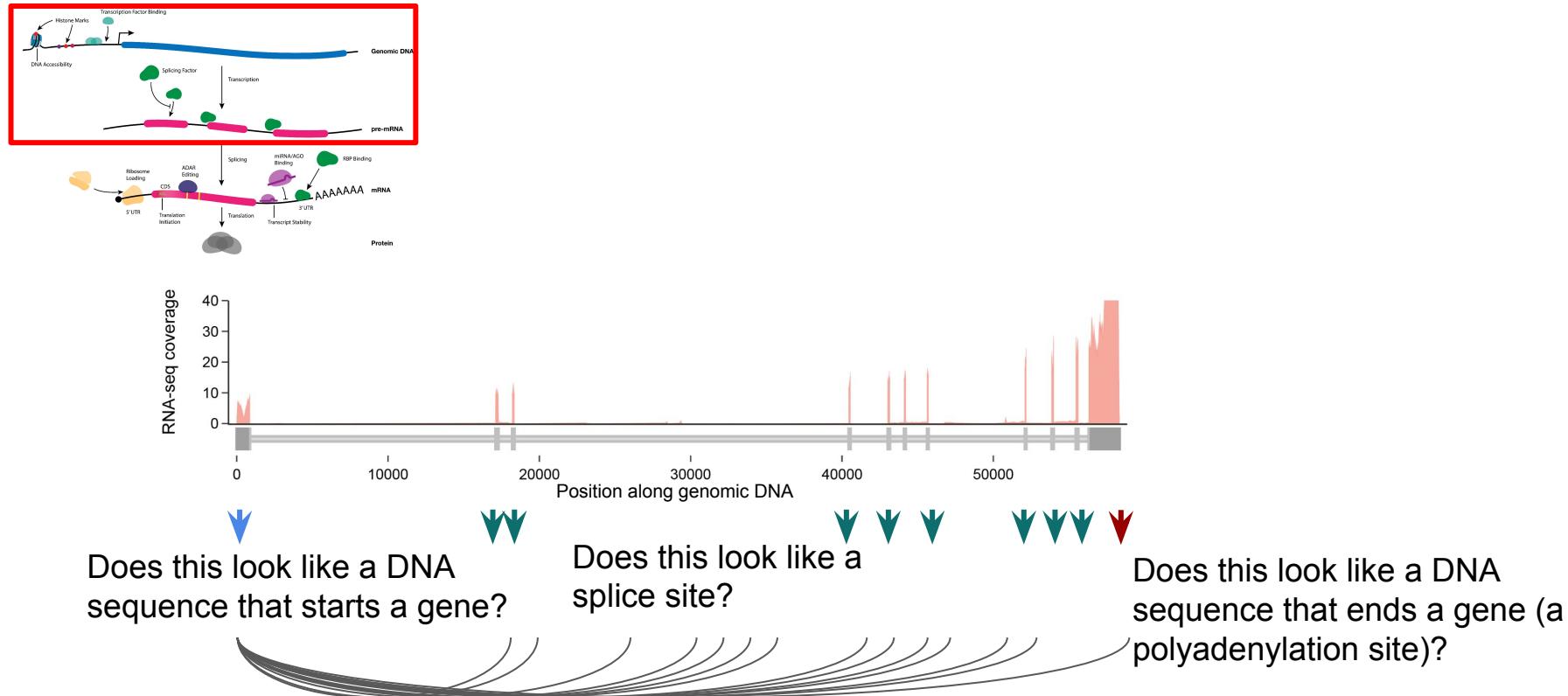


Transcription  
activity



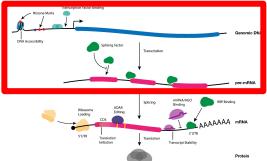


# Modeling some RNA genome tracks needs long-range information



# Attention mechanism for DNA sequence

CNNs can typically encode short range information about DNA



DNA-seq coverage

0 5 10 15 20



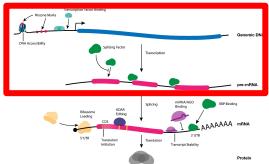
← Promoter TATAAA↑T

← Splice sites GTATGT↓

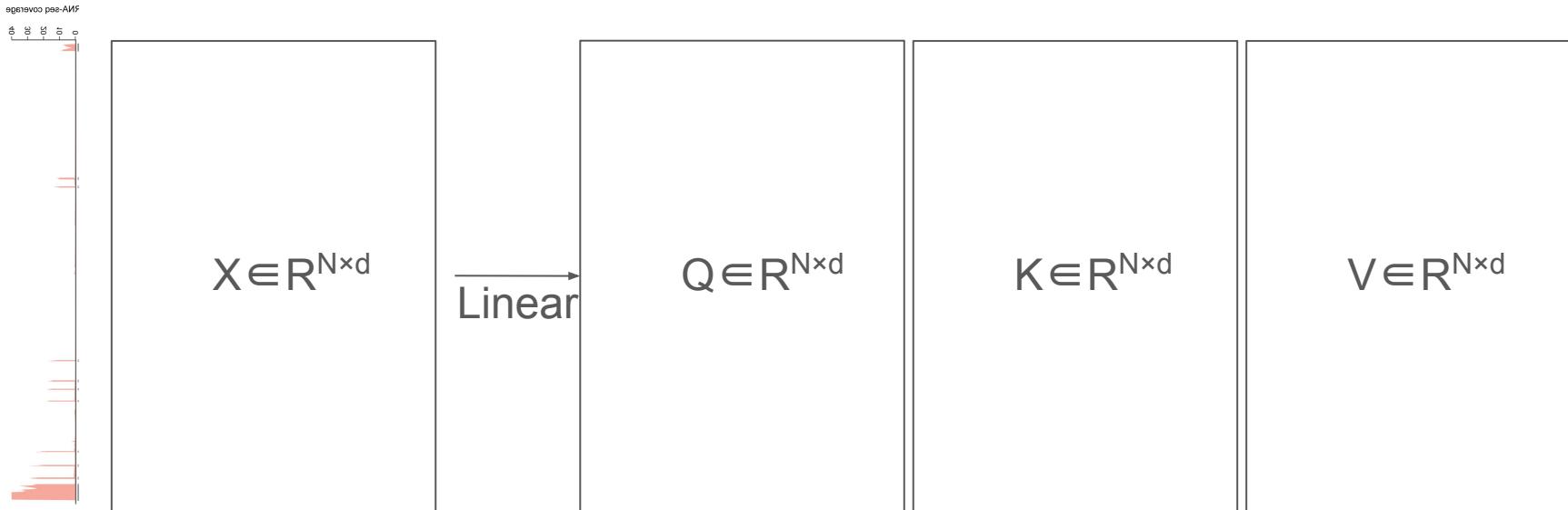
← Polyadenylation site AA↑AAA

Embedding matrix X - N DNA basepairs by d hidden channels

# Attention mechanism for DNA sequence



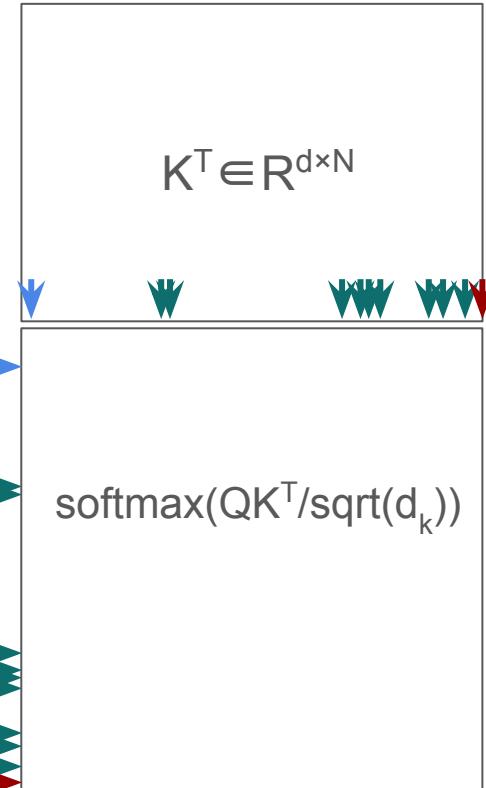
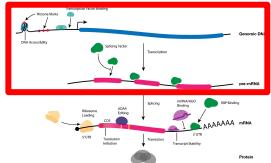
**Step 1:** Project DNA embeddings into Q, K, V (query, key, and value matrices)



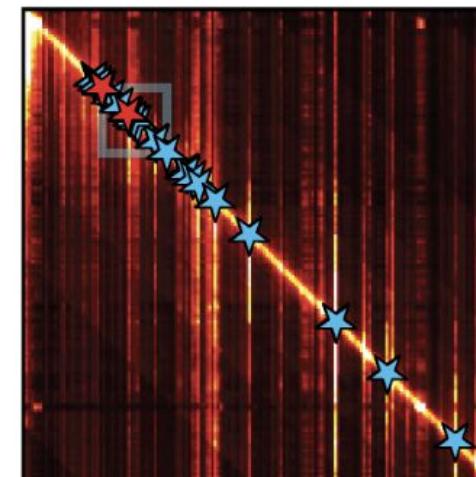
Embedding matrix  $X$  - N DNA basepairs by d hidden channels

# Attention mechanism for DNA sequence

Step 2: Allow each token to attend to each other token

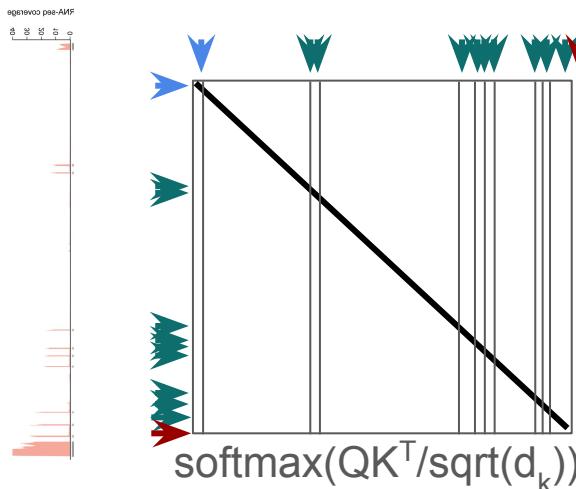
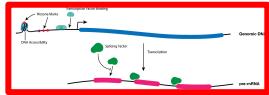


Linder et al 2023



# Attention mechanism for DNA sequence

Step 3\*: Transform the embeddings based on self-attention

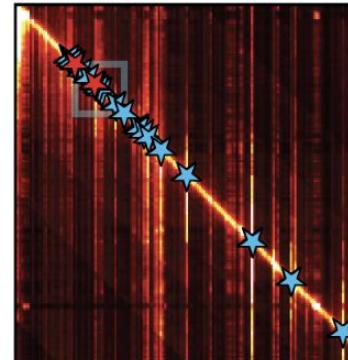


Embedding matrix  $X$  - N DNA basepairs by d hidden channels

$$X \in \mathbb{R}^{N \times d}$$

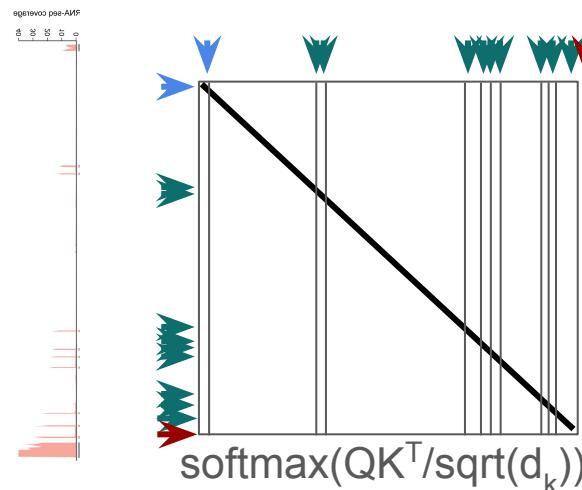
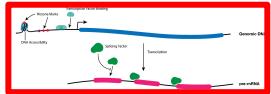
$$X' \in \mathbb{R}^{N \times d}$$

Linder et al 2023



# Attention mechanism for DNA sequence

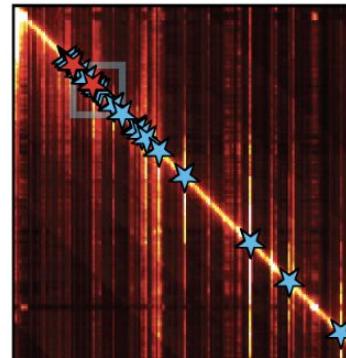
Step 3: Transform the embeddings based on self-attention



$$V \in \mathbb{R}^{N \times d}$$

$$X' \in \mathbb{R}^{N \times d}$$

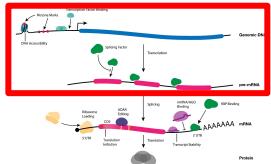
Linder et al 2023



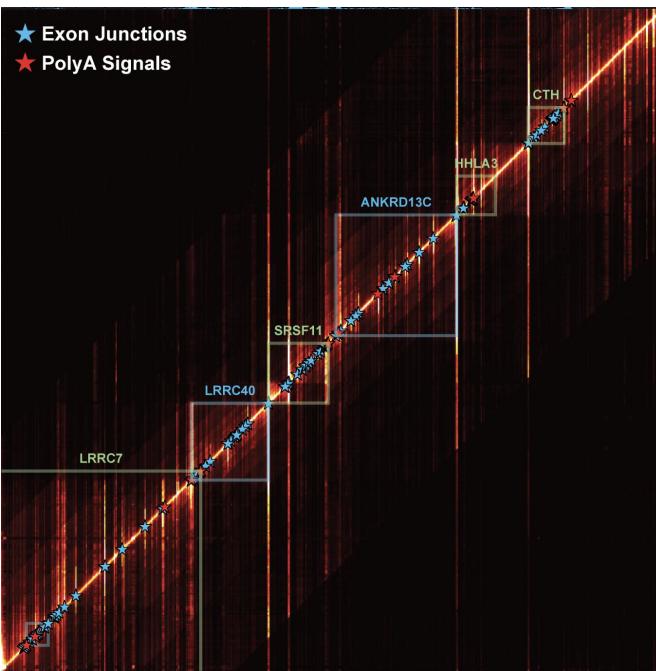
Embedding matrix X - N DNA basepairs by d hidden channels

# Attention mechanism for DNA sequence

Caveat: Also need to encode position of elements

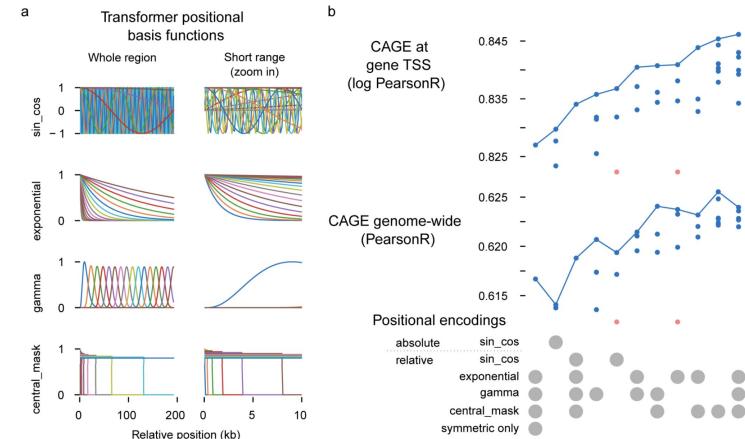


Linder et al 2023

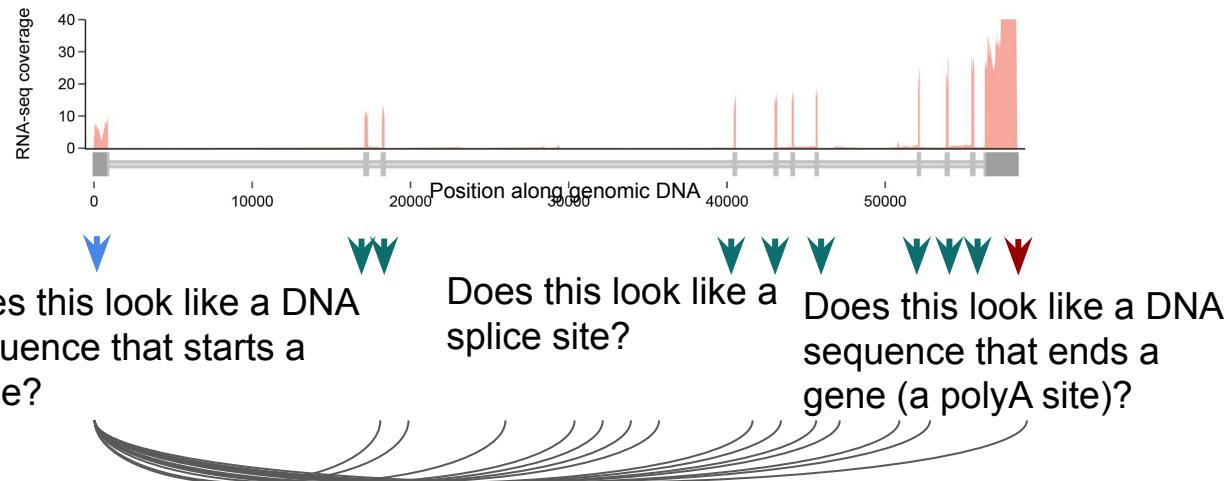
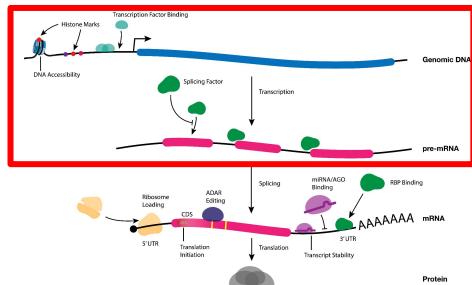


Avsec et al 2021

Extended Data Fig. 3: Custom relative positional encoding functions are required for good predictive performance



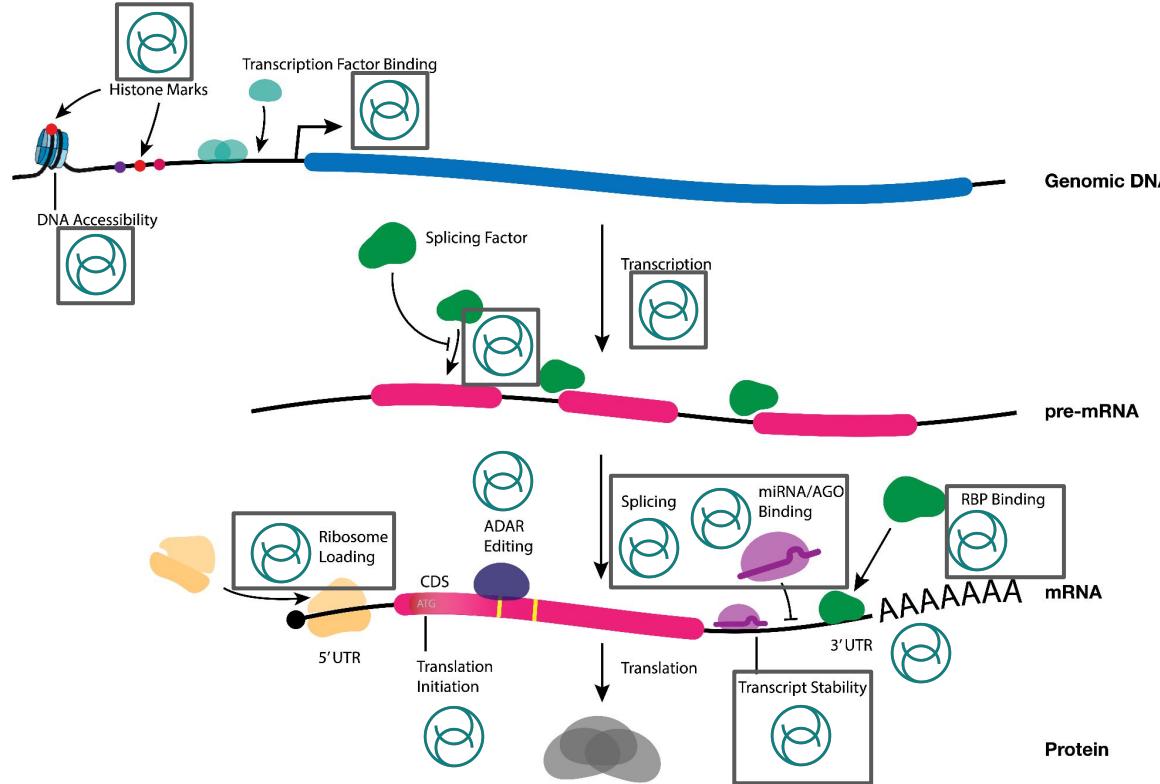
# Modeling some RNA genome tracks needs long-range information



This is a gene start and the region allows making an RNA. The RNA looks stable

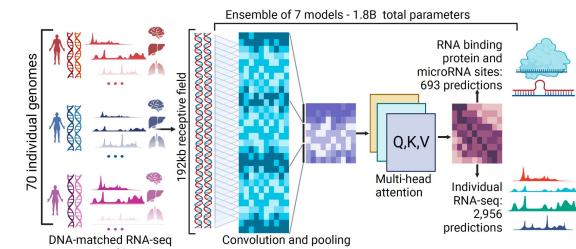
This is a splice site and we are within an RNA being transcribed. The RNA looks stable

# Using a foundation model to model multiple RNA processes



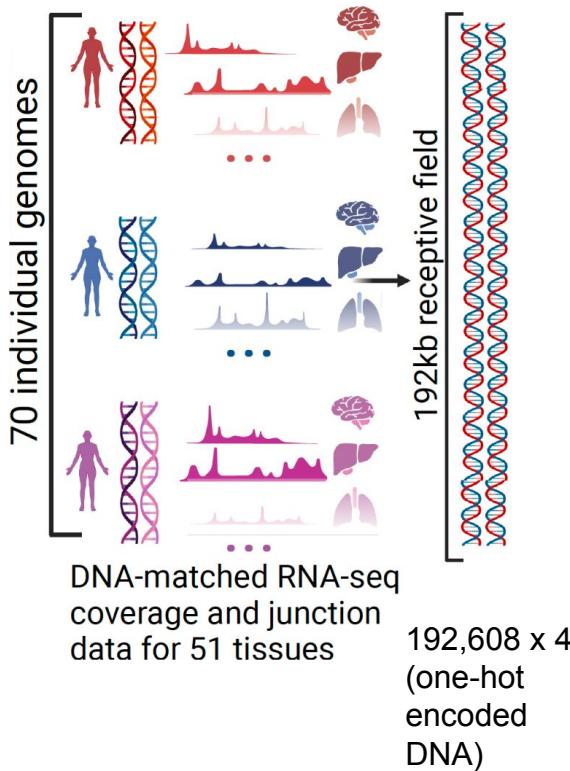
## BigRNA

Foundation model of RNA biology

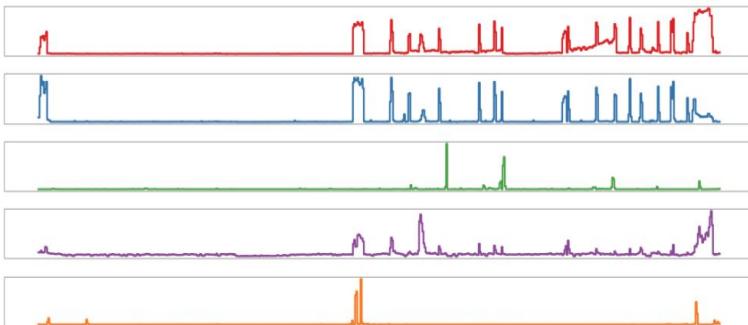
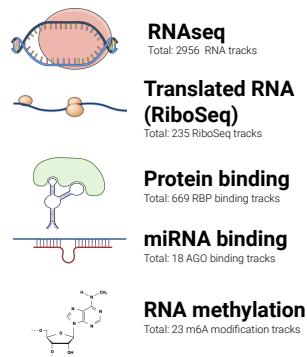
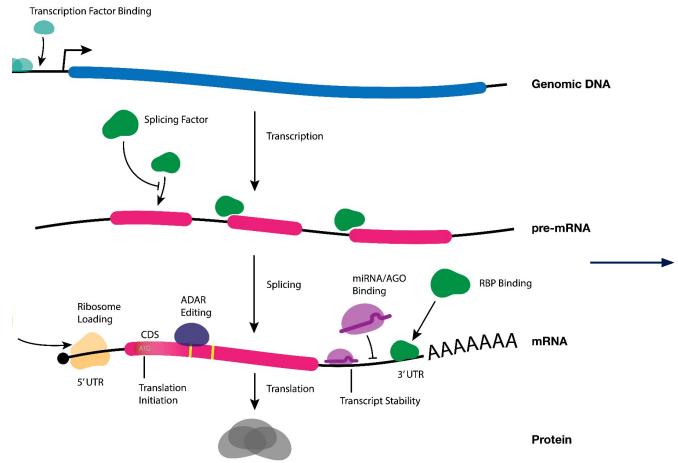


- Alston Lo  
Helen Zhu  
Vivian Chu  
Nicole Zhang

# BigRNA is an all-purpose “foundation model” for RNA genome biology

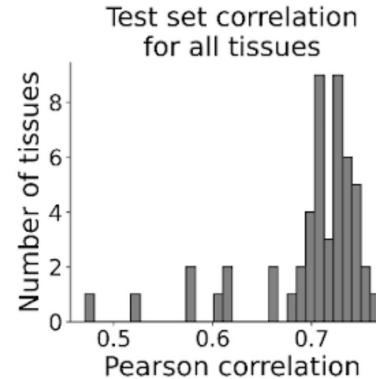
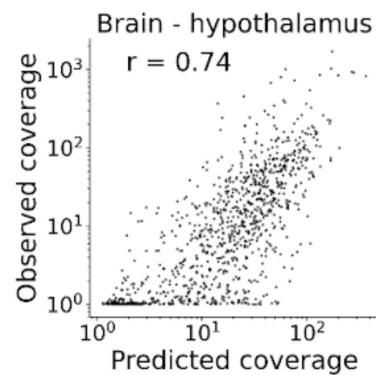
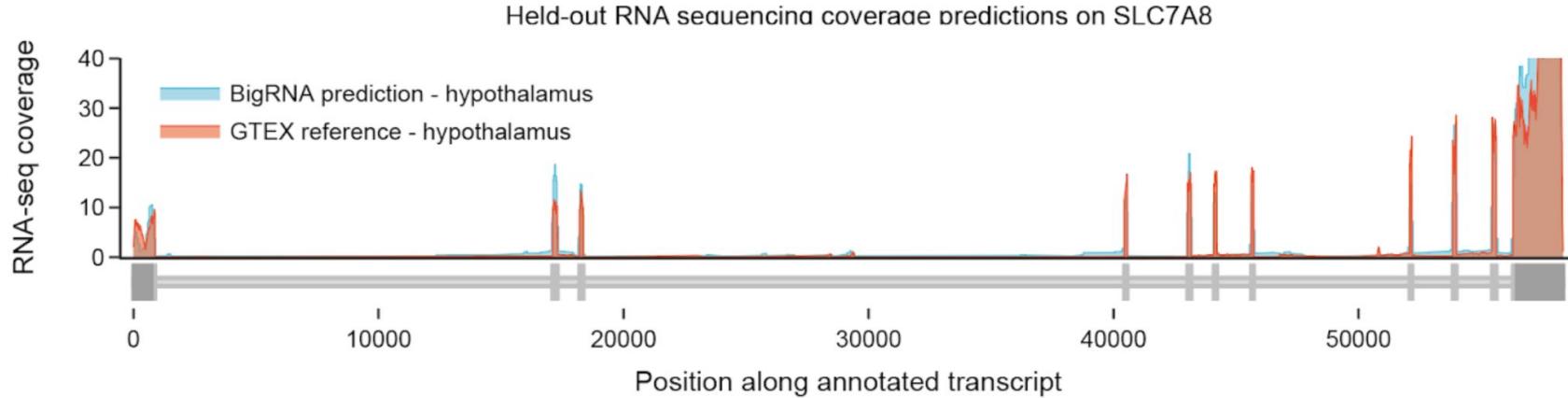


# BigRNA predicts many genome tracks from DNA sequence

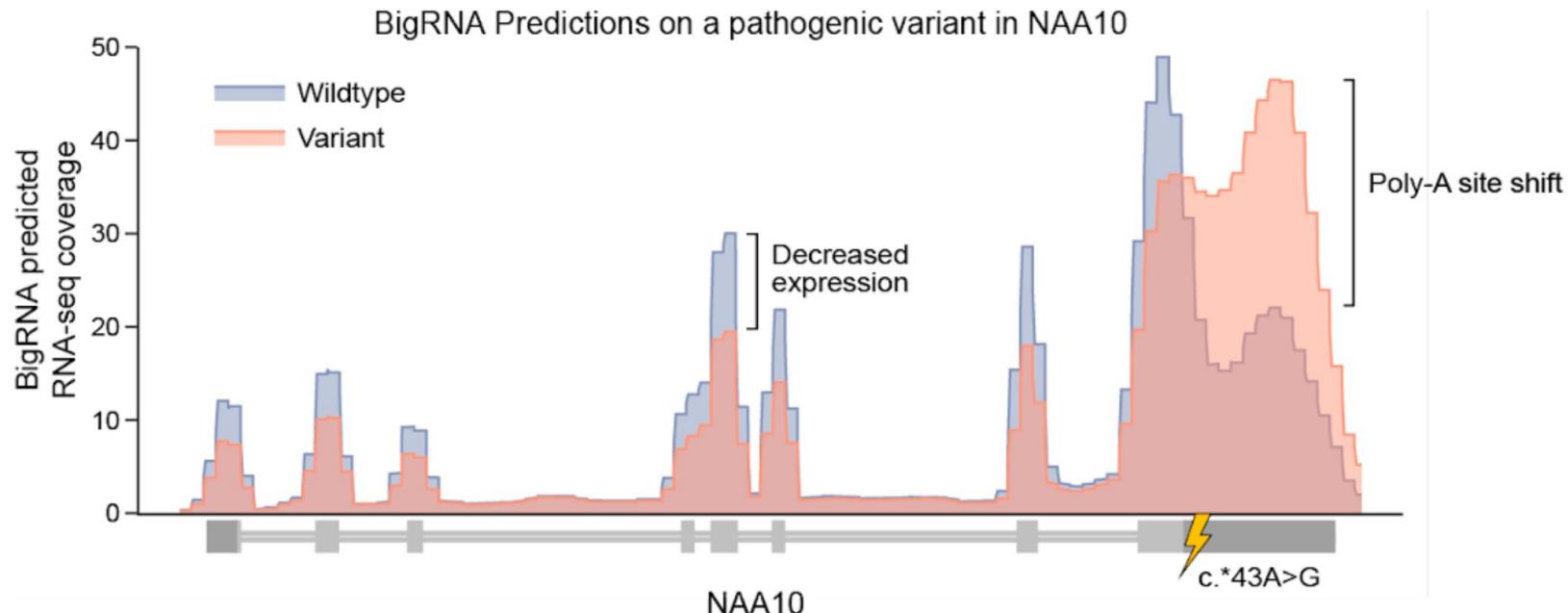


**Total: 3901 labels x 3 billion DNA base pairs**

# Predicting mRNA expression



# Predicting polyA-disrupting genetic variants

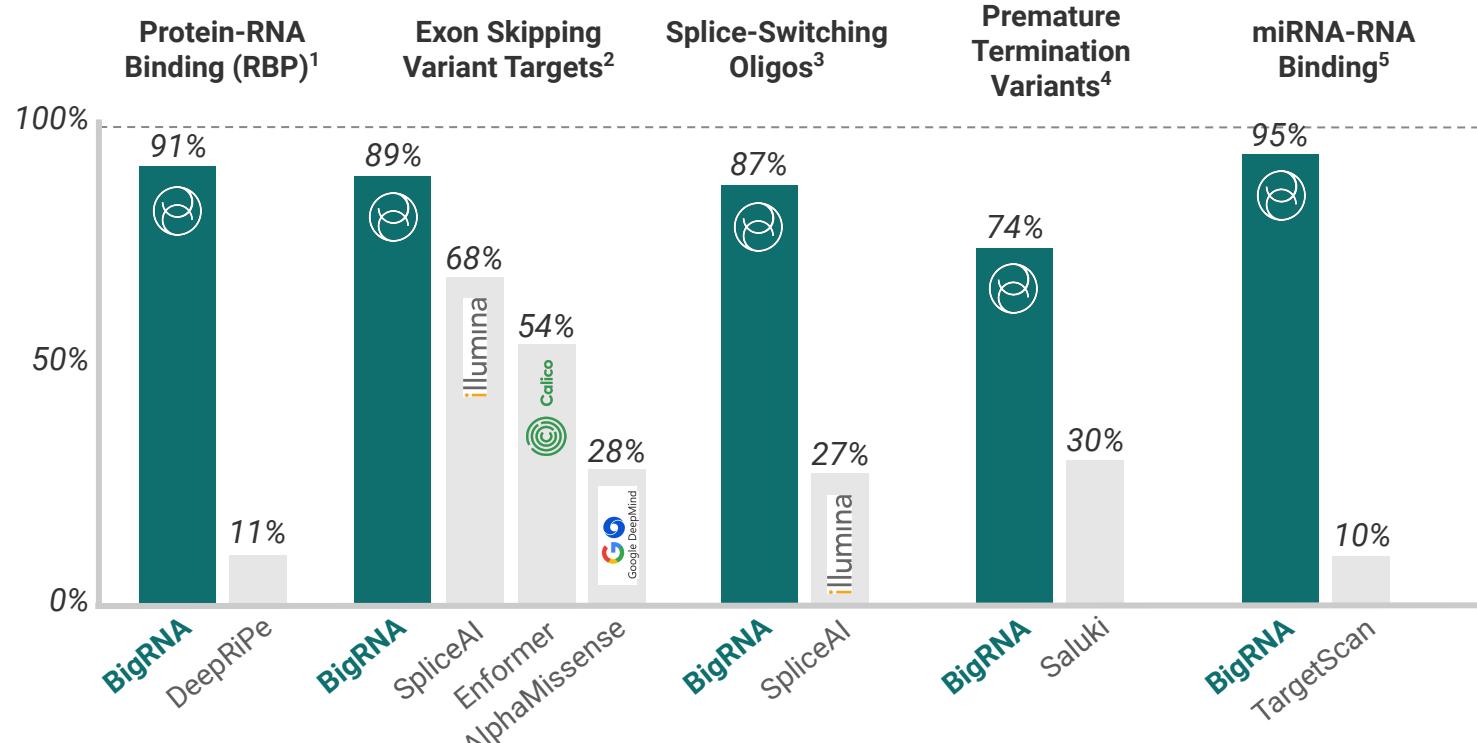


# BigRNA outperforms state-of-the-art models across a range of tasks



## Accuracy on Held Out Test Data (see notes)

1. Fraction of 141 proteins for which average precision for binding prediction  $\geq 0.1$ .
2. True positive rate at FPR=0.15 for classifying MaPSy variants (341 positives, 754 negatives).
3. Fraction of 15 exons (in 12 genes) for which predicted exon inclusion of oligos correlate with measured values by  $\geq 0.45$ .
4. Fraction of premature termination codon variants classified correctly at a 10% FPR vs common 5' UTR SNVs.
5. Fraction of 19 cell models for which AUC  $\geq 0.75$  for identification of CLIP miRNA binding sites.



For more, see: <https://www.biorxiv.org/content/10.1101/2023.09.20.558508v1>

# BigRNA acts as a zero-shot model of disease-causing mutations and RNA-blocking “oligo” drugs

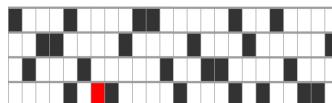


## Model input

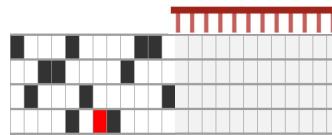
Healthy sequence



Disease sequence

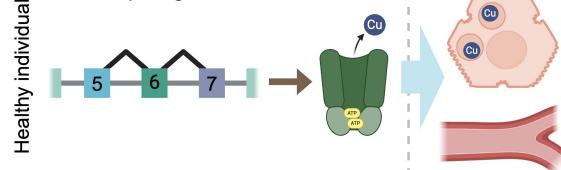


Disease sequence + RNA-blocking oligo



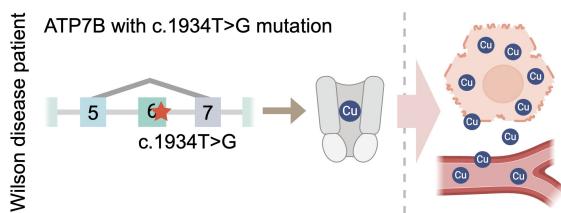
## Expected outcome

Normal splicing of ATP7B



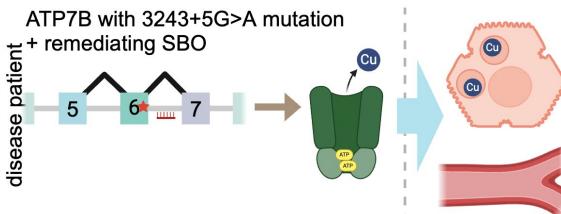
Healthy individual

ATP7B with c.1934T>G mutation

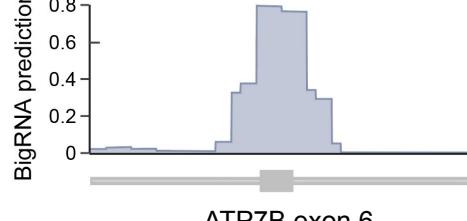


Wilson disease patient

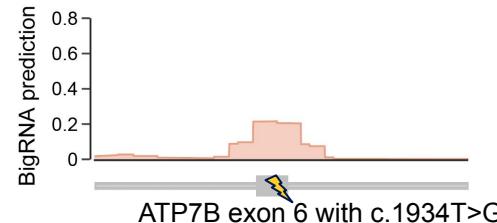
ATP7B with 3243+5G>A mutation + remediating SBO



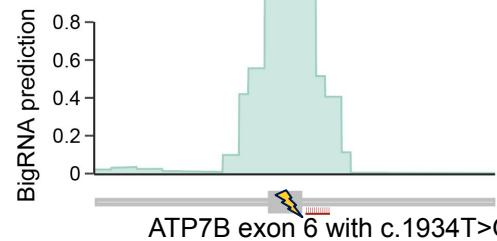
## Model Output



ATP7B exon 6



ATP7B exon 6 with c.1934T>G



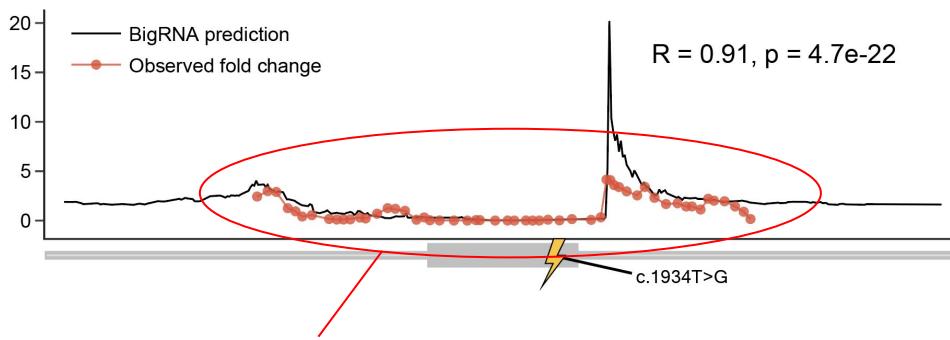
ATP7B exon 6 with c.1934T>G + oligo



Not actual size

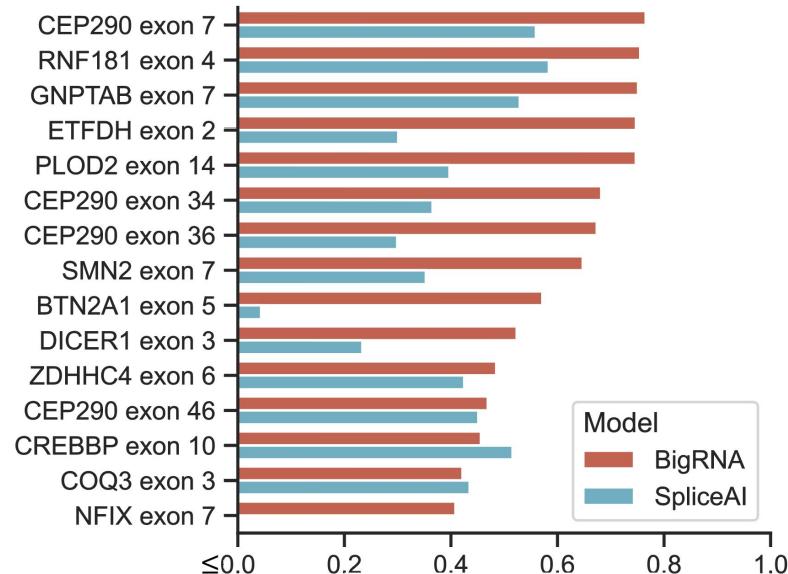
# BigRNA predicts the top compound, and correlates remarkably well with multiple experiments

Predicted versus observed RNA-blocking drug effects on ATP7B exon 6



BigRNA prediction achieves strong correlation with experimental measurements

Predicted versus observed RNA-blocking drug effects in fifteen different experiments



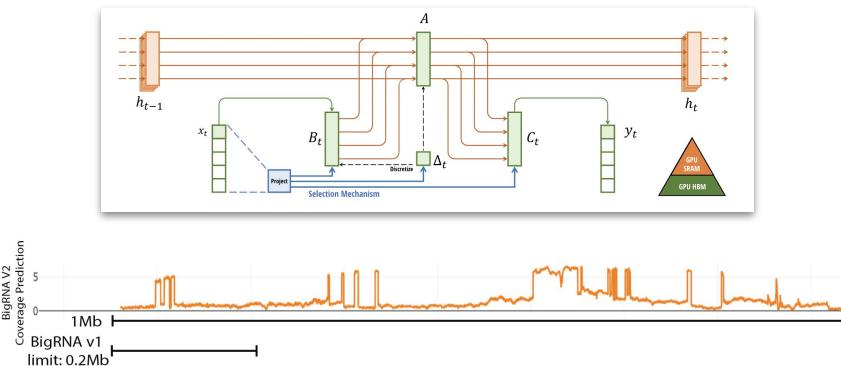
Correlation with experimentally-observed RNA-blocking drug effect

# Latest developments: Newest version of BigRNA is high resolution and high context



## From transformer ensemble to structured state space architecture (mamba)

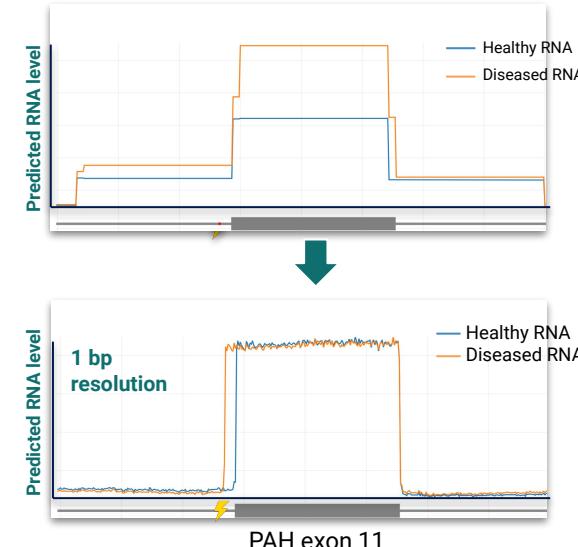
- Flexible input size **allows predictions for all genes**
- **Enables scaling** via higher memory efficiency
- Increased accuracy for **variant predictions** and **therapeutics**

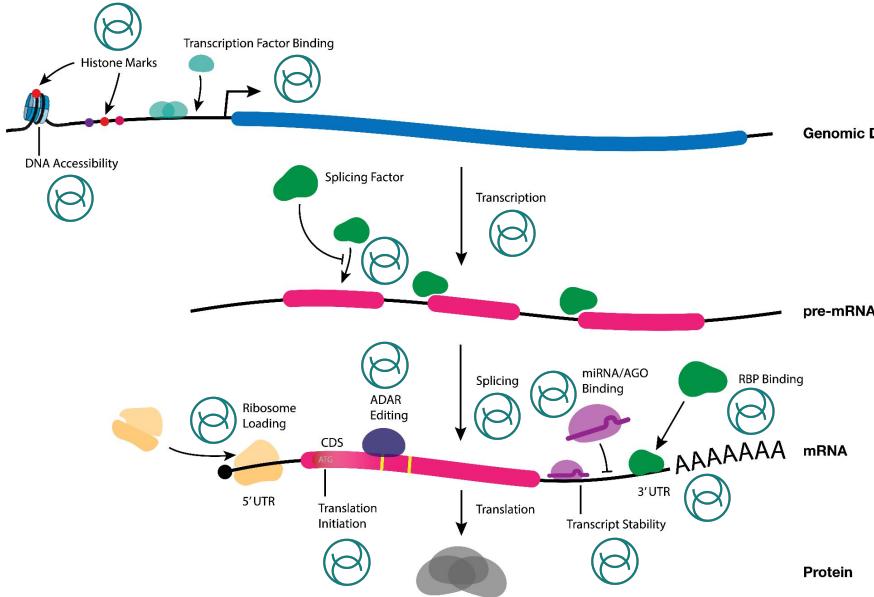


Alston Lo

## Increase resolution over many other models to 1-bp

- Contrast Enformer (128 bp), Borzoi (32 bp)
- Detects short RNA **changes that are drivers of disease**





Sara Pour



Andrew Jung



Vivian Chu



Rory Gao



Alston Lo



Nicole Zhang



Junru Lin



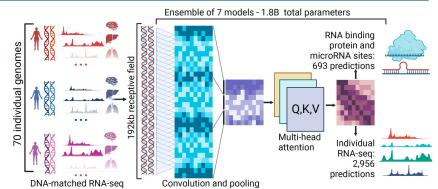
Helen Zhu



## Deep Genomics Models

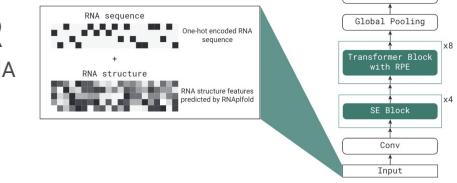
### BigRNA

Foundation model of RNA biology



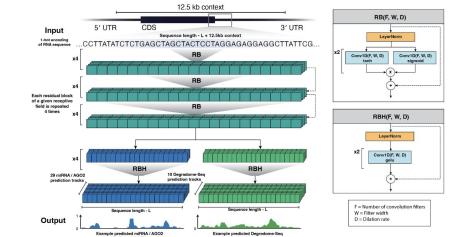
### DeepADAR

Deep learning for RNA editing oligo design



### REPRESS

Deep learning for repressive RNA elements



# What we're working on

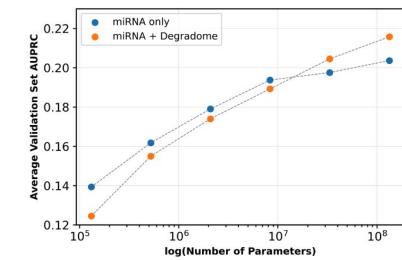
## ❑ Fine-grained predictions of individual cells

Single-cell data, individual cell types and cell states



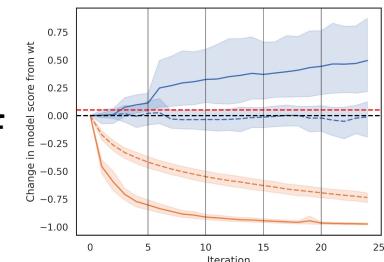
## ❑ Improving accuracy and understanding its limits

Type and quantity of training data, e.g. multiple species and learning from human variation, improvements to model architecture and training objective



## ❑ Beyond zero-shot learning: fine-tuning on drug screening data

Impressive zero-shot performance does not take advantage of internal data





# *Thank you!*

Presented by Albi Celaj for Deep Genomics

