



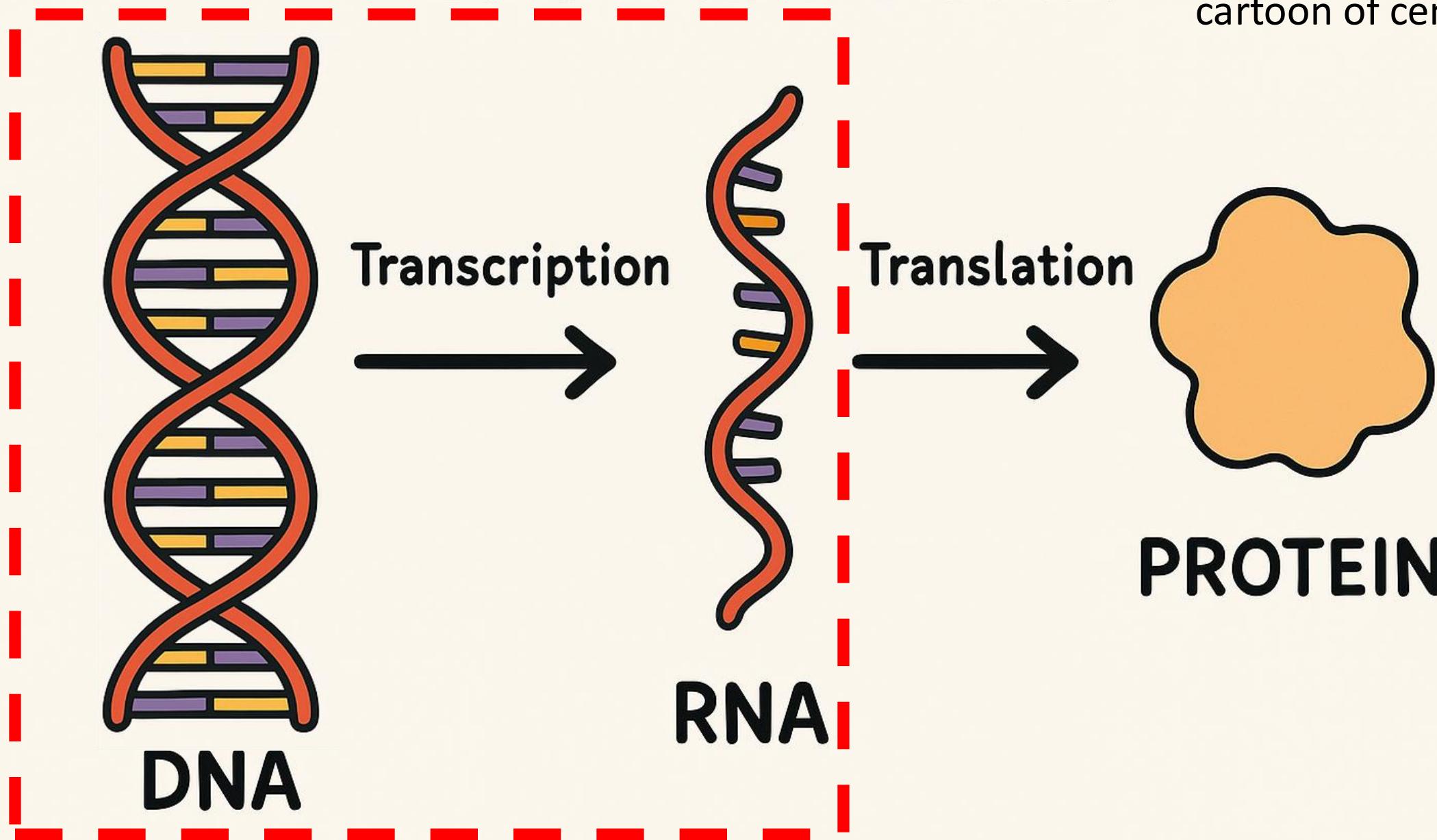
# Teaching machines to learn biology: splice site prediction and gene expression prediction

Kuan-Hao Chao

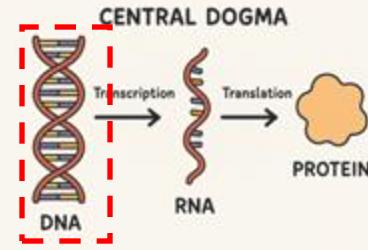
2025.04.02

 ChatGPT 4o:  
Generate schematic  
cartoon of central dogma

# CENTRAL DOGMA

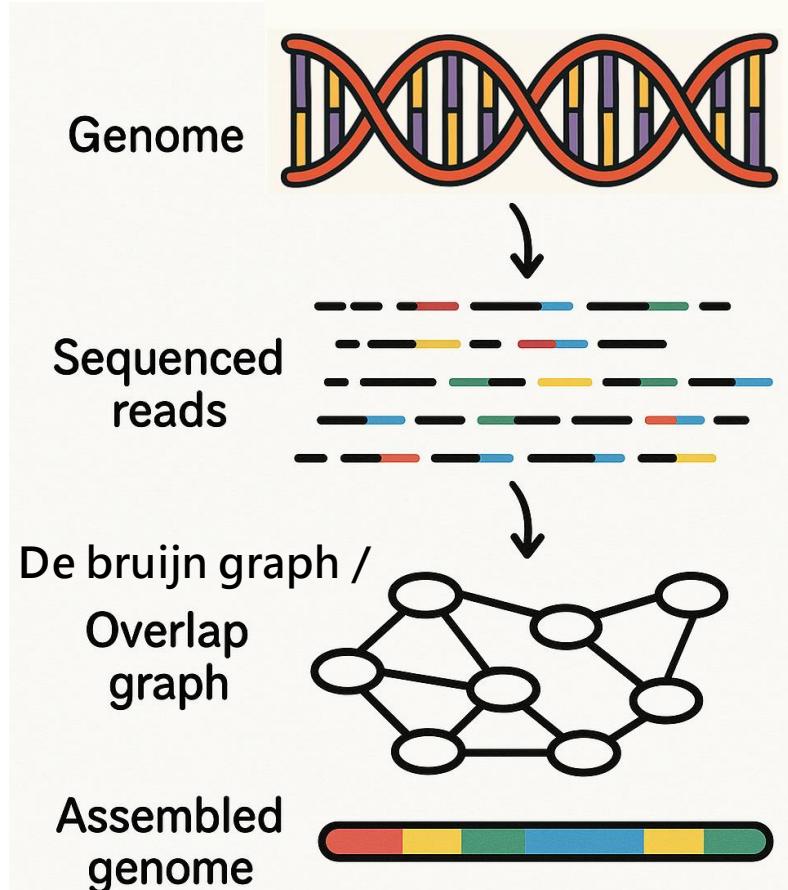


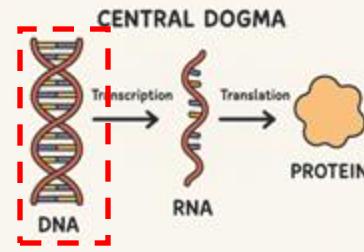
# Genome Assembly



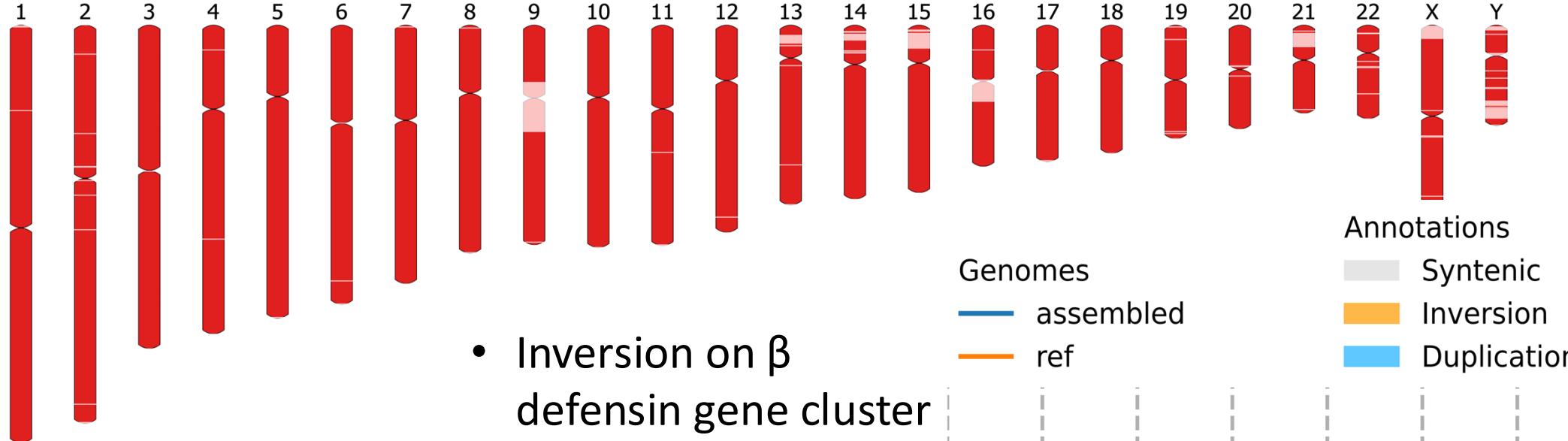
Wellcome genome bookcase

# 3 billions of ACGT!



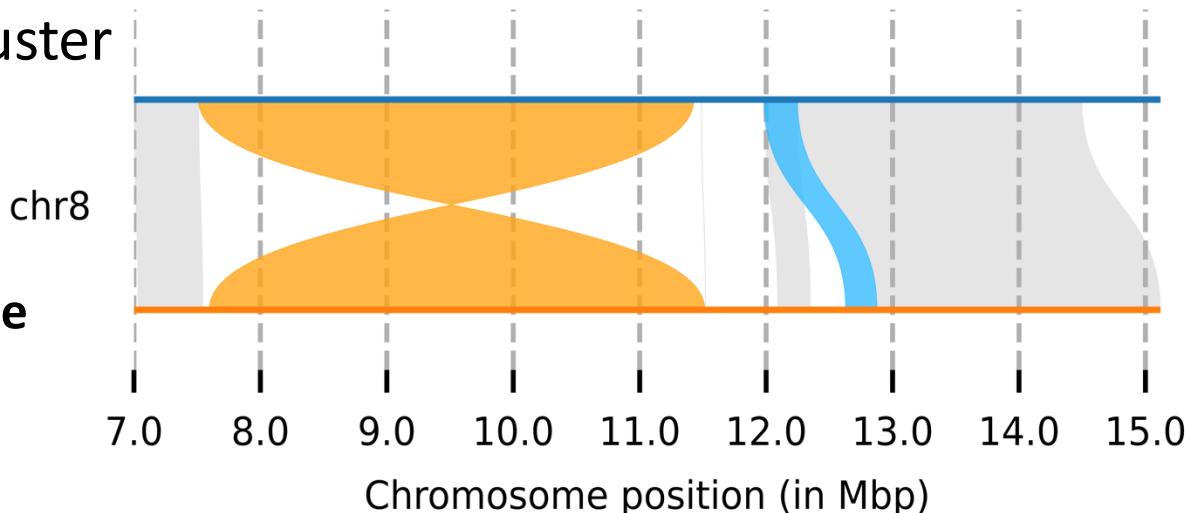


# Genome Assembly: Han1



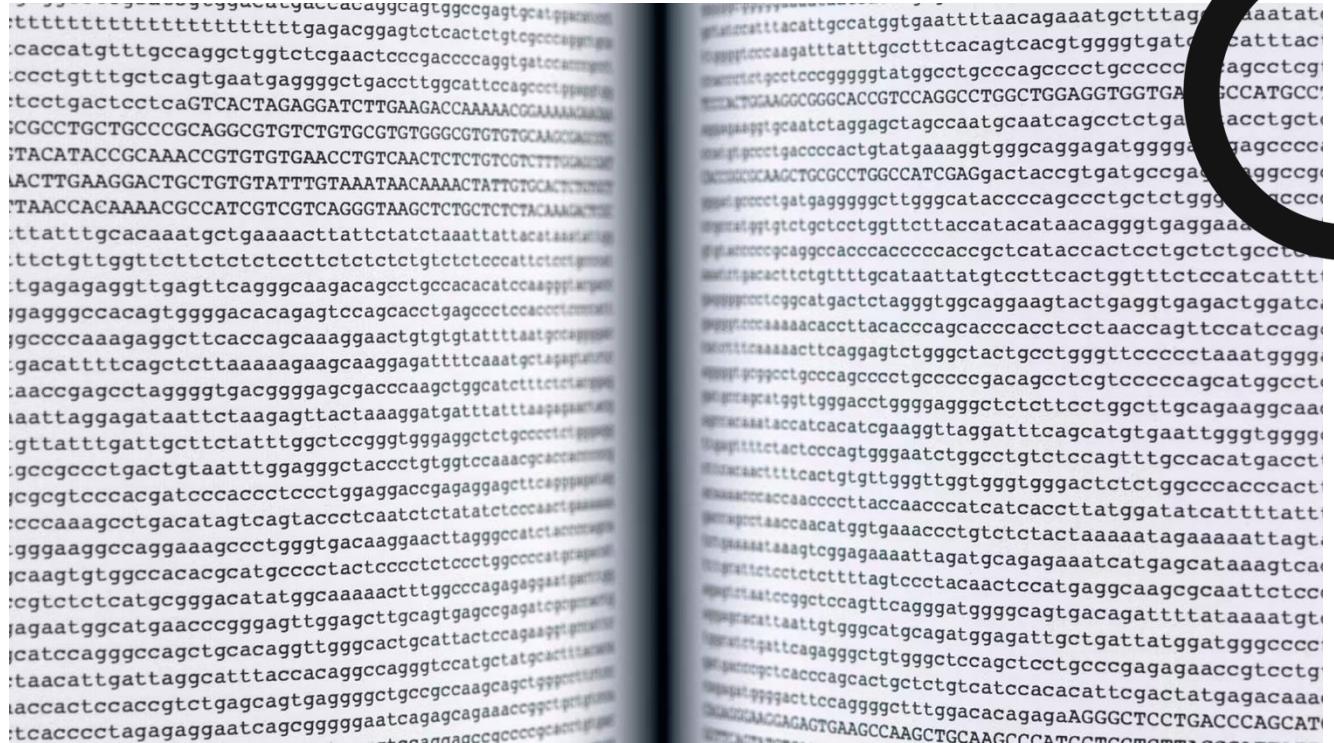
**Han1: 1<sup>st</sup> Gapless Southern Han Chinese**

Chao, K. H., Zimin, A. V., Pertea, M., & Salzberg, S. L. (2023). The first gapless, reference-quality, fully annotated genome from a Southern Han Chinese individual. *G3: Genes, Genomes, Genetics*, 13(3), jkac321.

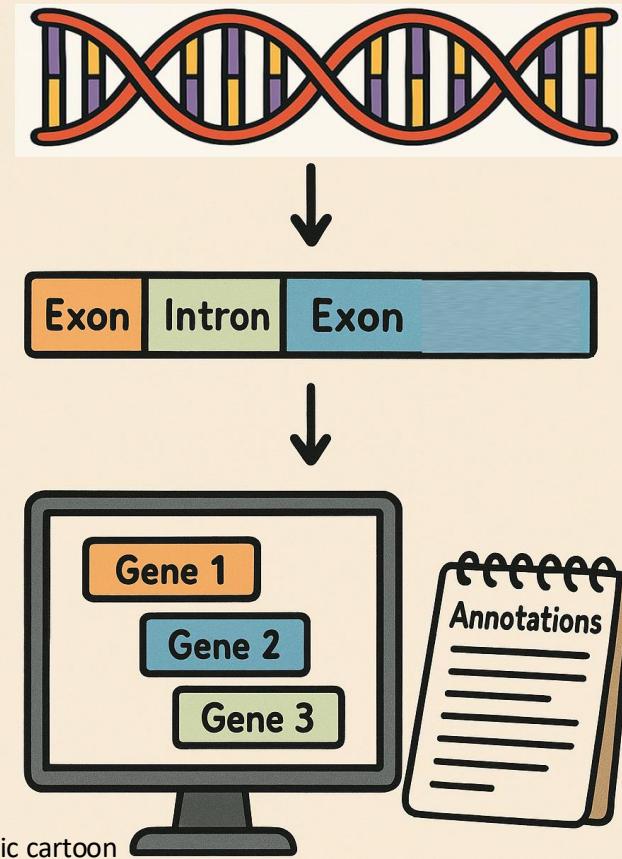


# Genome Annotation

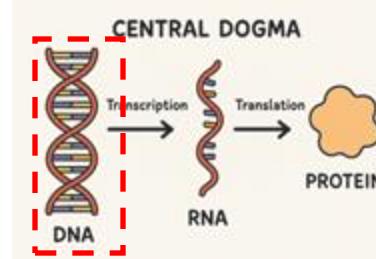
Where are the meaningful regions (Genes)?



## Genome Annotation



ChatGPT 4o:  
Generate schematic cartoon  
of genome annotation



# Genome Annotation

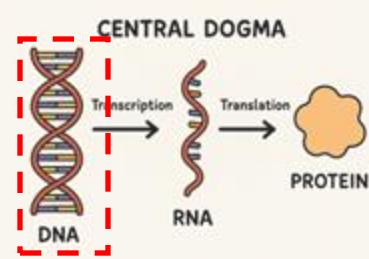


## Genome (FASTA)

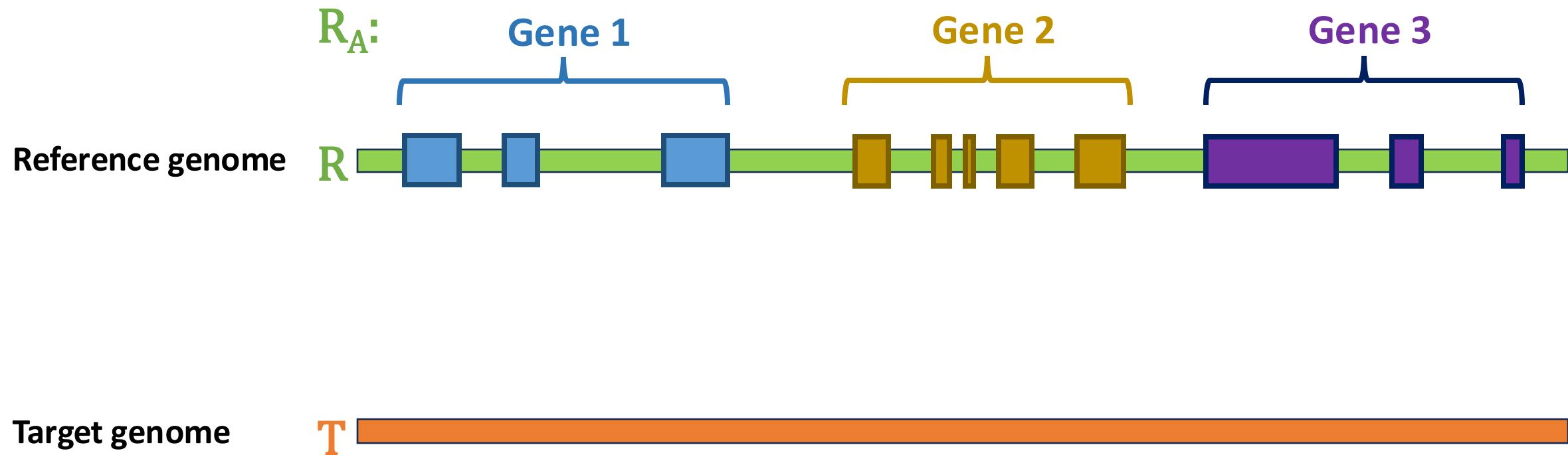
```
CAGCCCCGGAGACTtaaatacaggaagaaaaaggCAGGACAGAATTACAAGGTGCTGCCAGGGCGGGCAGCGGCCCT
GCCTCCTACCCTTGCCTCATGACCAGCTTGAAGAGATCCGACATCAAGTGCCACCTGGCTCGTGGCTCTCACT
GCAACGGAAAGCCACAGACTGGGGTGAAGAGATTCAAGTCACATGCGACCGGTgactccctgtccccaccccatgACACT
CCCCAGCCCTCCAAGGCCACTGTGTTCCAGTTAGCTCAGAGCCTCAGTCGATCCCTGACCCAGCACCGGGACTGATG
AGACAGCGGCTGTTGAGGAGccacctcccagccacctcggggccagggccagggtgtGCAGCACCAGTACAATGGGG
AAACTGGCCCAGAGAGGTGAGGCAGCTTGCCTGGGTACAGAGCAAGGCAAAAGCAGCGCTGGGTACAAGCTAAAACC
ATAGTGCCCAGGGCACTGCCGCTGCAGGGCGAGGCATCGCATCACACCAAGTGTCTGCCTCACAGCAGGCATCATCAGTA
```

## Annotation (GFF / GTF)

```
chr1 BestRefSeq gene 450740 451678 . - . ID=gene-OR4F29;
chr1 BestRefSeq mRNA 450740 451678 . - . ID=rna-NM_001005221.2;Parent=gene-OR4F29;
chr1 BestRefSeq exon 450740 451678 . - . ID=exon-NM_001005221.2-1;Parent=rna-NM_001005221.2;
chr1 BestRefSeq exon 452658 453675 . - . ID=exon-NM_001005221.2-2;Parent=rna-NM_001005221.2;
chr1 BestRefSeq exon 454672 459678 . - . ID=exon-NM_001005221.2-3;Parent=rna-NM_001005221.2;
chr1 BestRefSeq CDS 450740 451678 . - 0 ID=cds-NP_001005221.2-1;Parent=rna-NM_001005221.2;
chr1 BestRefSeq CDS 452658 453675 . - 0 ID=cds-NP_001005221.2-2;Parent=rna-NM_001005221.2;
```



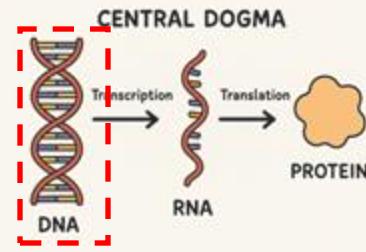
# Genome Annotation: Lift-over Problem



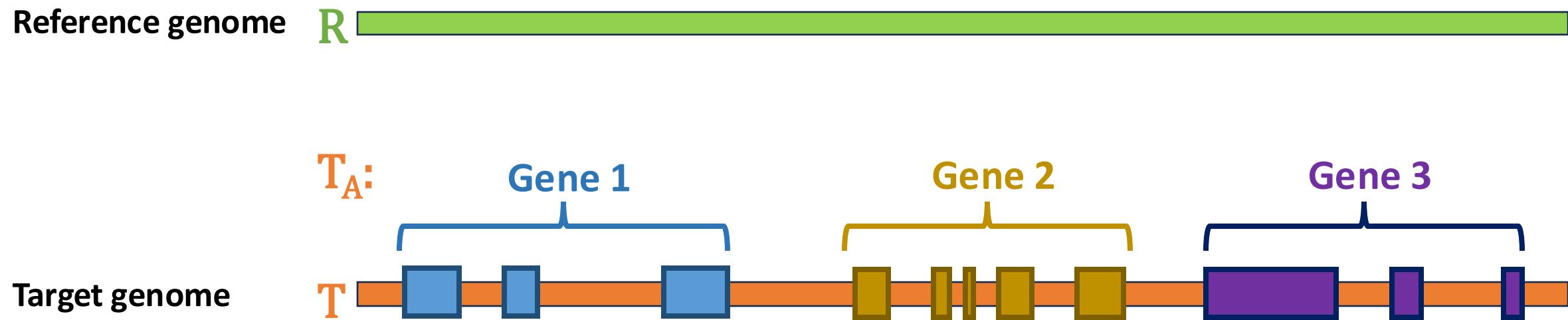
Chao, K. H., Heinz, J. M., Hoh, C., Mao, A., Shumate, A., Pertea, M., & Salzberg, S. L. (2025). Combining DNA and protein alignments to improve genome annotation with LiftOn. *Genome Research*, 35(2), 311-325.

GENOME  
RESEARCH

LiftOn



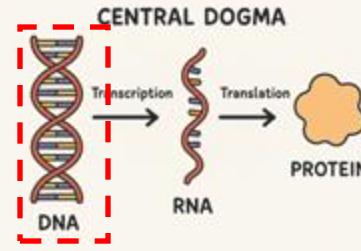
# Genome Annotation: Lift-over Problem



Chao, K. H., Heinz, J. M., Hoh, C., Mao, A., Shumate, A., Pertea, M., & Salzberg, S. L. (2025). Combining DNA and protein alignments to improve genome annotation with LiftOn. *Genome Research*, 35(2), 311-325.

GENOME  
RESEARCH

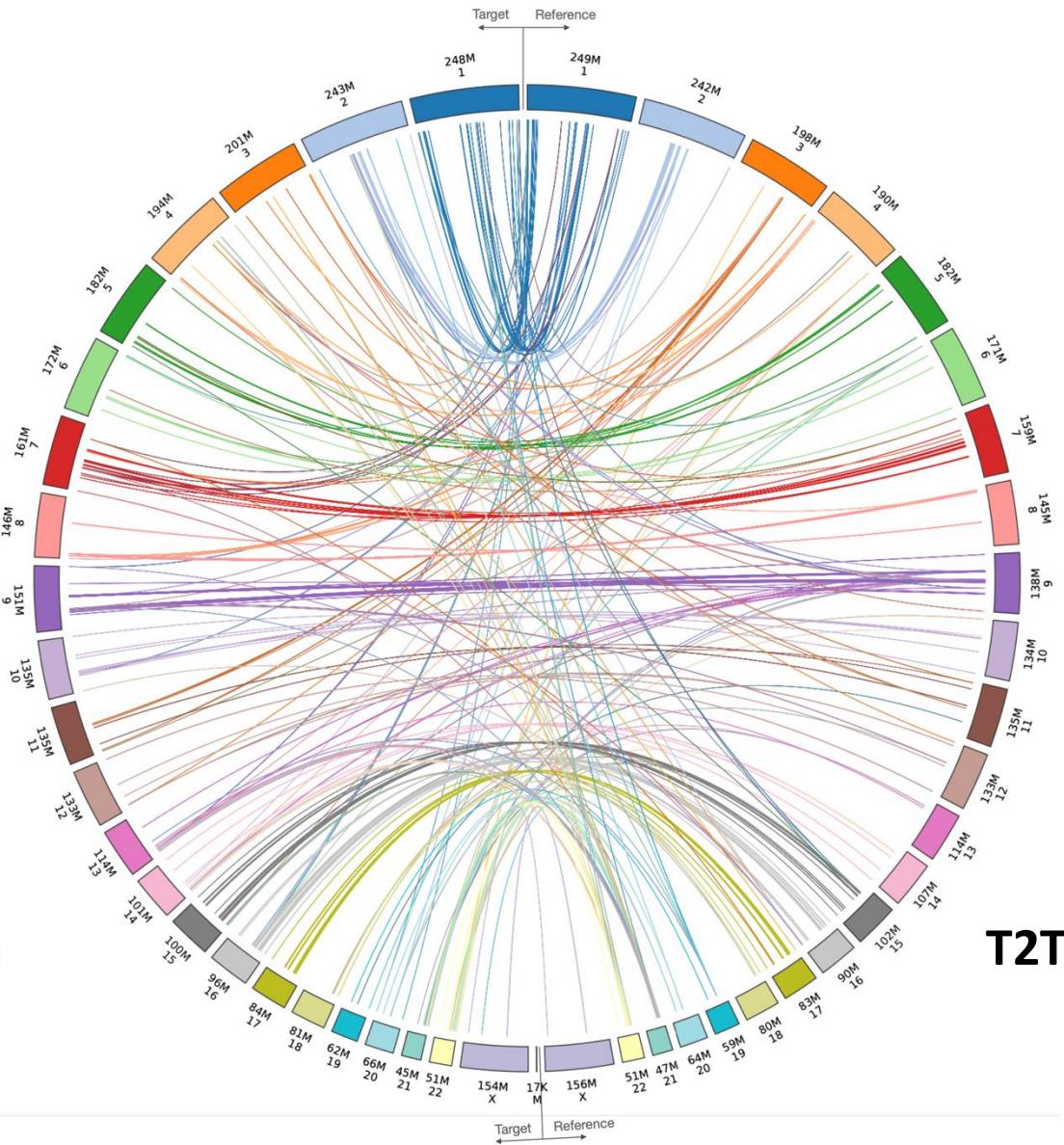
LiftOn



**Han1 – Southern Chinese Han Individual**

ChatGPT 4o:  
Generate my profile photo in Ghibli style.

Chao, K. H., Heinz, J. M., Hoh, C., Mao, A., Shumate, A., Pertea, M., & Salzberg, S. L. (2025). Combining DNA and protein alignments to improve genome annotation with LiftOn. *Genome Research*, 35(2), 311-325.



**T2T-CHM13 – Northern European Individual**

ChatGPT 4o:  
Generate a European origin individual icon in Ghibli style.

**GENOME  
RESEARCH**

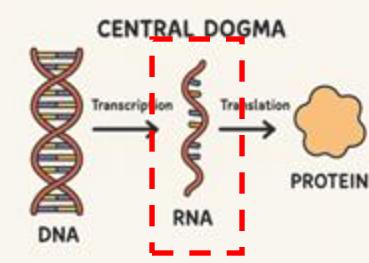
Genome Assembly

Genome Annotation

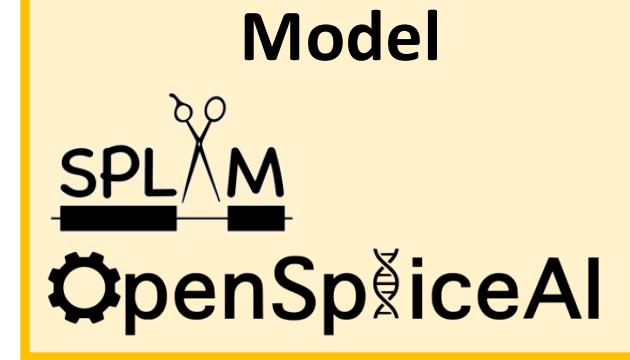
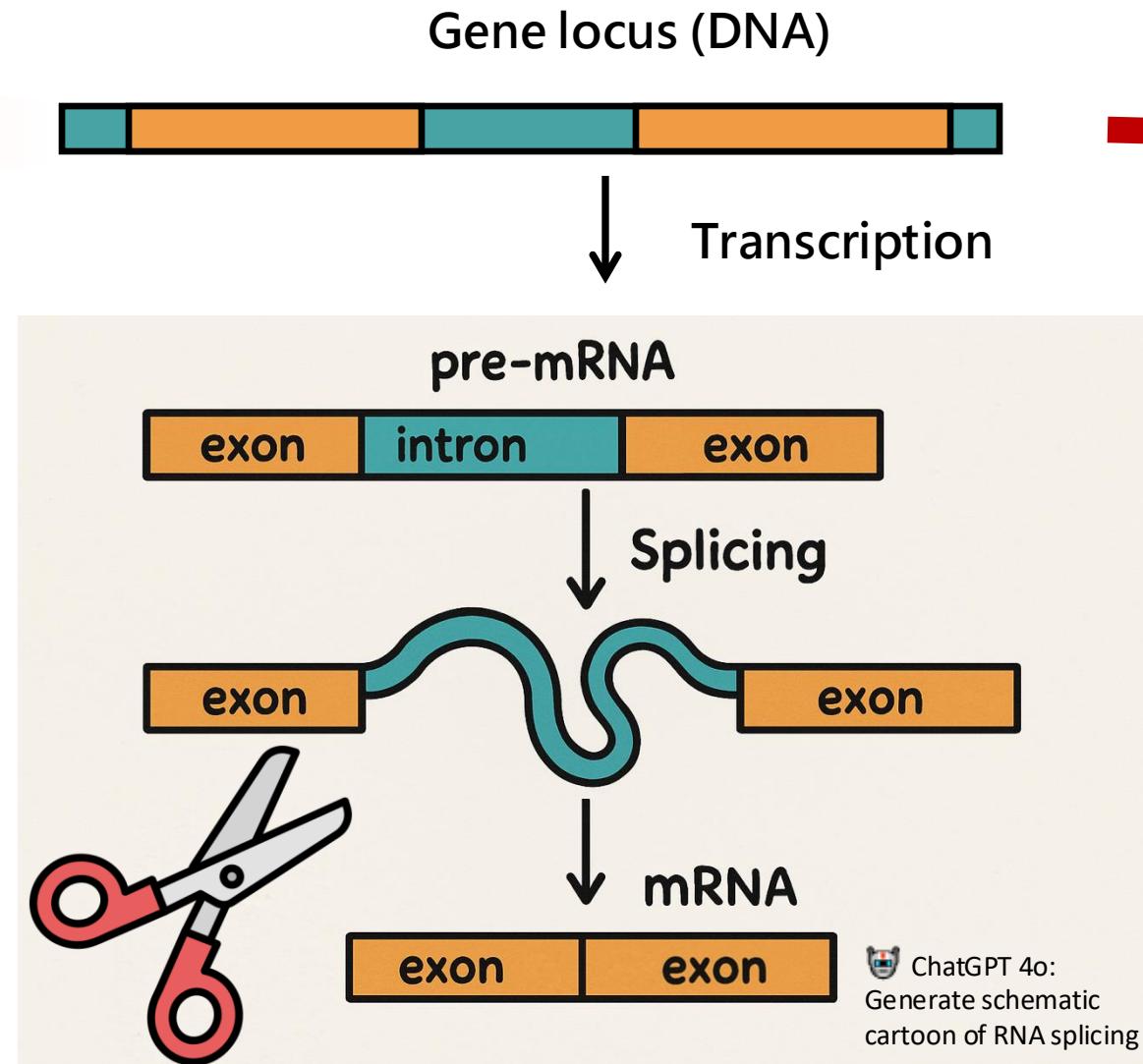
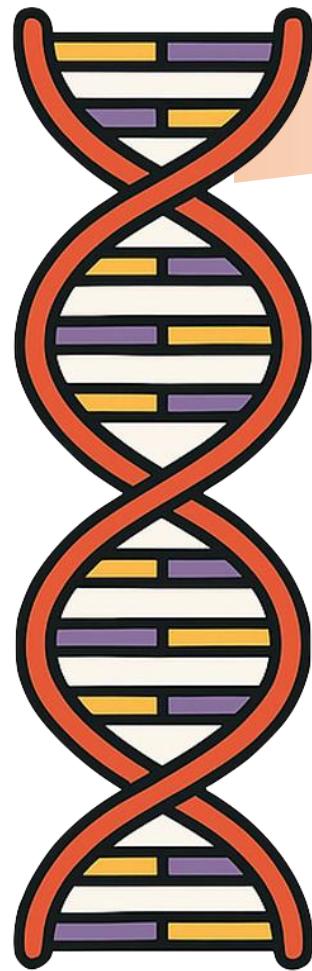
Splice Site Prediction

Gene Expression Prediction

**LiftOn**



# Splice Site Prediction

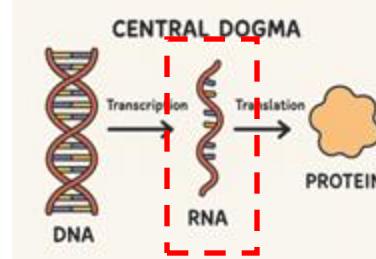


Where are the splice sites?

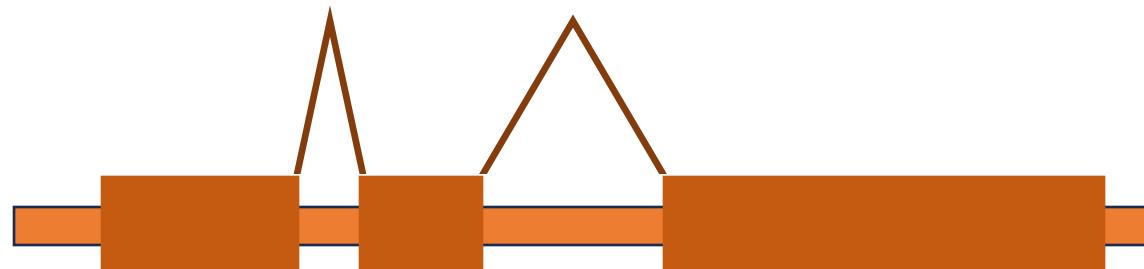
# Splice Site Prediction



openSpliceAI



Pre-mRNA



Genome Biology

**Chao, K. H., Mao, A., Salzberg, S. L., & Pertea, M. (2024).** Splam: a deep-learning-based splice site predictor that improves spliced alignments. *Genome biology*, 25(1), 243.

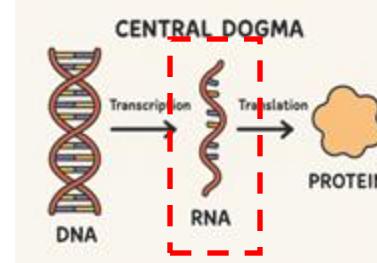
HUB

New AI tool pinpoints gene splicing with unmatched precision

bioRxiv  
THE PREPRINT SERVER FOR BIOLOGY

**Chao, K. H., Mao, A., Liu, A., Salzberg, S. L., & Pertea, M. (2025).** OpenSpliceAI: An efficient, modular implementation of SpliceAI enabling easy retraining on non-human species. *bioRxiv*, 2025-03.

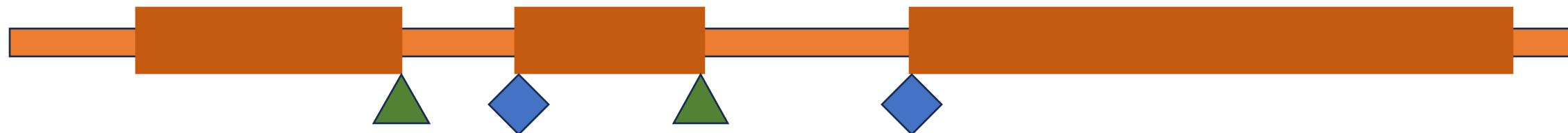
# Splice Site Prediction



Donor: 2

 Acceptor: 1

Neither: 0



X  
Y

AGACTCAGCCCCCGGAGACTTAGTTAGAGGAAGAAAAGGTAGGACAGAAGAAAAAGGCAGGACATAAGGTGCTGGCCCAGGGCGG



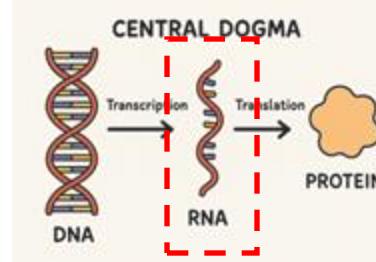
**Chao, K. H.**, Mao, A., Salzberg, S. L., & Pertea, M. (2024). Splam: a deep-learning-based splice site predictor that improves spliced alignments. *Genome biology*, 25(1), 243.



## New AI tool pinpoints gene splicing with unmatched precision



**Chao, K. H.**, Mao, A., Liu, A., Salzberg, S. L., & Pertea, M. (2025). OpenSpliceAI: An efficient, modular implementation of SpliceAI enabling easy retraining on non-human species. *bioRxiv*, 2025-03.



# Splice Site Prediction: data processing

~20k protein-coding genes

## Raw gene DNA sequence

Gene 1		$L = 33200$
Gene 2		$L = 14600$
...		
Gene n		$L = 25000$

X	Y
[7, 15000, 4]	[7, 5000, 3]
[3, 15000, 4]	[3, 5000, 3]
...	
[5, 15000, 4]	[5, 5000, 3]

Genome Biology

Chao, K. H., Mao, A., Salzberg, S. L., & Pertea, M. (2024). Splam: a deep-learning-based splice site predictor that improves spliced alignments. *Genome biology*, 25(1), 243.

HUB 

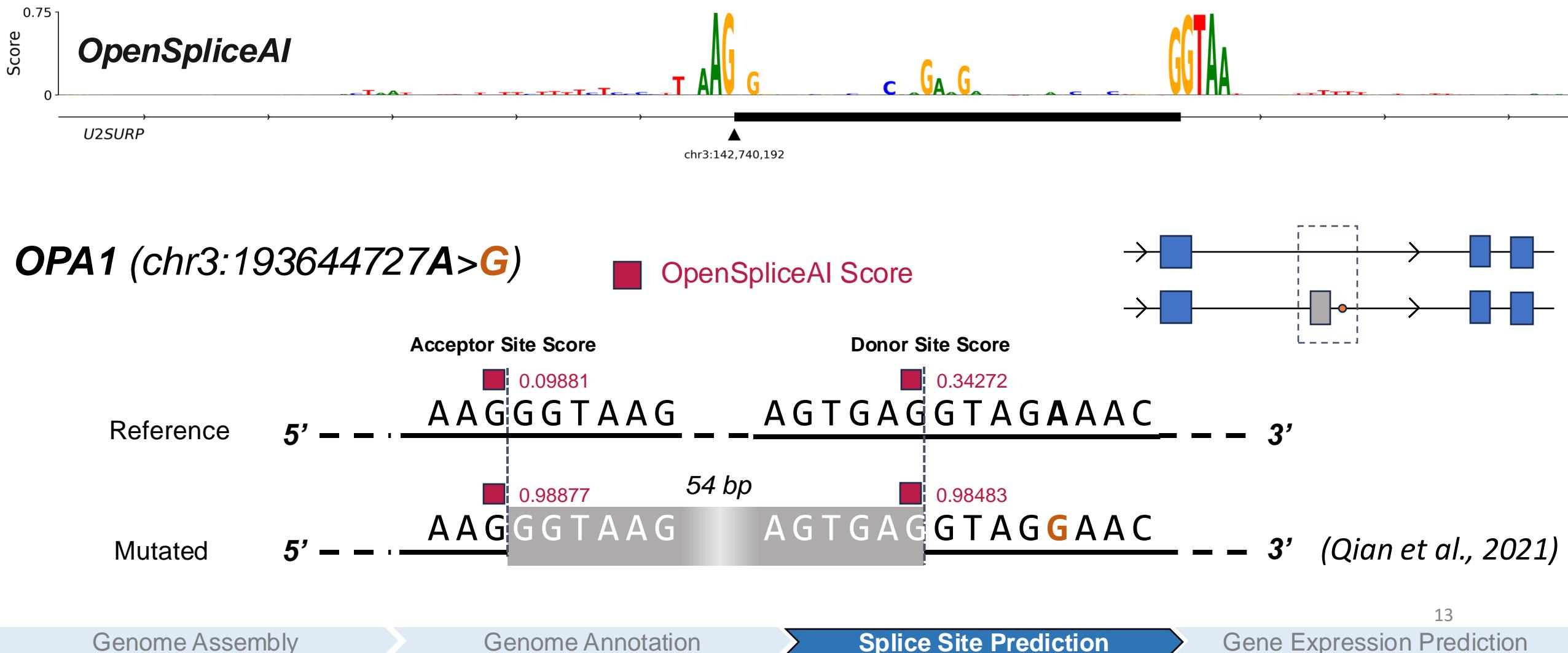
New AI tool pinpoints gene splicing with unmatched precision

bioRxiv

THE PREPRINT SERVER FOR BIOLOGY

Chao, K. H., Mao, A., Liu, A., Salzberg, S. L., & Pertea, M. (2025). OpenSpliceAI: An efficient, modular implementation of SpliceAI enabling easy retraining on non-human species. *bioRxiv*, 2025-03.

# Splice Site Prediction: what did model learn?

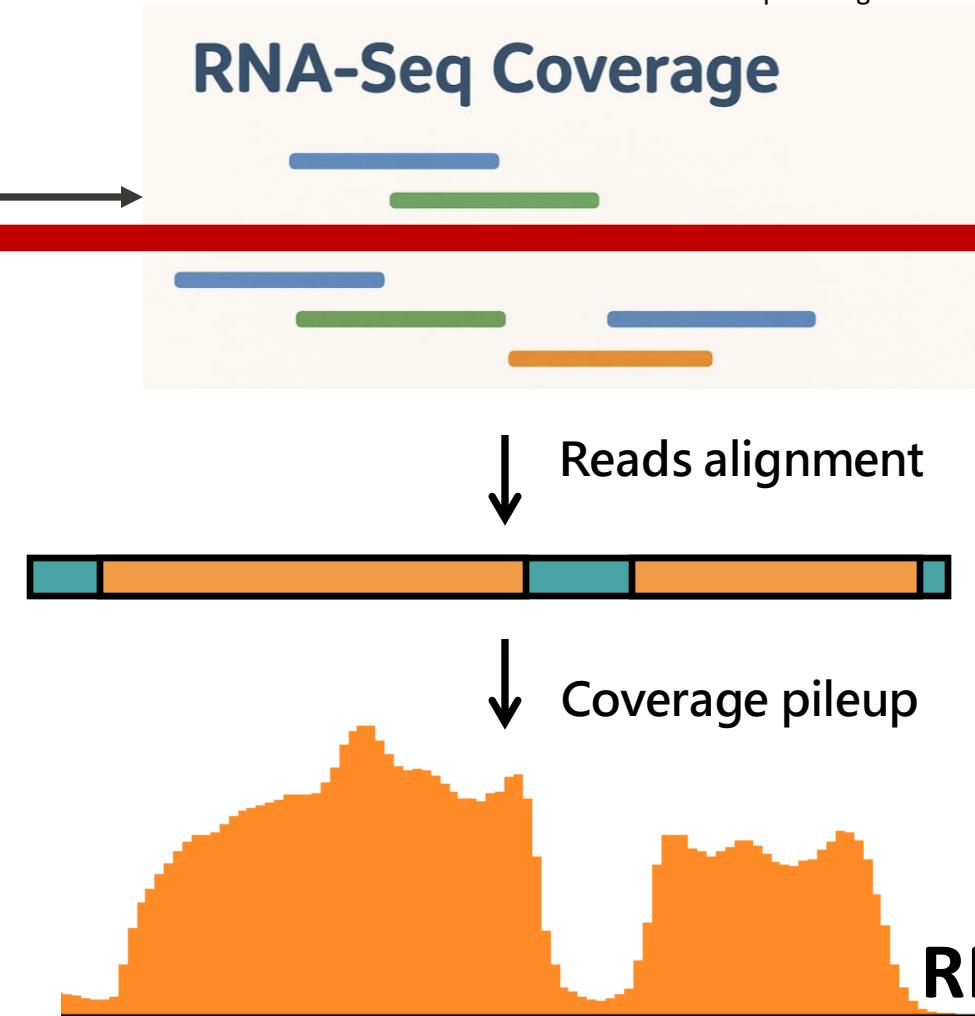
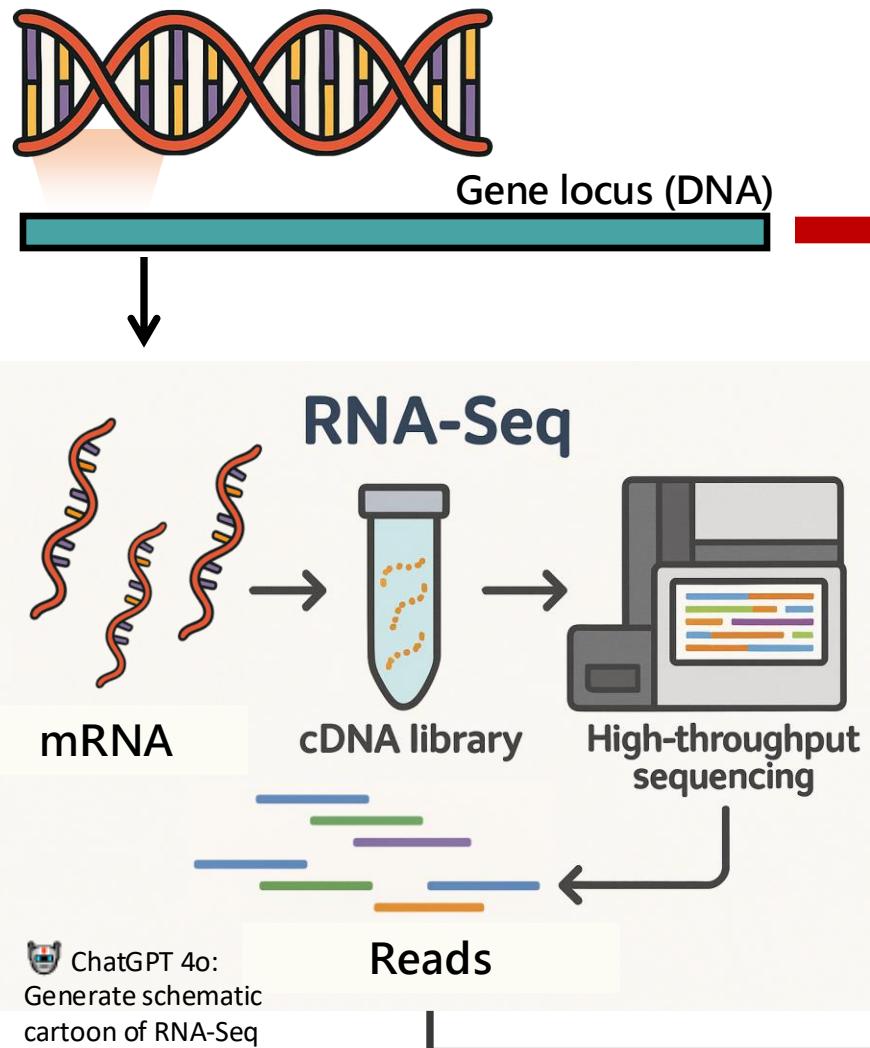
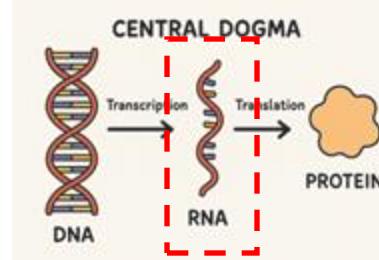




Calico

# RNA-Seq Prediction

ChatGPT 4o:  
Generate schematic cartoon  
of RNA-Seq coverage



**Model**  
**Fungal Model**

**RNA-Seq expression?**

Genome Assembly

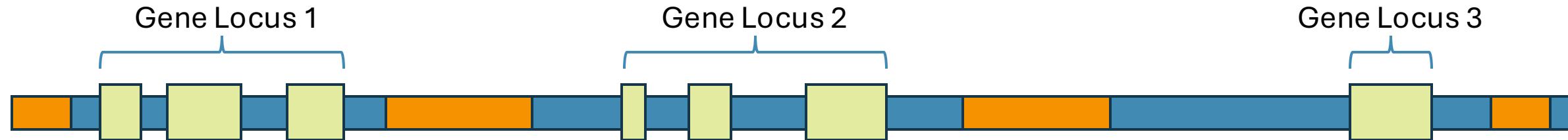
Genome Annotation

Splice Site Prediction

Gene Expression Prediction

Repeat regions  
Coding regions

# Processing 1K fungal genomes



# Processing 1K fungal genomes

Repeat regions  
Coding regions



16384



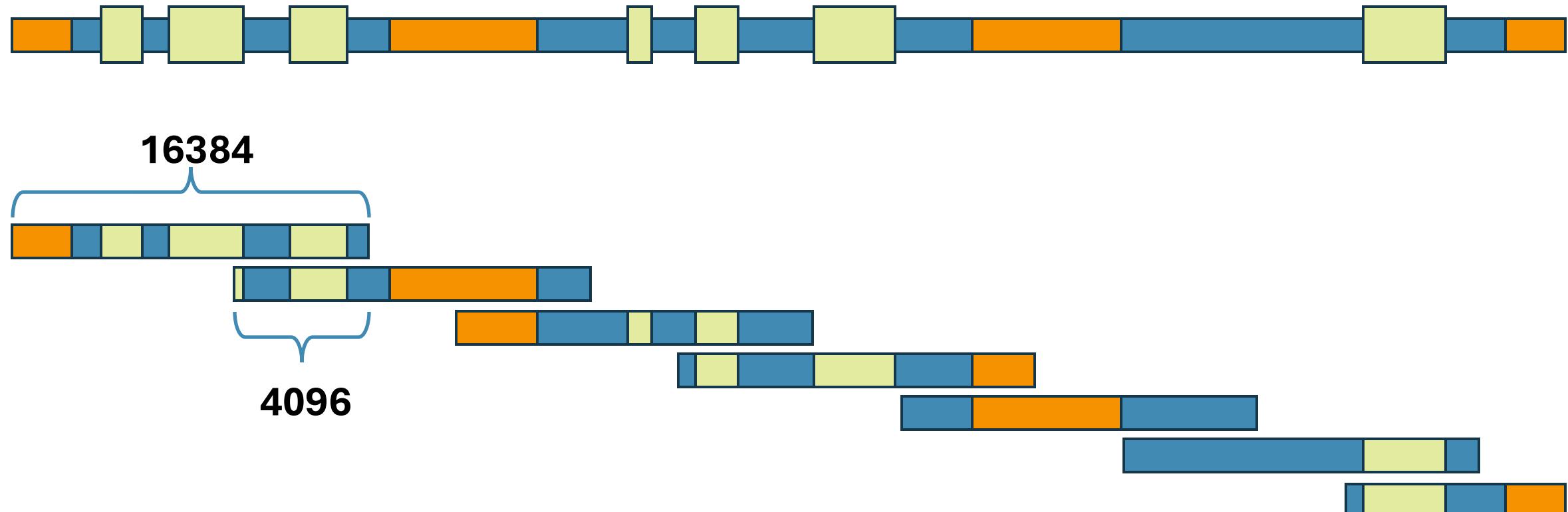
4096



~ 7 genes per window

# Processing 1K fungal genomes

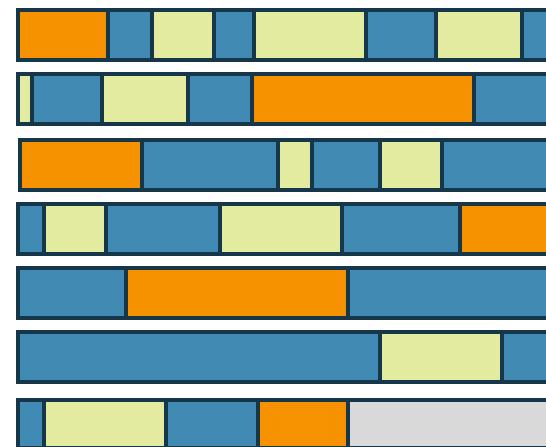
Repeat regions  
Coding regions



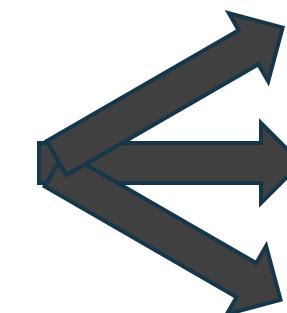
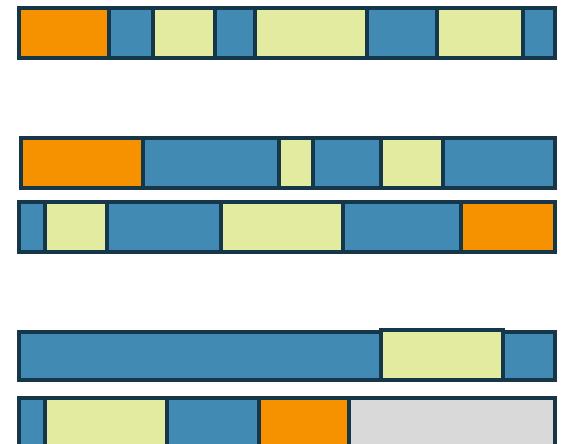
~ 7 genes per window

# Processing 1K fungal genomes

Repeat regions  
Coding regions



7% repeat threshold



Training

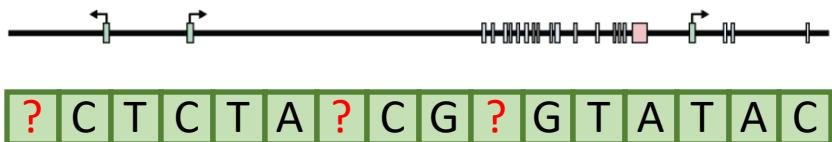
Validation  
(chrXI, chrXIII, chrXV)

Testing  
(chrXII, chrXIV, chrXVI)



Calico

16384bp



16384 \* 4

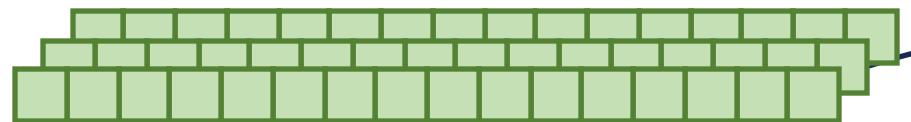
Masked language modeling loss

A	.8				.0			.0	
C	.1					.0		.7	
G	.1					.9		.1	
T	.0					.1		.2	

Reverse complementary



1bp res

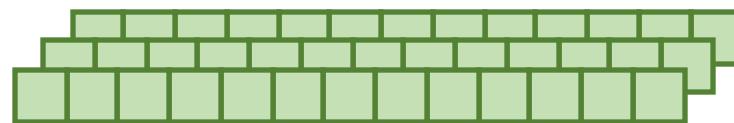


1bp res

...

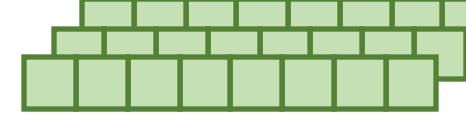
...

16bp res



16bp res

32bp res



32bp res

64bp res



64bp res

128bp res



128bp res

Transformer  
Blocks (8x)



Genome Assembly

Genome Annotation

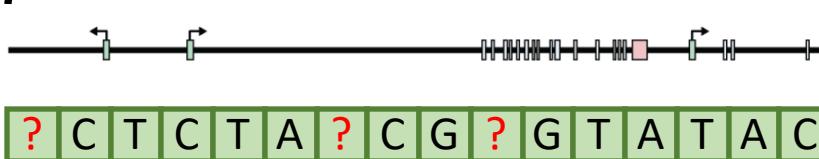
Splice Site Prediction

Gene Expression Prediction



Calico

16384bp



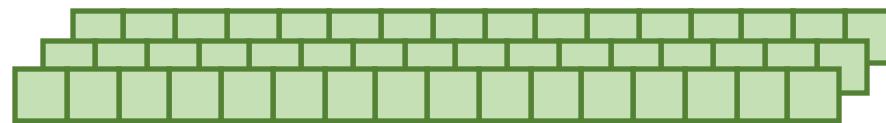
[ CHiP-exo (1128)  
Histone marks (20)  
RNA-Seq (1340) ]

### Coverage Tracks

Reverse complementary



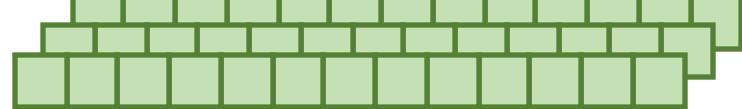
1bp res



...

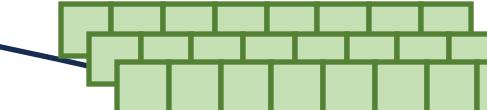
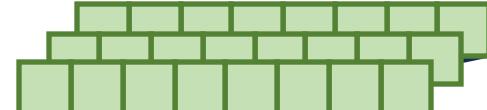
...

16bp res



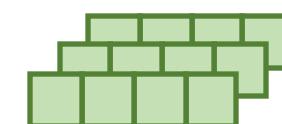
16bp res

32bp res



32bp res

64bp res



64bp res

128bp res



Transformer  
Blocks (8x)

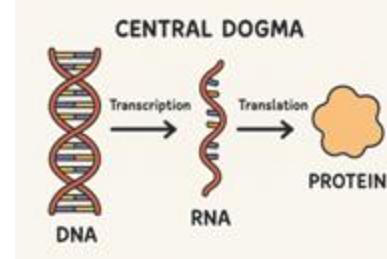


Genome Assembly

Genome Annotation

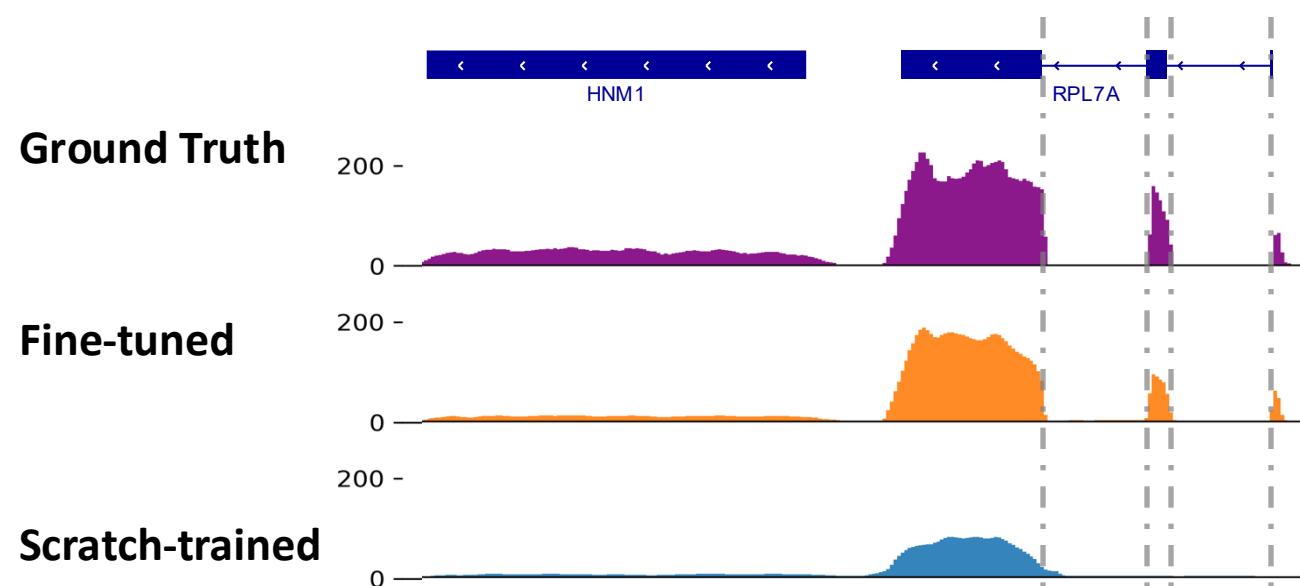
Splice Site Prediction

Gene Expression Prediction

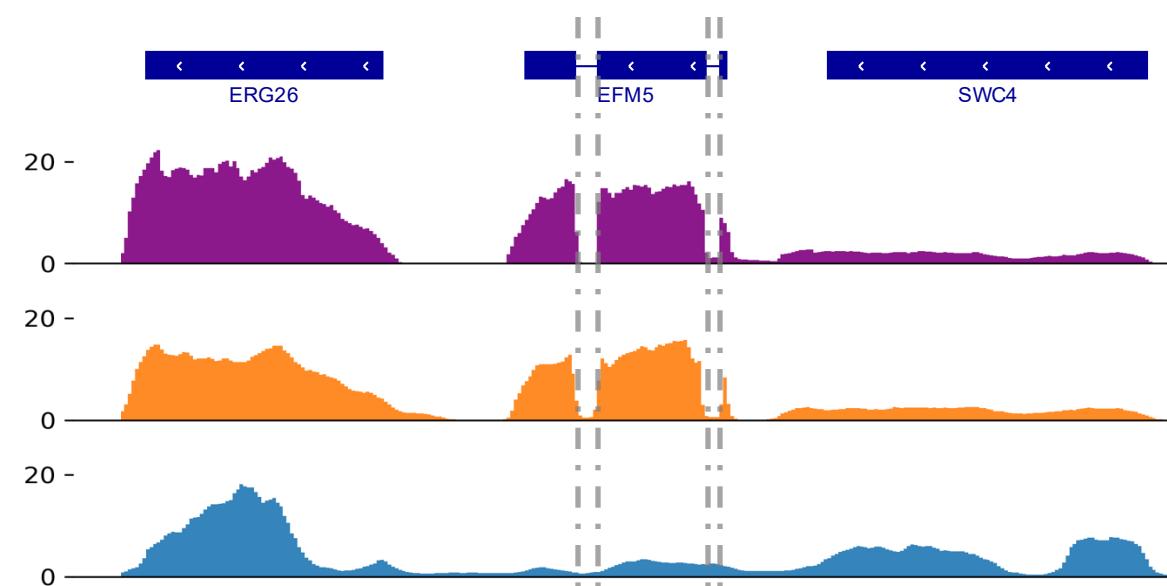


# RNA-Seq Coverage Prediction Examples

**chrVII:362,180-366,023 (RNA-Seq tracks)**

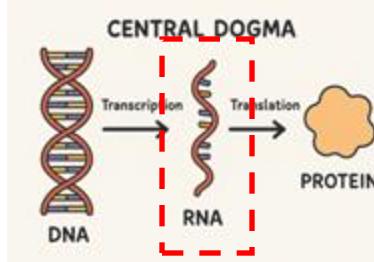


**chrVII:495,374-499,965 (RNA-Seq tracks)**



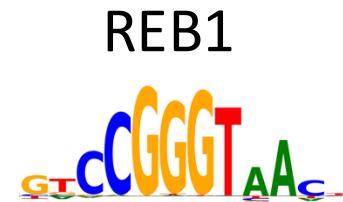
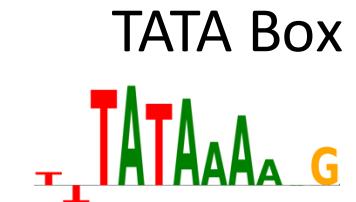
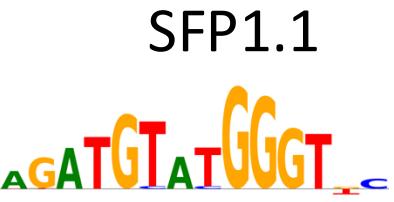
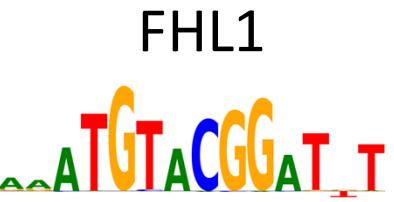


Calico

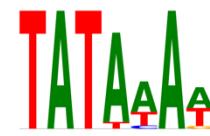


# Learned Regulatory Elements

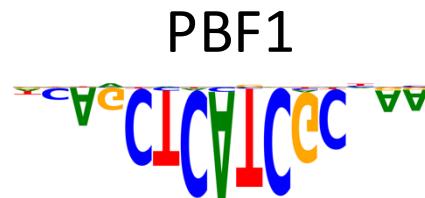
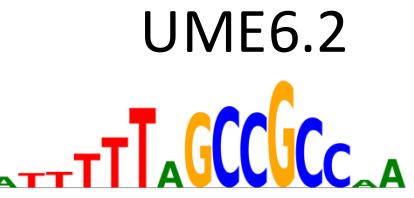
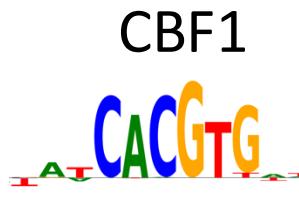
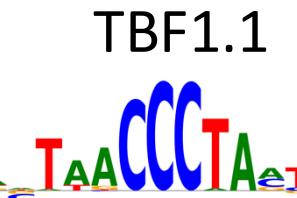
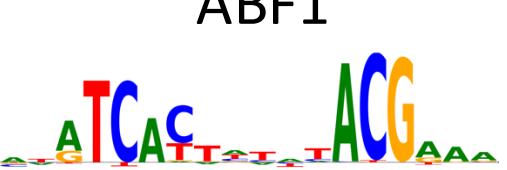
Fine-tuned motif



Motif DB

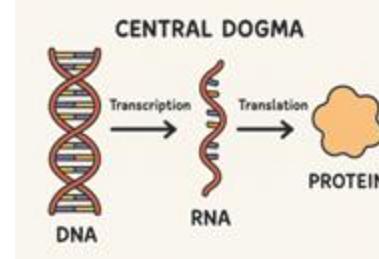


Fine-tuned motif



Motif DB





# Learned Splicing Motifs

Yeast splicing motifs  
(Schirman et. al.)

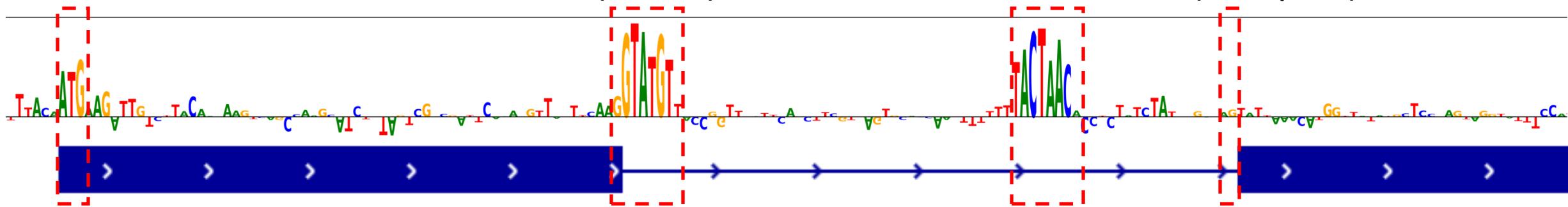


Start Codon

Splice site  
(donor)

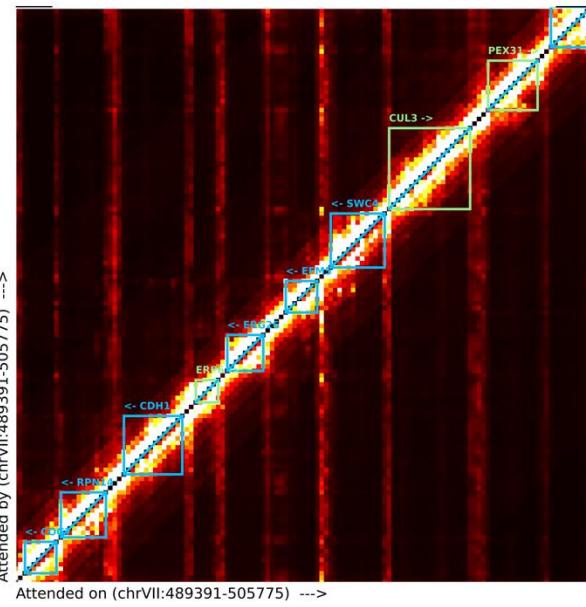
Branch point

Splice site  
(acceptor)

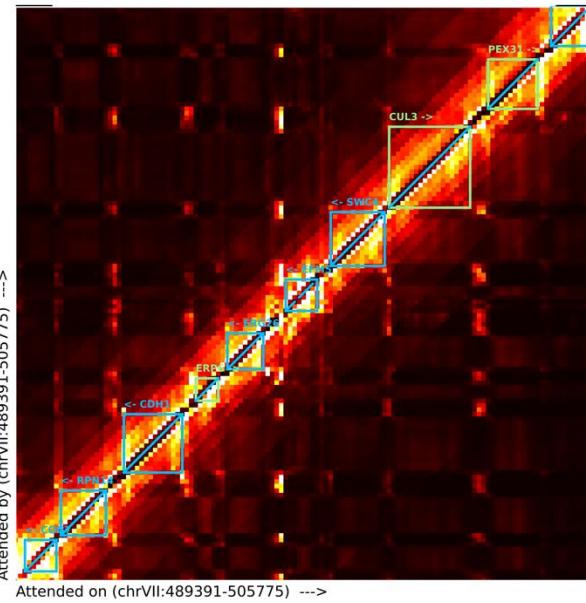


## Language model

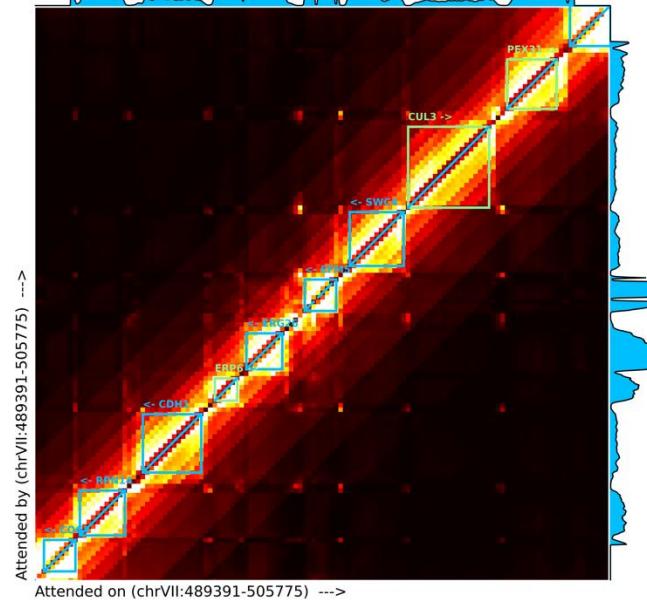
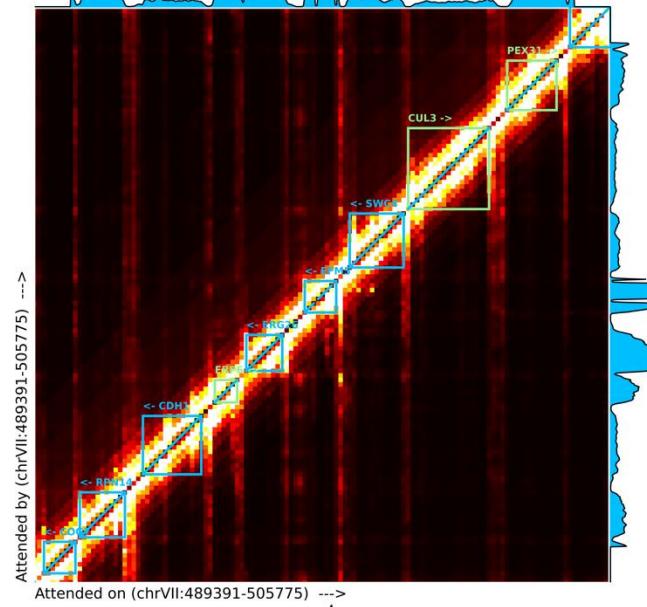
### First Self-attention Layer



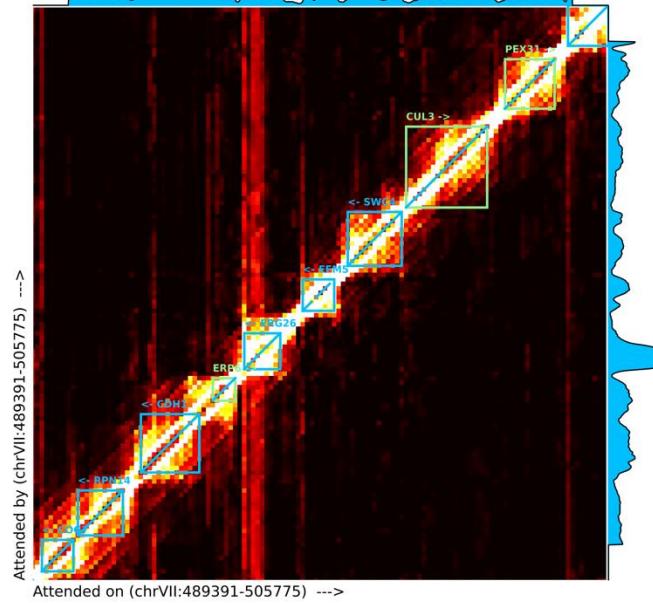
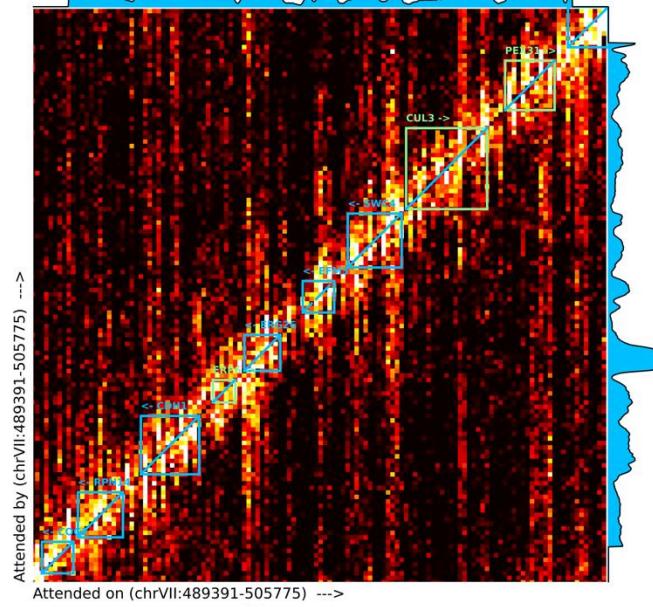
### Last Two Self-attention Layer



## Fine-tuned RNA-Seq model



## Scratch-trained RNA-Seq model

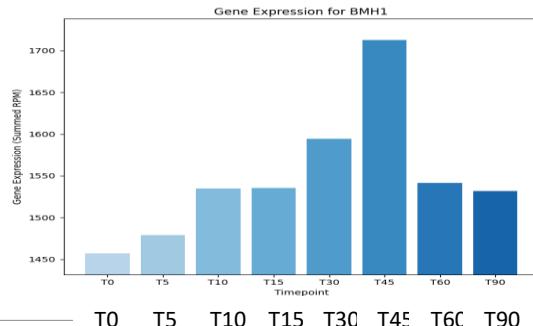




Calico

# MET4 Induction

## chrX:322935-323435 (Promoter region of BNA3)



YeTFaSCo DB motif



MET4

Genome Assembly

Genome Annotation

Splice Site Prediction

Gene Expression Prediction

*“If you think of mathematics as the perfect description language for physics, then AI might be the perfect one for biology.”*

Demis Hassabis, CEO of DeepMind, 2022



[khchao.com](http://khchao.com)



@kuanhaochao.bsky.social



@KuanHaoChao



Kuanhao-Chao

# Acknowledgement

